



ProtoNet 6.0: Organizing 10 million protein sequences in a compact hierarchical family tree

Journal:	<i>Nucleic Acids Research</i>
Manuscript ID:	NAR-02320-DATA-E-2011.R1
Manuscript Type:	6 Database Issue
Key Words:	Protein family, Annotation, Clustering , Protein structure, Functional inference

SCHOLARONE™
Manuscripts

Review

1
2
3
4 **ProtoNet 6.0: Organizing 10 million protein sequences in a**
5
6
7 **compact hierarchical family tree**
8
9
10

11
12
13 Nadav Rappoport¹, Solange Karsenty¹, Amos Stein¹, Nathan Linial¹ and Michal
14 Linial²
15

16
17 ¹School of Computer Science and Engineering, ²Department of Biological
18 Chemistry, Institute of Life Sciences, The Sudarsky Center for Computational
19 Biology, The Hebrew University of Jerusalem, Israel
20
21

22
23
24 *Corresponding author
25

26 Corresponding author details:

27 Michal Linial

28 E-Mail: michall@cc.huji.ac.il
29
30

31
32 Department of Biological Chemistry, Institute of Life Sciences
33 The Hebrew University
34 Givat Ram Campus
35 Jerusalem, 91904
36 Israel
37

38 Phone: 972-2-6585425

39 FAX: 972-2-6586448
40
41
42
43
44
45

46 Table 1

47 Figures 3
48
49
50
51

52 Running title: ProtoNet 6.0: Clustering the protein sequence space
53
54
55
56
57
58
59
60

ABSTRACT

ProtoNet 6.0 (<http://www.protonet.cs.huji.ac.il>) is a data structure of protein families that cover the protein sequence space. These families are generated through an unsupervised bottom-up clustering algorithm. This algorithm organizes large sets of proteins in a hierarchical tree that yields high quality protein families. The 2012 ProtoNet (Version 6.0) tree includes over 9 million proteins of which 5.5% come from UniProtKB/SwissProt and the rest from UniProtKB/TrEMBL. The hierarchical tree structure is based on an all-against-all comparison of 2.5 million representatives of UniRef50. Rigorous annotation-based quality tests prune the tree to most informative 162,088 clusters. Every high quality cluster is assigned a ProtoName that reflects the most significant annotation of its proteins. These annotations are dominated by GO terms, UniProt/Swiss-Prot keywords and InterPro. ProtoNet 6.0 operates in a default mode. When used in the advanced mode, this data structure offers the user a view of the family tree at any desired level of resolution. Systematic comparisons with previous versions of ProtoNet are carried out. They show how our view of protein families evolves, as larger parts of the sequence space become known. ProtoNet 6.0 provides numerous tools to navigate the hierarchy of clusters.

INTRODUCTION

ProtoNet (1) was launched in 2002. The goal of this system was to achieve an automatic hierarchical clustering of the protein sequences space. It covered 94,000 protein sequences from Swiss-Prot. Now, almost 10 years later, our census of proteins has grown tremendously. Thus, the UniProtKB database of protein sequences (2) includes over 17 millions proteins (UniProt, August 2011) of which 0.53 million proteins form the UniProtKB/Swiss-Prot section. While the size of UniProtKB/Swiss-Prot grew from 2002 by a factor of 5 (SwissProt release 40.0, October 2001), the TrEMBL section (TrEMBL Release 18.0) went from 550,000 to 16.5 million sequences, a 30-fold increase during the same period.

Notably, even in the curated high quality UniProtKB/Swiss-Prot section, only 25% of the proteins carry evidence at the protein or transcript levels, while 70% of the sequences are inferred from homology and about 3% remain questionable and marked as predicted or even uncertain proteins. The situation with the millions of sequences from UniProtKB/TrEMBL is far less satisfying. Only 3% carry some experimental supporting evidence and the majority of sequences (74%) are only based on prediction. With this immense growth in the number of protein sequences, it is clear that only unsupervised methods can cope with this data set. We need algorithms that can automatically trace the functional and evolutionary relatedness among protein sequences (3).

Assigning biological functions to proteins is a major obstacle and a challenging task (4,5). Despite important progress in structural genomics, enzyme classifications and

1
2
3 phylogenomics, the goal of automatic functional inference is far from being reached
4
5 (3,6-8). Numerous motif recognition algorithms, statistical model-based and
6
7 clustering methods were developed during the last two decades for the purpose of
8
9 handling the growing number of sequences. These methods differ in their coverage,
10
11 the level of manual curation involved and even in the basic definition of a domain
12
13 family. For example, Pfam (9), SMART (10), EVEREST (11), PANTHER (12) and
14
15 Gene3D (13) are based on thousands of profile hidden Markov models (profile
16
17 HMMs). New sequences that pass a pre-determined threshold of similarity are
18
19 assigned to the corresponding model domain family. Additional resources are based
20
21 on algorithms that search for signature, regular expressions or Position Specific
22
23 Scoring Matrix (PSSM) fingerprints. Representative databases that follow this
24
25 paradigm include PROSITE (14), PRINTS (15), ProDom (16), BLOCKS (17). The
26
27 above resources are based on sequence data.
28
29
30
31

32
33 In addition, integrative resources such as PIRSF (18), CDD (19) and InterPro (20)
34
35 take a different approach to the end of attaining higher coverage of the protein space.
36
37 They accomplish this by merging a variety of external sources with a focus on protein
38
39 families, domains and functional sites. The classifications of SCOP (21), CATH (22)
40
41 and SUPERFAMILY (23) rest on 3D-structural information. A functional perspective
42
43 is offered by the ontology-based resource of Gene Ontology (GO) (24).
44
45

46
47 The available data is highly redundant, which creates a major difficulty in this area.
48
49 Thus the main archive of UniProt database contains 25 million sequences (25) which
50
51 represent about 17 million unique proteins. The UniRef50 with only 4 million
52
53 sequence is created by grouping together proteins with >50% identical amino acids.
54
55
56
57
58
59
60

1
2
3 However, in order to study sequence homologies and the evolution of protein
4 families, they must be viewed at a much finer level of granularity.
5
6

7
8 In order to deal with the enormous number of known protein sequences, ProtoNet 6.0
9 generates automatically, with no supervision a consistent classification tree. This
10 system covers over 9 million proteins from UniProtKB. To address the expected
11 future growth in the number of protein sequences, the system is equipped with a
12 protocol for maintenance and updating. A system-provided confidence parameter
13 quantifies the quality of every cluster in ProtoNet 6.0. Additional tools for analysis
14 and visualization enhance the user's navigation options through the ProtoNet tree.
15 These tools provide a rich biological context for the observed parts of the tree.
16
17

18 We describe here the newly introduced capabilities and improvements compared with
19 the previous version (26) where one million proteins were classified (1,072,911
20 sequences, UniProt Release, Feb. 2005, ProtoNet Version 4.0).
21
22

23 24 25 26 27 28 29 30 31 32 33 34 35 **PROTEIN SEQUENCES DATABASE**

36 All database sources used in ProtoNet 6.0 has been thoroughly updated. The most
37 critical aspect is the use of UniRef50 clusters as our basic objects. On average a
38 UniRef50 cluster contains 4 proteins. Thus, the 2,478,328 UniRef50 proteins that are
39 included in ProtoNet 6.0 represent over 9 million sequences. In comparison the
40 number of protein sequences in ProtoNet 4.0 is 1,072,911.
41
42
43
44
45
46
47
48

49 50 51 52 53 54 55 56 57 58 59 60 **PROTONET TREE CONSTRUCTION**

The basic algorithm of ProtoNet was previously described (1,27). It starts by pre-
calculating an all-against-all BLAST similarity score (28) for all protein

1
2
3 representatives from the UniRef50 resource (called cluster seed proteins). The
4 similarities' E-scores were used to produce a continuous hierarchical bottom-up
5 clustering process. At each step, the two most similar protein clusters are joined (the
6 exact algorithm is described (29)). Importantly, BLAST E-score with an extremely
7 relaxed threshold is considered throughout the ProtoNet construction (E-score=100).
8
9 The bottom-up agglomerative clustering of the ProtoNet algorithm benefits from such
10 relaxed E-score distances in constructing a robust family tree. A key ingredient of
11 ProtoNet 6.0 that is essential for handling such a large number of proteins is the
12 Constrained Memory-ProtoNet algorithm (29).
13
14
15
16
17
18
19
20
21
22
23

24 The result is a hierarchy of protein clusters at various degrees of biological
25 granularity. This hierarchy is structured as a collection of trees that form what we call
26 ProtoNet Tree (actually it is a ProtoNet forest). The root clusters contain all the
27 proteins of the tree while other clusters represent subdivisions of proteins into smaller
28 groups. The hierarchical definitions allow the user to navigate from a protein to the
29 sub-family and the super-family levels in order to discover specific functions and
30 evolutionary signals.
31
32
33
34
35
36
37
38
39

40 THE HIERARCHY'S QUALITY

41
42
43
44 The entire protocol to construct ProtoNet is unsupervised and therefore no annotations
45 are included. However, measuring the correspondence between a given cluster and
46 specific annotations that are provided by external expert systems is essential for the
47 *supervised validation* of the automatically generated ProtoNet clusters.
48
49
50
51
52

53
54 We thus define the notion of a correspondence score (CS). The CS for a specific
55 cluster and a given keyword is a measure of correlation between two. Formally, let us
56
57
58
59
60

1
2
3 fix a cluster C in the ProtoNet tree and a keyword K (from a specific source such as
4 InterPro). Let c be the set of proteins in cluster C and let k be the set of proteins in the
5 system annotated by keyword K. We define the CS as:
6
7
8

9
10 Correspondence Score (cluster C for keyword K) = $CS(C, K) = \frac{|c \cap k|}{|c \cup k|}$
11
12

13
14 The cluster receiving the maximal score for keyword K (called K's *best cluster*) is
15 considered the cluster that best represents K within the ProtoNet tree. The score for a
16 given cluster on keyword K ranges from 0 (no correspondence) to 1 (the cluster C is
17 comprised of all the proteins with keyword K). The CS values are used as a quality
18 measure for the ProtoNet tree. For example, we may consider the distribution of CS
19 value over all ProtoNet clusters or over clusters of size that exceeds some cutoff
20 threshold. In order to obtain a biologically relevant view of the hierarchy, we applied
21 several tests that allow us to focus only on the clusters that are enriched with some
22 coherent biological information.
23
24
25
26
27
28
29
30
31
32

33
34 The main algorithmic difference between ProtoNet 6.0 and the earlier version
35 ProtoNet 4.0 (26) is the use of CM-ProtoNet (29). We refined the clusters' quality test
36 by evaluating the CM-ProtoNet method over a single-linkage performance (that is
37 implemented in ClusTr (30)). The tests were carried out on 3.2 million proteins from
38 ProtoNet 5.1 (Table 1). In addition, we tested the impact of selecting UniRef50 as
39 cluster seed proteins for ProtoNet. It can be seen that CM-ProtoNet outperforms the
40 other methods that were applied to the same set of proteins. Notice that the main
41 improvement of MC-ProtoNet comes from enhanced sensitivity. The performance of
42 the *Single linkage* algorithm drops drastically due to a low sensitivity. We tested three
43 choices of cluster seeds: UniRef50 representatives (the choice that we finally
44 adopted), UniRef90 (proteins sharing >90% sequence identity) and the complete
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 redundant protein sequences. It is remarkable that the quality of clustering with
4
5 respect to all three choices remains essentially unchanged.
6
7

8
9 The same tests with respect to a set of keywords from Pfam Clan (9) validated the
10
11 high performance of the MC-ProtoNet algorithm over other clustering methods (not
12
13 shown). We confirmed that the protocol that was applied to construct the ProtoNet 6.0
14
15 produces a stable tree with a collection of biologically coherent families and super-
16
17 families.
18
19

20 **SELECTING STABLE CLUSTERS**

21
22
23

24 The ProtoNet tree is huge, and the immense number of its protein clusters makes it
25
26 quite impractical to navigate the tree. In order to deal with this difficulty, we pruned
27
28 the tree. The basic idea is that many clusters that are created along the process of
29
30 generating the tree are biologically irrelevant and uninteresting. For example, a root
31
32 cluster in the ProtoNet forest contains thousands of unrelated proteins.
33
34
35

36 A process of repeated pair-wise merging yields a tree of size roughly twice the
37
38 number of leaves (see illustration in Figure 1A). Therefore, starting with the 2.5
39
40 million UniRef50 seed proteins we obtain 5.0 million clusters. We applied several
41
42 computational procedures that are aimed at reducing this number. Our aim is to
43
44 simplify the navigation in the system while maintaining the hierarchical structure and
45
46 with essentially no loss in clusters' quality.
47
48
49

50 To this end we sought intrinsic parameters of ProtoNet that measure the *stability* of a
51
52 cluster. One such parameter is Life Time (LT), which is the difference between the
53
54 time (i.e. merging steps) in which a cluster is created and the time it is merged to a
55
56 larger cluster. This number reflects the relative height of a cluster in the merging tree.
57
58
59
60

1
2
3 The level of the tree (called ProtoLevel, PL) is used as an internal monotonic *timer for*
4 *merging*, along the clustering process (which is reflected by the index of the cluster,
5
6
7 Figure 1A). Individual protein sequences have PL=0 and for the root of the ProtoNet
8
9
10 tree PL=100. The idea is that *stable clusters* tend to be more relevant biologically. We
11
12 thus used a tradeoff between the number of clusters that are retained and the reduction
13
14 in the performance of the clusters, measured by the average of the CS for all clusters.
15
16 A minimal reduction in the average CS score for the InterPro keyword annotations
17
18 was attained for $LT < 1.0$. We thus set the $LT=1.0$ as a default parameter (see
19
20 “*Advanced Navigation*”).
21
22

23
24 Figure 1 illustrates the pruning process at different LT cutoffs (marked x, y).
25
26 Evidently, fewer valid clusters (colored red) remain as LT is increased. Figure 1B
27
28 shows a cluster summary at different LT cutoffs. Note that the statistical parameters
29
30 of the analyzed clusters depend on the choice of LT values.
31
32

33
34 The pruned version of the ProtoNet 6.0 tree at a $LT=1.0$ and $PL=90$ has 162,088 high
35
36 quality stable clusters. With these parameters the original number of 5 million clusters
37
38 (including leaves) is reduced by a factor of about 30.
39
40

41 ANNOTATION INFERENCE

42
43
44

45 Functional inference in ProtoNet 6.0 is done by an automatic high-confidence method
46
47 that infers the functional annotation of a cluster by integrating the annotations of its
48
49 individual proteins. The method builds on functional annotations from multiple
50
51 resources including InterPro, Gene Ontology (GO) (24), UniProt keywords (2),
52
53 ENZYME (31) and more. We consider all the annotations that cover >1% of the
54
55 proteins and focus on those that best fit the proteins of the cluster.
56
57
58
59
60

1
2
3 Evidently, automatic inference cannot be error-free. Thus, a predetermined specificity
4 threshold is calculated for the keywords associated with the cluster's proteins. Such
5 annotation is assigned as the ProtoName (Figure 2). To avoid faulty inference, we
6 calculated ProtoName for clusters in which >20% of the proteins share the specific
7 annotation where this annotation shows an enrichment of p-value <0.005. Recall that
8 presenting additional names for a cluster often hints at a novel overlooked function or
9 the presence of multi-domain proteins that exhibit multi-functionality.
10
11

12 Each of the ~162,000 stable clusters was assigned a ProtoName. On average, a cluster
13 is associated with 9.7 possible names. Most names are derived from Taxonomy
14 (33%), UniProt (19%), GO (18%), InterPro (17%) and the rest includes information
15 from structural classifications (e.g., SCOP (21) and CATH (22)) or ENZYME-based
16 annotations (31). A partition of the unique clusters according to their annotation types
17 is shown (Figure 2). Notably, most annotation types contribute to some ProtoName.
18 This suggests that the integration of knowledge from diverse annotation sources
19 substantially improves the performance of the ProtoNet tree.
20
21

22 **GENOMIC VIEW ON PROTEIN CLUSTERS**

23 A huge number of organisms are represented in UniProtKB (Figure 2B). Still, a third
24 of the protein sequences originate from a relatively small number of organisms that
25 were completely sequenced. A substantial number of all these sequences (mostly from
26 multi-cellular organisms) also serve as genetic model organisms. Therefore, we
27 included a selected list of over 30 organisms on which the user can choose to focus.
28 These organisms represent all superkingdoms.
29
30

31 **WEBSITE PROPERTIES**

1
2
3 Several added features in the ProtoNet 6.0 website make it easier to reach an in-depth
4 analysis of the ProtoNet tree. We describe these new features in ‘simplified mode’
5 and in ‘advanced mode’ (Figure 3).
6
7
8
9

10 (i) BROWSING CLUSTER NAMES

11
12
13 Cluster names are now available for browsing. One can choose a keyword of interest
14 and view clusters that are named by it. Note that a keyword of low statistical
15 significance will be absent in ProtoName. Figure 2 shows the contribution of the
16 major annotation resources that are included in determining the ProtoName.
17
18
19
20
21
22

23 (ii) HYPOTHETICAL & PUTATIVE PROTEINS

24
25
26 The assignment of a biological function to clusters suggests a safe scheme of
27 assigning function to proteins with unknown function. Naively, the protein can be
28 assigned the function of all clusters to which it belongs. This can be applied for
29 ‘hypothetical’ and ‘putative’ proteins within the clusters. It can also be used for a new
30 user-provided protein sequence (with the "*Classify your protein*" option). We provide
31 a list of all the proteins that are marked as hypothetical and putative proteins in the
32 summary table (Figure 3).
33
34
35
36
37
38
39
40
41
42

43 (iii) PROTONET TREE RESOLUTION

44
45
46 Following the pruning process described above, ProtoNet is no longer a binary tree.
47 To cope with this non-binary condensed version, we introduced the *ProtoBrowser*
48 page that zooms in on the tree only in the vicinity of the cluster that is being analyzed.
49 A selected branch is shown in the context of related neighboring branches. The user
50 hovers the mouse over a cluster to display essential information such as the cluster
51
52
53
54
55
56
57
58
59
60

1
2
3 size, the number of proteins according to selected species (if a ‘genomic view’ was
4 activated). An example of such ProtoBrowser tree views is shown (Figures 3).
5
6
7

8 **(iv) INTEGRATION OF ANNOTATION SOURCES**

9

10
11 The functional analysis of a cluster is performed using PANDORA (32) visualization,
12 which allows in-depth analysis of large protein sets. The system allows direct export
13 from the cluster page to PANDORA. Using PANDORA it is possible to assess the
14 functional relevance of the proteins in the clusters from numerous biological aspects.
15 The annotation sources used by PANDORA were updated, and now offer ~200,000
16 different annotations, spanning several different biological domains.
17
18
19
20
21
22
23

24
25 Specifically, PANDORA extracts most of the annotations from UniProtKB. For
26 structural annotations CATH (22), SCOP and Gene3D (13) are considered. The
27 functional domain is covered by the 4 layers of the ENZYME classification (31) and
28 the Gene Ontology (GO) structure with the 3 main functional branches: Cellular
29 component, Biochemical function and Biological process (24).
30
31
32
33
34
35
36

37
38 The protein families are forwarded to PANDORA analysis tool that statistically
39 analyzes a given cluster by means of the annotations that are assigned to its proteins
40 (32). On average, each protein sequence in ProtoNet is associated with 6.6 different
41 annotation types (11 and 10 annotations for human and mouse, respectively).
42
43
44
45
46
47 PANDORA supports also each of the dozen domain and family resources of the
48 InterPro collection.
49
50

51
52 In a typical application of PANDORA the user concentrates on any of the 200,000
53 annotations with the query “*Get clusters containing proteins with a given keyword*”
54
55
56 (e.g., InterPro domain: GTPase-binding/formin homology 3). In response one receives
57
58
59
60

1
2
3 an integrated view of all proteins that are associated with this annotation, not only
4
5 those that belong to the UniRef50 seed proteins (see below).
6
7

8 **(v) EXPANDED PROTEINS**

9
10
11 The ProtoNet tree is started with the representative proteins of UniRef50. The cluster
12
13 view offers a list of the proteins of the cluster. Two levels of expansion are provided:
14
15 the list of proteins according to the UniRef representatives and the complete
16
17 UniProtKB list. On average, the passage from UniRef50 to UniRef90 and from
18
19 UniRef90 to the UniProtKB full list results in a 2.5 folds and 1.8 fold expansion,
20
21 respectively. Cluster A4686503 contains 487 proteins that have a maximal CS for the
22
23 keyword *Cadherins* of CATH homologous superfamily (CS=0.767). This cluster is
24
25 expanded to 2349 proteins. Similarly, the expanded list of proteins can be
26
27 conveniently viewed via PANDORA (see “*Integration of Annotation Sources*”). For
28
29 example, 557 proteins in the ProtoNet 6.0 database are annotated *Cadherins*
30
31 according to the CATH homologous superfamily, but using PANDORA, this list is
32
33 expanded to a total of 2298 proteins.
34
35
36
37
38

39 **(vi) PHYLOGENETIC TREE VIEWER**

40
41
42 The user can select one or several organisms and have the branches in the ProtoNet
43
44 tree that include the selected organisms highlighted. Navigation through the selection
45
46 of complete proteomes is illustrated in Figure 3. It is shown for a few selected
47
48 mammals (human, mouse, rat). Only branches that include proteins from the *selected*
49
50 organisms become visible, though all ‘faded’ clusters can still be analyzed. In Figure
51
52 3, the indicated cluster (Cluster ID 4201544) contains 310 proteins. The number of
53
54 proteins that is covered by the selected proteins is listed (Figure 3). At any stage the
55
56 user can reset or remove or change the taxonomical based selection.
57
58
59
60

(vii) COMPARING VERSIONS

The user may select to navigate each of the main releases of ProtoNet. Maintenance of the different versions allows assessing the changes in the clusters along the constant growth in protein sequences. For example, with the same threshold of PL=90 and LT=1 there are 5,245 and 74,446 stable clusters by ProtoNet versions of SwissProt 40.28 and UniProt 8.1, respectively.

(viii) ADVANCED NAVIGATION

The advance mode provides additional control for the user on the parameters of the visualization that concern: (i) the ProtoBrowser. (ii) ProtoNet tree condensation. The user can choose to activate the ProtoBrowser at a different resolution. While the simplified mode (Figure 3, upper panel) shows two levels above and below the observed cluster (marked in red font in the tree, Figure 3), in 'advanced mode', the number of presented surrounding tree layers is a user-selected parameter. By moving up the tree one observes how the cluster grows in size and becomes more diverse.

The user can change the tree resolution by modifying the parameters of the tree condensation protocol (see "*Selecting Stable Clusters*"). Such change of parameters turns a binary tree to a non-binary tree, and some browsing options help the user in following this modification.

Other capacities of the 'advanced mode' reflect certain intrinsic properties of the ProtoNet Tree. The user can retrieve the ProtoNet clusters at a specific ProtoLevel (PL) (Figure 3, lower panel). This determines the number of clusters to be presented but it also (indirectly) allows the user to focus on a PL that is maximally enriched by proteins with unknown function. While a careful biological interpretation of the

1
2
3 ProtoNet 6.0 clusters is beyond the scope of this paper, we should note a significant
4
5 explosion of proteins of unknown function that appears at PL above 90.
6
7

8
9 Additional queries address the connectivity of selected proteins in the tree. In
10
11 particular, one can get the lowest common cluster of any two proteins. Search for the
12
13 appearance of a specific protein within a cluster, search for all the clusters that are
14
15 associated with a selected keyword and more.
16
17

18 **A TEST CASE - METAGENOME TO FUNCTION**

19

20
21 Global Ocean Sampling (GOS) sequences is a huge collection of (mostly)
22
23 unidentified marine metagenome sequences that covers nearly all known prokaryotic
24
25 protein families (33). We now illustrate a test case of one of hypothetical protein
26
27 GOS_6351915.
28
29

30
31 Applying the ProtoNet option '*Paste your new sequence*' in basic mode with default
32
33 parameters finds this sequence in cluster 4033656 (26 proteins, 5 named 'predicted
34
35 protein' and additional 2 proteins named 'putative') all of which belong to InterPro
36
37 entry of "Longin" and to additional keywords that specifies the relevance to SNARE-
38
39 like (based on SCOP). However, upon moving up the tree to a larger cluster with 107
40
41 proteins (Cluster ID 4312270), the dominating keyword (ProtoName) is changed to
42
43 InterPro IPR016444: Synaptobrevin that metazoa/fungi. The taxonomy of the merged
44
45 cluster includes only metazoa and fungi (excluding green plans).
46
47
48

49
50 Activating the 'advanced mode' for a condensed tree (LT threshold=10) indicates that
51
52 GOS_6351915 sequence belongs to a larger cluster (213 proteins, Root) where the
53
54 most significant annotation (Cluster ID 4446624 and CS=0.965) is from CATH
55
56 topology of Beta-Lactamase and homologous group of CATH 3.30.450.50. Analyzing
57
58
59
60

1
2
3 this very stable cluster via PANDORA shows that the dominating features are
4
5 *membrane* and *coiled coil*. The significant p-value for other functional annotations
6
7 such as v-SNARE, trafficking, synaptic vesicles, ER and golgi confirm that
8
9 GOS_6351915 sequence is a genuine member of the SNARE family. We postulate
10
11 with high confidence that this sequence is a Synaptobrevin-like protein that is
12
13 probably derived from the unicellular species of marine-centric diatom.
14
15

16 17 **MAINTENANCE AND UPDATING**

18
19
20 ProtoNet will be updated once a year. A partition in UniProt to the sections of
21
22 UniProt/Swiss-Prot and UniProt/TrEMBL will be implemented. This will allow users
23
24 to focus, as needed, on each UniProt section, separately. Future ProtoNet releases will
25
26 incorporate additional annotation resources from KEGG, STRING, OMIM and GO
27
28 evidence codes. To provide the user with control over the confidence level,
29
30 annotations evidence (e.g., experimental, inferred) will be added for each protein in
31
32 our database.
33
34
35

36
37 ProtoNet 6.0 had also incorporated few fundamental technical improvements in the
38
39 automation, database design and technologies. These improvements concern the
40
41 automation for the future updates and releases.
42
43
44

45 **ACKNOWLEDGEMENTS**

46
47
48 We thank the ProtoNet team and especially Yaniv Loewenstein for his support in
49
50 establishing the system's performance. We thank Eden Dror for his help in
51
52 maintaining and improving the database of ProtoNet. We thank the fellowship support
53
54 by the Sudarsky Center for Computational Biology (SCCB) to NR. This study is
55
56
57
58
59
60

1
2
3 partially funded by the EU FRVII Prospects consortium and the Israel Science
4
5 Foundation ISF 592/07.
6
7

8 REFERENCES

- 9
10
11
12
13
14 1. Sasson, O., Vaaknin, A., Fleischer, H., Portugaly, E., Bilu, Y., Linial, N. and
15 Linial, M. (2003) ProtoNet: hierarchical classification of the protein space. *Nucleic*
16 *Acids Res*, 31, 348-352.
17
18
19
20
21 2. The UniProt Consortium (2011) The Universal Protein Resource (UniProt) in
22 2010. *Nucleic Acids Res*, 38, D142-148.
23
24
25
26 3. Loewenstein, Y., Raimondo, D., Redfern, O.C., Watson, J., Frishman, D.,
27 Linial, M., Orengo, C., Thornton, J. and Tramontano, A. (2009) Protein function
28 annotation by homology-based inference. *Genome Biol*, 10, 207.
29
30
31
32 4. Fleischmann, W., Moller, S., Gateau, A. and Apweiler, R. (1999) A novel
33 method for automatic functional annotation of proteins. *Bioinformatics*, 15, 228-233.
34
35
36
37 5. Friedberg, I. (2006) Automated protein function prediction--the genomic
38 challenge. *Brief Bioinform*, 7, 225-242.
39
40
41
42 6. Brown, D.P., Krishnamurthy, N. and Sjolander, K. (2007) Automated protein
43 subfamily identification and classification. *PLoS Comput Biol*, 3, e160.
44
45
46
47 7. Watson, J.D., Sanderson, S., Ezersky, A., Savchenko, A., Edwards, A.,
48 Orengo, C., Joachimiak, A., Laskowski, R.A. and Thornton, J.M. (2007) Towards
49 fully automated structure-based function prediction in structural genomics: a case
50 study. *J Mol Biol*, 367, 1511-1522.
51
52
53
54
55
56
57
58
59
60

- 1
2
3 8. Pazos, F. and Sternberg, M.J. (2004) Automated prediction of protein function
4 and detection of functional sites from structure. Proc Natl Acad Sci U S A, 101,
5 14754-14759.
6
7
- 8
9
10 9. Finn, R.D., Mistry, J., Tate, J., Coggill, P., Heger, A., Pollington, J.E., Gavin,
11 O.L., Gunasekaran, P., Ceric, G., Forslund, K. et al. (2010) The Pfam protein families
12 database. Nucleic Acids Res, 38, D211-222.
13
14
- 15
16 10. Letunic, I., Doerks, T. and Bork, P. (2009) SMART 6: recent updates and new
17 developments. Nucleic Acids Res, 37, D229-232.
18 <http://www.ncbi.nlm.nih.gov/pubmed/18978020>
19
20
- 21
22 11. Portugaly, E., Linial, N. and Linial, M. (2007) EVEREST: a collection of
23 evolutionary conserved protein domains. Nucleic Acids Res, 35, D241-246.
24
25
- 26
27 12. Mi, H., Dong, Q., Muruganujan, A., Gaudet, P., Lewis, S. and Thomas, P.D.
28 (2010) PANTHER version 7: improved phylogenetic trees, orthologs and
29 collaboration with the Gene Ontology Consortium. Nucleic Acids Res, 38, D204-210.
30
31
- 32
33 13. Yeats, C., Lees, J., Carter, P., Sillitoe, I. and Orengo, C. (2011) The Gene3D
34 Web Services: a platform for identifying, annotating and comparing structural
35 domains in protein sequences. Nucleic Acids Res, 39, W546-550.
36
37
- 38
39 14. Sigrist, C.J., Cerutti, L., de Castro, E., Langendijk-Genevaux, P.S., Bulliard,
40 V., Bairoch, A. and Hulo, N. (2010) PROSITE, a protein domain database for
41 functional characterization and annotation. Nucleic Acids Res, 38, D161-166.
42
43
- 44
45 15. Attwood, T.K., Bradley, P., Flower, D.R., Gaulton, A., Maudling, N.,
46 Mitchell, A.L., Moulton, G., Nordle, A., Paine, K., Taylor, P. et al. (2003) PRINTS
47 and its automatic supplement, prePRINTS. Nucleic Acids Res, 31, 400-402.
48
49
- 50
51 16. Bru, C., Courcelle, E., Carrere, S., Beausse, Y., Dalmar, S. and Kahn, D.
52 (2005) The ProDom database of protein domain families: more emphasis on 3D.
53 Nucleic Acids Res, 33, D212-215.
54
55
56
57
58
59
60

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
17. Henikoff, J.G., Greene, E.A., Pietrokovski, S. and Henikoff, S. (2000) Increased coverage of protein families with the blocks database servers. *Nucleic Acids Res*, 28, 228-230.
18. Wu, C.H., Nikolskaya, A., Huang, H., Yeh, L.S., Natale, D.A., Vinayaka, C.R., Hu, Z.Z., Mazumder, R., Kumar, S., Kourtesis, P. et al. (2004) PIRSF: family classification system at the Protein Information Resource. *Nucleic Acids Res*, 32, D112-114.
19. Marchler-Bauer, A., Lu, S., Anderson, J.B., Chitsaz, F., Derbyshire, M.K., DeWeese-Scott, C., Fong, J.H., Geer, L.Y., Geer, R.C., Gonzales, N.R. et al. (2010) CDD: a Conserved Domain Database for the functional annotation of proteins. *Nucleic Acids Res*, 39, D225-229.
20. Hunter, S., Apweiler, R., Attwood, T.K., Bairoch, A., Bateman, A., Binns, D., Bork, P., Das, U., Daugherty, L., Duquenne, L. et al. (2009) InterPro: the integrative protein signature database. *Nucleic Acids Res*, 37, D211-215.
21. Andreeva, A., Howorth, D., Chandonia, J.M., Brenner, S.E., Hubbard, T.J., Chothia, C. and Murzin, A.G. (2008) Data growth and its impact on the SCOP database: new developments. *Nucleic Acids Res*, 36, D419-425.
22. Cuff, A.L., Sillitoe, I., Lewis, T., Redfern, O.C., Garratt, R., Thornton, J. and Orengo, C.A. (2009) The CATH classification revisited--architectures reviewed and new ways to characterize structural divergence in superfamilies. *Nucleic Acids Res*, 37, D310-314.
23. Pearl, F., Todd, A., Sillitoe, I., Dibley, M., Redfern, O., Lewis, T., Bennett, C., Marsden, R., Grant, A., Lee, D. et al. (2005) The CATH Domain Structure Database and related resources Gene3D and DHS provide comprehensive domain family information for genome analysis. *Nucleic Acids Res*, 33, D247-251.

- 1
2
3 24. Barrell, D., Dimmer, E., Huntley, R.P., Binns, D., O'Donovan, C. and
4 Apweiler, R. (2009) The GOA database in 2009--an integrated Gene Ontology
5 Annotation resource. *Nucleic Acids Res*, 37, D396-403.
6
7
8
9
10 25. Suzek, B.E., Huang, H., McGarvey, P., Mazumder, R. and Wu, C.H. (2007)
11 UniRef: comprehensive and non-redundant UniProt reference clusters.
12 *Bioinformatics*, 23, 1282-1288.
13
14
15
16 26. Kaplan, N., Sasson, O., Inbar, U., Friedlich, M., Fromer, M., Fleischer, H.,
17 Portugaly, E., Linial, N. and Linial, M. (2005) ProtoNet 4.0: a hierarchical
18 classification of one million protein sequences. *Nucleic Acids Res*, 33, D216-218.
19
20
21
22 27. Sasson, O., Kaplan, N. and Linial, M. (2006) Functional annotation prediction:
23 all for one and one for all. *Protein Sci*, 15, 1557-1562.
24
25
26
27
28 28. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W.
29 and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of
30 protein database search programs. *Nucleic Acids Res*, 25, 3389-3402.
31
32
33
34 29. Loewenstein, Y., Portugaly, E., Fromer, M. and Linial, M. (2008) Efficient
35 algorithms for accurate hierarchical clustering of huge datasets: tackling the entire
36 protein space. *Bioinformatics*, 24, i41-49.
37
38
39
40 30. Petryszak, R., Kretschmann, E., Wieser, D. and Apweiler, R. (2005) The
41 predictive power of the CluSTr database. *Bioinformatics*, 21, 3604-3609.
42
43
44
45
46 31. Bairoch, A. (2000) The ENZYME database in 2000. *Nucleic Acids Res*, 28,
47 304-305.
48
49
50
51 32. Rappoport, N., Fromer, M., Schweiger, R. and Linial, M. (2010) PANDORA:
52 analysis of protein and peptide sets through the hierarchical integration of annotations.
53 *Nucleic Acids Res*, 38, W84-89.
54
55
56
57
58
59
60

1
2
3 33. Yooseph, S., Sutton, G., Rusch, D.B., Halpern, A.L., Williamson, S.J.,
4 Remington, K., Eisen, J.A., Heidelberg, K.B., Manning, G., Li, W. et al. (2007) The
5 Sorcerer II Global Ocean Sampling expedition: expanding the universe of protein
6 families. PLoS Biol, 5, e16.
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For Peer Review

TABLES

Table 1.

Database	Clustering	CS	Specificity	Sensitivity
UniRef90	MC-ProtoNet	0.89	0.96	0.92
	Single Linkage	0.78	0.93	0.24
	ProtoNet 4.0	0.75	0.94	0.79
UniRef50	MC-ProtoNet	0.88	0.96	0.91
	Single Linkage	0.72	0.91	0.79
SwissProt	MC-ProtoNet	0.90	0.96	0.94
	Single Linkage	0.81	0.90	0.91

Clustering performance evaluation based on Pfam keywords. Tests were performed on UniRef90 (1.8M), UniRef50 (960K) and SwissProt (220K). CS, correspondence score.

FIGURES LEGENDS

Figure 1

ProtoNet clusters following pruning at selected thresholds. (A) A scheme of the binary tree following low and high condensations ($LT \geq x$ and $LT \geq y$). The high level of compression ($LT=5$) results in a smaller number of stable clusters. (B) Each panel represents a cluster summary according to a selected threshold (LT). Low ($LT=0.2$) and high condensation level ($LT=5.0$) differ in their cluster size and other statistical properties. Details on the cluster size, depth (by ProtoLevel, PL), the number of hypothetical proteins, solved structures in the PDB database and more are shown.

Figure 2

The contribution of annotation types to ProtoNet clusters. (A) About 40 annotation types that cover different aspects of function are included. Some of the minor annotation sources were combined and depicted as 'others'. (B) The major annotation types and their coverage as measured by the fraction of proteins that are assigned with the indicated annotation type are listed. In ProtoNet 6.0, a total of 143,849,828 annotations (74,416,565 without taxonomy) is associated with the ~9 million protein sequences.

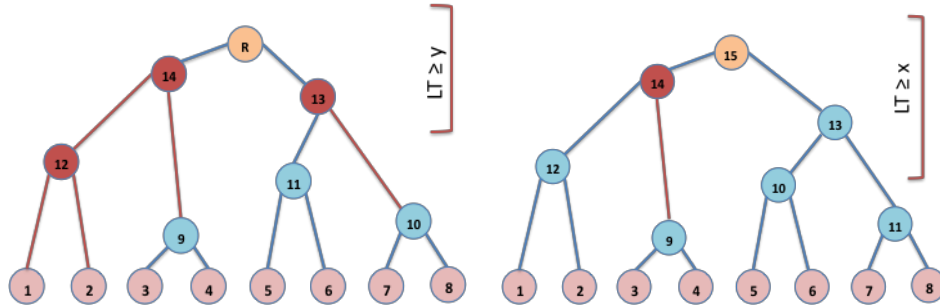
Figure 3

ProtoNet cluster page and a tree viewer in simplified and advanced modes. (Top) From the cluster page (Cluster ID 4201544) the user can focus on the ProtoName and

1
2
3 the collection of additional high quality annotations that are associated with this
4 cluster. The number of proteins from the selected organisms is indicated with a
5 framed T-symbol (for Taxonomy). Similarity, clusters that include proteins with 3D
6 solved structures as marked by a symbol for PDB. Each cluster provides a short
7 summary as a popup box with the number of proteins and the appearance of pre-
8 selected organisms. The red edges in the tree indicate the branches that include the
9 selected organisms. All other branched are faded. (Bottom) Using the advanced mode,
10 the number of clusters in the ProtoNet tree is listed according to the predetermined LT
11 and PL values. There are several sorting options according to the cluster size and the
12 properties of the tree. An interactive use of the condensation levels allows inspecting
13 the near vicinity of a subjected cluster in the ProtoNet hierarchy.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Figure 1

A



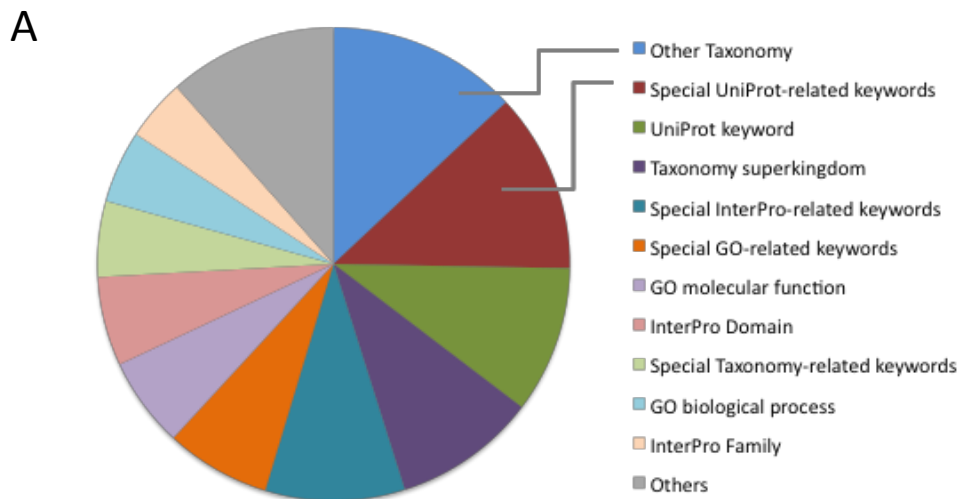
B

LT=0.2

LT=5

Size	387 proteins	636 proteins
Fraction of UniProtKB/Swiss-Prot proteins	46.2%	32.2%
Fraction of UniProtKB/TrEMBL proteins	53.7%	67.7%
Number of Solved structures (in PDB)	16 proteins, 70 PDB's	17 proteins, 74 PDB's
Number of proteins without Prosite ID	68	176
Number of Hypothetical proteins	6	23
Number of Fragments	0	0
Average length and standard deviation of proteins	398 ± 63.2	376.4 ± 132.1
Number of clusters at ProtoLevel at creation	351879	get data
Number of leaves at ProtoLevel at creation	195678	get data
ProtoLevel at creation	0.798639	48.9337
ProtoLevel at termination	1.54009	60.7214
Lifetime	0.313586	6.3886

Figure 2



B

Major annotation sources	Coverage (% of proteins)	Amount of different annotations
ENZYME (6/09)	11	5,190
SMART (6.0)	14	720
GENE3D (6.1.0)	27	1,024
CATH (3.2.0)	27	3,338
SCOP (1.75)	35	7,821
GO (Gene Ontology) (1.7)	52	27,050
UniProtKB/SWP (15.4)	63	949
PFAM (23.0)	73	10,640
InterPro (21.0)	77	18,638
NCBI Taxonomy (6/09)	100	442,867

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58

Cluster 4201544
 Cluster Name: Cadherins
 Additional names (23):
 • CATH architecture: Sandwich
 • CATH class: Mainly Beta
 • CATH topology: Immunoglobulin-like
 • GO biological process: Biological adhesion
 • GO cellular component: Cell part
 • GO molecular function: Binding
 • InterPro Domain: Cadherin-like
 • SCOP superfamily: Cadherin-like
 • Taxonomy kingdom: Metazoa
 • Taxonomy superkingdom: Eukaryota
 • UniProt keyword: Calcium

General Info
 Size: 310 proteins
 Fraction of UniProtKB/Swiss-Prot proteins: 17.4%
 Fraction of UniProtKB/TrEMBL proteins: 82.5%
 Number of Solved structures (in PDB): 4 proteins, 5 PDB's
 Number of proteins without Prosite ID: 0
 Number of Hypothetical proteins: 64
 Number of Fragments: 0
 Average length and standard deviation of proteins: 1499.7 ± 1197.2
 Fraction of "parent" cluster: 94%

Cluster tree
 Cluster: 4118312
 Cluster name: Cadherins (CATH homologous superfamily)
 Cluster contains 1 selected org.:
 1) Homo sapiens (2)
 Size: 16
 Parent: 4201544

Version 6.0 Type of classification "Arithmetic UniProt 15.4"

Clusters at ProtoLevel "90" (without leaves).
 Method of Condensation: lifetime. Threshold: 5.

Amount of clusters: 154134
 Sorted by: Cluster Size (descending)
 Back to sorting by default

View: Get leaves at this ProtoLevel [view]

No.	Cluster ID	Size	ProtoLevel at creation (birthtime)	ProtoLevel at termination (deathtime)	Lifetime
	A4653887	5660	88.6527	93.9894	5.3367
	A4667801	2576	89.857	94.9953	5.1383
	A4653210	2280	88.50	94.827	6.2670
4	A4532761	2131	73.8441	93.9587	20.1146
5	A4625535	1847	85.3756	91.4418	6.0662
6	A4623007	1839	85.1443	90.9315	5.7872
7	A4637348	1725	86.7717	93.5379	6.7662