**Supplementary Figures and Tables**

**The EFI Web Resource for Genomic Enzymology Web Tools: Leveraging Protein, Genome, and Metagenome Databases to Discover Novel Enzymes and Metabolic Pathways**

Rémi Zallot[1], Nils O. Oberg[1], and John A. Gerlt[*,1, 2, 3]

[1]Institute for Genomic Biology, [2]Department of Biochemistry, and [3]Department of Chemistry, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States.

Corresponding author: John A. Gerlt (j-gerlt@illinois.edu).
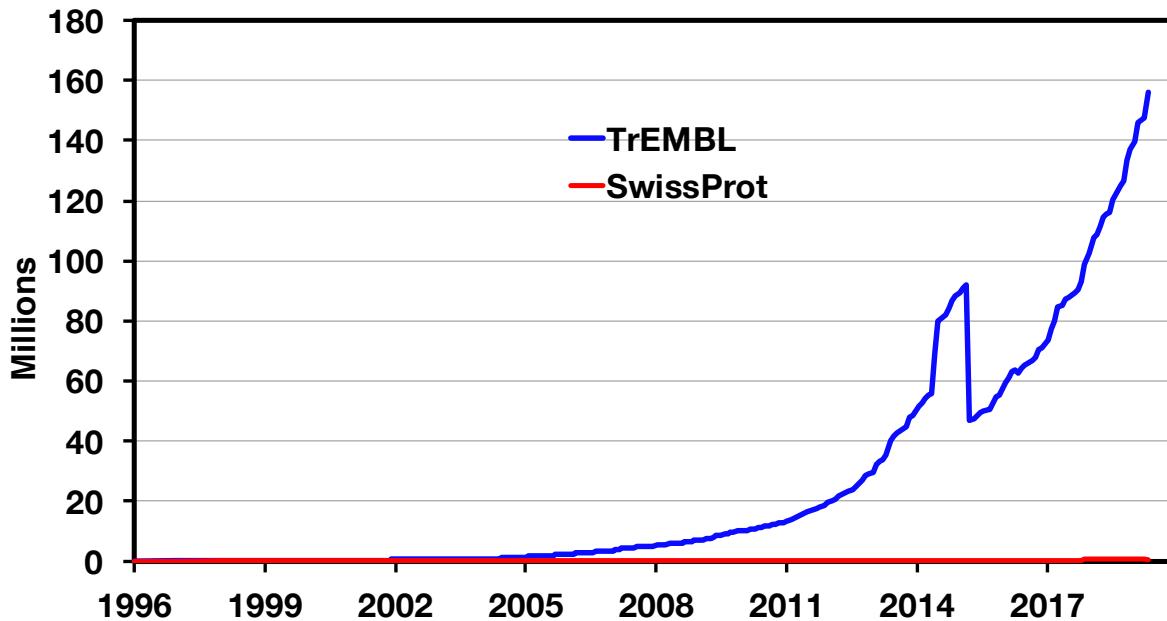
Phone: +1-217-244-7414

**Figure S1**. **The growth of the UniProtKB database**. The web tools use the UniProtKB database for sequences used to generate SSNs and bioinformatic data included as node attributes in the SSNs. UniProtKB is the aggregate of the UniProtKB/TrEMBL database that contains computationally annotated entries [156,077,686 in Release 2019_04 (08-May-2019)] and the UniProtKB/SwissProt database that contains manually curated entries (560,118 in Release 2019_04). The decrease in 2015 is the result of archiving sequences from redundant proteomes in UniParc to manage the growth of the database.

**Figure S2**. **The home pages for the web tools**. **Panel A**, EFI-EST for generating sequence similarity networks (SSNs), accessible at https://efi.igb.illinois.edu/efi-est/. **Panel B**, EFI-GNT for collecting and analyzing genome context for bacterial, archaeal, and fungal proteins in SSNs, accessible at https://efi.igb.illinois.edu/efi-gnt/. **Panel C**, EFI-CGFP for mapping metagenome abundance to SSNs clusters, accessible at https://efi.igb.illinois.edu/efi-cgfp/.

**Figure S3**. **The EFI-EST pages for generating SSNs**. **Panel A**, Option A, a user-provided sequence is used as the query for a BLAST search of the UniProt database to collect homologues. **Panel B**, Option B, one or more user-specified protein families (Pfam, InterPro, and/or Pfam clans) is used to generate the SSN. **Panel C**, Option C, a user-provided FASTA file provides the sequences to generate the SSN. **Panel D**, Option C, a user-provided list of accession IDs (UniProt and/or NCBI) specifies the sequences used to generate the SSN. **Panel E**, Color SSN Utility, unique colors and numbers are assigned to the clusters in a user-provided SSN.

**Figure S4**. **The sequence of steps in generating an SSN with EFI-EST, using the glycyl radical enzyme superfamily (IPR004184) as an example**. **Panel A**, the family identifier is specified, and the user chooses the database (UniProt, UniRef90, or UniRef50) for generating the SSN (red arrow). **Panel B**, the "Dataset Completed" page that provides histograms and boxplots that are used for selecting in minimum alignment score threshold for generating the initial SSN; see text for how the histograms and plots are used to select the alignment score. **Panel C**, the SSN "Finalization" tab for entering the "Alignment Score Threshold" (blue arrow) and "Minimum"/"Maximum" length filters (green arrow). **Panel D**, the "Download Network Files" page that provides access to the full and representative node SSNs. The details are described in the text.

**Figure S5**. **A comparison of the SSNs for IPR004184 generated using the UniProt (Panel A), UniRef90 (Panel B), and UniRef50 (Panel C) databases**. The SSNs were colored with the Color SSNs Utility that assigns a unique color and number to the SSN clusters, with the numbers assigned in order of decreasing number of UniProt IDs in the clusters. Some singleton nodes are colored in the UniRef90 and UniRef50 SSNs—these contain multiple UniProt IDs so are considered clusters. The UniProt SSN (Panel A) is the highest resolution; the UniRef90 SSN (Panel B) provides similar resolution but the file size is significantly less than that for the UniProt SSN. The UniRef50 SSN (Panel B) should be used only for the largest families, with subsequent generation of daughter SSNs for sequences in individual clusters allowing higher resolution analyses. In this example, the UniProt and UniRef90 SSNs contain the same number of clusters.

```
Cluster 1 and Cluster 5, PFL
sp|P09373|PFLB_ECOLI              PDAYGRGRIIGDY      DYA-IACCvSPmIVGK
sp|Q5HJF4|PFLB_STAAC              PDAYGRGRIIGDY      DYG-IACCvSAmTIGK
sp|Q7A7X6|PFLB_STAAN              PDAYGRGRIIGDY      DYG-IACCvSAmTIGK
tr|D6XC58|D6XC58_9ACTN            PDAYGRGRIIGDY      DTA-IACCvSAmAVGR
tr|A0A1Y1WTI3|A0A1Y1WTI3_9FUNG    PDGYGRGRIIGDY      DYG-IACCvSAmRIGK
tr|D8UHK4|D8UHK4_VOLCA            PDGYGRGRIIGDY      DYS-IACCvSAmRVGK
tr|A0A2P6TKX7|A0A2P6TKX7_CHLSO    PDGYGRGRIIGDY      DYG-IACCvSAmRIGK
tr|A0A0A2VZ23|A0A0A2VZ23_BEABA    PDAYGRGRIIGDY      DYA-IACCvSPmIVGK
tr|A0A084AER5|A0A084AER5_LACLC    PDAYSRGRIIGVY      MSC-ISCCvSPLDPEN
tr|J7M1Y6|J7M1Y6_STRP1            PDAYSRGRIIGVY      MSC-ISCCvSPLDPEN
tr|B4U135|B4U135_STREM            PDAYSRGRIIGVY      MSC-ISCCvSPLDPEN
tr|A0A0T8QTS3|A0A0T8QTS3_STREE    PDAYSRGRIIGVY      MSC-ISCCvSPLDPEN
tr|A0A133S3J7|A0A133S3J7_STRMT    PDAYSRGRIIGVY      MSC-ISCCvSPLDPEN


Cluster 7, Choline Trimethylamine Lyase (CutC)
tr|B8JOI2|B8JOI2_DESDA            HALNGGGDSNPGY      DYC-LMGCvEPQKSGR
tr|R4Y5E4|R4Y5E4_KLEPR            HQINGGGDTCPGY      DYC-LMGCvEPQKSGR
tr|A0A0H2QDC9|A0A0H2QDC9_9GAMM    HQINGGGDTCPGY      DYC-LMGCvEPQKSGR
tr|A0A1B7JWB3|A0A1B7JWB3_9GAMM    HQINGGGDTCPGY      DYC-LMGCvEPQKSGR
tr|D1P2A0|D1P2A0_9GAMM            HQINGGGDTCPGY      DYC-LMGCvEPQKSGR
tr|B6XDY0|B6XDY0_9GAMM            HQINGGGDTCPGY      DYC-LMGCvEPQKSGR


Cluster 2, "PFL"
tr|A0A0A2VZ96|A0A0A2VZ96_BEABA    NMTSGDAHLAVNF      DYA-AIGCIETAVGGK
tr|A0A0U0K2C1|A0A0U0K2C1_STREE    KMNSGDAHLAVNY      DYS-AIGCvETAVPGK
tr|A0A1L8WPN6|A0A1L8WPN6_9ENTE    KMNSGDAHLAVNY      DYS-AIGCvETAVPGK
tr|A0A1T4P4Y5|A0A1T4P4Y5_9ENTE    NITSGDAHIAVSY      NYS-AIGCvETAVPGK
tr|A0A1H7XPM2|A0A1H7XPM2_9LACO    NITSGDGHIAVNY      NYS-AIGCvETAIPGK


Cluster 3, "PFL"
tr|A0A381GG45|A0A381GG45_CITAM    QTDKGQGHIIIDY      DYA-VVGCvELSIPGR
tr|A0A0R2FPV6|A0A0R2FPV6_9LACO    QTDKGQGHIIMDF      DYG-VVGCvETTIPGK
tr|A0A239SQG5|A0A239SQG5_9STRE    QTDKGQGHIIMDF      DYG-TVGCvEISIPGR
tr|A0A1L8X0W5|A0A1L8X0W5_9ENTE    QTDKGQGHIIMDF      DYA-TVGCvETSIPGK
tr|A0A1A7T0L3|A0A1A7T0L3_ENTFC    QTDKGQGHIIMDF      DYA-TVGCvETSIPGK


Cluster 4, 4-OH Proline Dehydratase
sp|A0A031WDE4|HYPD_CLODI          MEQRAPGHTVCG-      LGG-TSGCvETGCFGK
tr|A0A101F1Q5|A0A101F1Q5_9EURY    MEQRSPGHTAGG-      TSG-VSGCvETGAFGK
tr|A0A2N2ZKH4|A0A2N2ZKH4_9BACT    MEQRAPGHTALD-      EGG-CSGCIEVGAFGK
tr|A0A1W1HBH2|A0A1W1HBH2_9DELT    MEQRAPGHTALD-      EGG-CSGCIETGAFGK
tr|A0A087E582|A0A087E582_9BIFI    MAQRGPGHTVAD-      ESGIASGCvETGTAGK
tr|R7D6G9|R7D6G9_9ACTN            YEQRAGGHTCLGS      HGG-SSGCvETGCWGY


Cluster 6, Glycerol Dehydratase and 1,2-Propanediol Dehydratase
tr|Q8GEZ8|Q8GEZ8_CLOBU            YYYNGVGHVSVDY      DYG-IIGCvEPQKPGK
tr|Q1A666|Q1A666_9FIRM            YFYNGVGHVTVQY      NYN-IIGCvEPQVPGK
tr|A0A1M6ZJ05|A0A1M6ZJ05_9FIRM    YFYNGVGHVTVAY      EYN-IIGCvEPQKAGK
tr|E6MIX1|E6MIX1_9FIRM            YYFGGIGHVCVDY      DWL-PIGCvEPQPQHK
tr|A0A425W4N8|A0A425W4N8_9FIRM    YYYNGIGHVCVDY      DWL-PIGCvEPQPQHK
tr|A8S5K2|A8S5K2_CLOBW            YFYGGVGHVCVDY      SYC-IIGCvEPQCPHK
tr|A0A1H4EE94|A0A1H4EE94_9FIRM    YYYGGVGHVCVDY      NYC-IIGCvEPQCPHK
                                       .              *.
```

**Figure S6**. Partial multiple sequence alignments (MSAs) for the largest seven clusters in SSN$_{240}$ for IPR004184. These first region includes a conserved His in the active sites of characterized dehydratases; the second region includes the conserved Cys-Cys motif that is characteristic of the PFL function. The bold sequence entries are either SwissProt-curated ("sp"; PFL and 4-OH Pro dehydratase) or from the literature (choline trimethylamine lyase, glycerol dehydratase, and 1,2-propanediol dehydratase). Based on their sequences, clusters 2 and 3, that have SwissProt-curated PFL functions (inferred from homology), are predicted to be dehydratases, consistent with their

colocation with 4-OH Pro dehydratase, glycerol dehydratase, and 1,2-propanediol dehydratase in SSN$_{185}$. Choline trimethylamine lyase is neither PFL nor a dehydratase so it is lacking the conserved motifs for these functions.

**Figure S7**. **The sequence of steps in generating GNNs and GNDs with EFI-GNT, using the glycyl radical enzyme superfamily (IPR004184) as an example**. **Panel A**, the SSN is uploaded, with the user specifying the neighborhood size (±10 orfs is the default; red arrow) and query-neighborhood family co-occurrence frequency (20% is the default; the example in the text uses 10%; blue arrow) for generating GNNs. The scripts collect genome neighbors in a ±20 orf window, but the user-specified value is used to generate the GNNs; the GNNs can be recalculated using different neighborhood sizes and co-occurrence frequencies. **Panel B**, the "Results" page that provides the colored SSN (with unique cluster colors and numbers and "Neighbor Pfam Families" and "Neighbor InterPro Families" node attributes; green arrow), the GNNs (SSN cluster-hub nodes with Pfam family spoke nodes, with the SSN cluster-hub nodes colored/numbered; and Pfam family-hub nodes with SSN cluster spoke nodes, with the SSN cluster-spoke nodes colored/numbered; magenta arrow); and access to the GND viewer (orange arrow).

**Figure S8**. The sequence of steps in generating CGFP heatmaps and boxplots with EFI-CGFP, using the glycyl radical enzyme superfamily (IPR004184) as an example. **Panel A**, a colored SSN is uploaded (with unique cluster and singleton numbers to enable mapping of

metagenome abundance to clusters and singletons; red arrow); minimum and maximum length filters are recommended to ensure that the consensus sequences for ShortBRED families used for marker identification are not biased by the presence of fragments (blue arrow). **Panel B**, the "Markers Computation Results" page that allows the user to choose metagenomes for abundance mapping from a library of 380 metagenomes from six body sites from healthy individuals (green arrow). **Panel C**, the "Quantify Results" page that provides (top) heatmaps for metagenome abundance for clusters and singletons in the input SSN and (bottom) boxplots showing quantitative analyses of the metagenome abundances for selected clusters.

**Table S1.  Node Attributes for SSNs Generated by EFI-EST**

| Node Attribute | Options A, B, C with FASTA header reading, D |
|---|---|
| Name | Full network - UniProt or UniRef ID; Rep Node network - UniProt or UniRef ID for the longest sequence in the representative node (seed sequence for CD-Hit).  For domain SSNs, ID:N-terminus:C-terminus |
| Shared name | Full network - UniProt or UniRef ID; Rep Node network - UniProt or UniRef ID for the longest sequence in the representative node (seed sequence for CD-Hit).  For domain SSNs, ID:N-terminus:C-terminus |
| UniRef90 Cluster Size | Number of UniProt IDs in UniRef90 cluster |
| UniRef90 Cluster IDs | List of UniProt IDs in the UniRef90 cluster |
| UniRef50 Cluster Size | Number of UniProt IDs in UniRef50 cluster |
| UniRef50 Cluster IDs | List of UniProt IDs in the UniRef50 cluster |
| Number of IDs in Rep Node[1] | Number of UniProt IDs in the representative node |
| List of IDs in Rep Node[1] | List of UniProt IDs in the representative node |
| Sequence Source | Option A, "INPUT" if input sequence, "BLASTHIT" if identified in BLAST, "FAMILY" if from user-specified user-specified Pfam/InterPro family, "USER+BLASTHIT" if from BLAST and family |
| | Options B, C, and D, "USER" if from user-supplied file, "FAMILY" if from user-specified Pfam/InterPro family, "USER+FAMILY" if from both |
| Query IDs | Options C and D, Input Query ID(s) that identified a UniProt match in the idmapping file |
| Other IDs | Option C, headers for FASTA sequences that could not identify a UniProt match in the idmapping file |
| User IDs in Cluster | Options A, B, and C with UniRef family added and/or rep node SSNs, UniProt IDs for BLASTHITs or user-supplied sequences in metanode |
| Organism | organism genus/genera and species, from UniProt taxonomy.xml |
| Taxonomy ID | NCBI taxonomy identifier(s), from UniProt |
| UniProt Annotation Status | SwissProt - manually annotated; TrEMBL - automatically annotated; from UniProt |
| Description | protein name(s)/annotation(s), from UniProtKB |
| SwissProt Description | protein name(s)/annotation(s), from UniProtKB for SwissProt reviewed entries |
| Sequence Length | number(s) of amino acid residues, from UniProt |
| Cluster ID Sequence Length | Sequence length for Cluster ID in UniRef SSNs ("most informative" sequence in cluster, as designated by UniProt) |
| Gene name | gene name(s) |
| NCBI IDs | RefSeq/GenBank IDs and GI numbers, from UniProt idmapping |
| Superkingdom | domain of life of the organism, from UniProt taxonomy.xml |
| Kingdom | kingdom of the organism, from UniProt taxonomy.xml |

| | |
|---|---|
| Phylum | Phylogenetic phylum of the organism, from UniProt taxonomy.xml |
| Class | Phylogenetic class of the organism, from UniProt taxonomy.xml |
| Order | Phylogenetic order of the organism, from UniProt taxonomy.xml |
| Family | Phylogenetic family of the organism, from UniProt taxonomy.xml |
| Genus | Phylogenetic genus of the organism, from UniProt taxonomy.xml |
| Species | Phylogenetic species of the organism, from UniProt taxonomy.xml |
| EC | EC number, from UniProt |
| PFAM | Pfam family, from UniProt |
| PDB | Protein Data Bank entry, from UniProt |
| InterPro (Domain) | InterPro domain(s), from InterPro |
| InterPro (Family) | InterPro family(ies), from InterPro |
| InterPro (Homologous Superfamily) | InterPro homologous superfamily(ies), from InterPro |
| InterPro (Other) | Other InterPro classes (repeat, site), from InterPro |
| BRENDA ID | BRENDA Database ID, from UniProt |
| CAZY Name | Carbohydrate-Active enZYmes (CAZy) family name(s), from UniProt |
| GO Term | Gene Ontology classification(s), from UniProt |
| KEGG ID | KEGG Database ID, from UniProt |
| PATRIC ID | PATRIC Database ID, from UniProt |
| STRING ID | STRING Database ID, from UniProt |
| HMP Body Site | location(s) of organism(s) in/on the body, if human microbiome organism, spreadsheet from HMP |
| HMP Oxygen | oxygen requirement(s), if human microbiome organism, from HMP |
| P01 gDNA | availability of gDNA(s) at EFI Protein Core, custom |
| Sequence | Option C, Sequence from UniProt database if ID can be located |

| Node Attribute | Option C without FASTA header reading |
|---|---|
| Name | zzznnn, where nnn = number of the sequence in FASTA file |
| Shared Name | zzznnn, where nnn = number of the sequence in FASTA file |
| Description | FASTA Header |
| Sequence Length | Length of sequence in FASTA entry |
| Sequence | Sequence from FASTA entry |

| Additional Node Attributes | Colored SSN (from Colored SSNs utility) |
|---|---|
| Cluster Number | Number assigned to cluster, in order of decreasing number of sequences in the clusters |
| Cluster Sequence Count | Number of sequences in the cluster |
| Node.fillColor | Unique color assigned to cluster, in hexadecimal |
| Singleton Number | Number assigned to singleton |

**Table S2.  Formats for UniProt, NCBI and PDB IDs; FASTA Headers for Option C**

**A.  Formats for UniProt IDs, NCBI IDs, and PDB IDs**

**UniProt IDs**
UniProtKB ID is 6 or 10 alphanumerical characters in the following formats:

```
    1        2        3         4          5          6       7       8          9          10
[O,P,Q]    [0-9] [A-Z,0-9] [A-Z,0-9] [A-Z,0-9] [0-9]
[A-N,R-Z]  [0-9] [A-Z]     [A-Z,0-9] [A-Z,0-9] [0-9]
[A-N,R-Z]  [0-9] [A-Z]     [A-Z,0-9] [A-Z,0-9] [0-9]   [A-Z] [A-Z,0-9] [A-Z,0-9] [0-9]
```

For example:
```
P11444
T2HDW6
A0A0A7PVN6
```

**NCBI RefSeq IDs**
An NCBI RefSeq ID is 2 letters followed by an underscore followed by a series of digits, a period, and one or more digits for the sequence version number, e.g.,
```
WP_016501748.1
NP_708575.1
YP_002409124.1
```

**NCBI UniProt/Swiss-Prot IDs**
An NCBI UniProt/Swiss-Prot ID is the UniProt ID followed by a period and one or more digits for the sequence version number, e.g.,
```
Q31XL1.1
B7LEJ8.1
C4ZZT2.1
```

**NCBI GenBank IDs**
The format for NCBI GenBank IDs is 3 letters followed by five digits, a period, and one or more digits for the sequence version number, e.g.,
```
BAN56663.1
AAC15504.1
BAM38409.1
```

**PDB IDs**
The format for PDB IDs is one digit followed by two letters and a digit/letter:
```
1MDL
1MRA
3UXL
```

**NCBI GI Numbers**
An NCBI GI number (now retired) is a series of digits.

**B. Formats for FASTA headers for Option C**

**UniProt (TrEMBL and SwissProt, respectively; from UniProt BLAST)**
>tr|**R9RJF1**|R9RJF1_PSEAI Mandelate racemase OS=Pseudomonas aeruginosa PE=4 SV=1
>sp|**P11444**|MANR_PSEPU Mandelate racemase OS=Pseudomonas putida GN=mdlA PE=1 SV=1

**NCBI RefSeq (from NCBI BLAST)**
>**WP_016501748.1** mandelate racemase [Pseudomonas putida]

**NCBI UniProt/Swiss-Prot ID (from NCBI BLAST)**
>**Q0TE80.1** RecName: Full=Enolase; AltName: Full=2-phospho-D-glycerate hydro-lyase; AltName: Full=2-phosphoglycerate dehydratase

**NCBI GenBank ID (from NCBI BLAST)**
>**AAA25887.1** mandelate racemase (EC 5.1.2.2) [Pseudomonas putida]

**NCBI PDB ID (from NCBI BLAST)**
>pdb|**1MDR**|A Chain A, The Role Of Lysine 166 In The Mechanism Of Mandelate Racemase From Pseudomonas Putida: Mechanistic And Crystallographic Evidence For Stereospecific Alkylation By (r)-alpha-phenylglycidate

**NCBI GI Number (from NCBI BLAST; now retired)**
>gi|**347012980**| 4-O-methyl-glucuronoyl methylesterase [Myceliophthora thermophila ATCC 42464]

**Option C also accepts FASTA headers in which the IDs (formats described in Option D) immediately follow the ">" symbol**, e.g., the following headers abbreviated from those shown above:

**UniProt**
>**R9RJF1**
>**P11444**

**NCBI RefSeq**
>**WP_016501748.1**

**NCBI UniProt/Swiss-Prot ID**
>**Q0TE80.1**

**NCBI GenBank ID**
>**AAA25887.1**

**NCBI PDB ID**
>**1MDR**

**NCBI GI Number (now retired)**
>**347012980**

**Table S3. Node Attributes for Colored SSNs Generated by EFI-GNT**

| Node Attribute | Options A, B, C with FASTA header reading, D |
|---|---|
| Name | Full network - UniProt or UniRef ID; Rep Node network - UniProt or UniRef ID for the longest sequence in the representative node (seed sequence for CD-Hit). For domain SSNs, ID:N-terminus:C-terminus |
| Shared name | Full network - UniProt or UniRef ID; Rep Node network - UniProt or UniRef ID for the longest sequence in the representative node (seed sequence for CD-Hit). For domain SSNs, ID:N-terminus:C-terminus |
| UniRef90 Cluster Size | Number of UniProt IDs in UniRef90 cluster |
| UniRef90 Cluster IDs | List of UniProt IDs in the UniRef90 cluster |
| UniRef50 Cluster Size | Number of UniProt IDs in UniRef50 cluster |
| UniRef50 Cluster IDs | List of UniProt IDs in the UniRef50 cluster |
| Number of IDs in Rep Node | Number of UniProt IDs in the representative node |
| List of IDs in Rep Node | List of UniProt IDs in the representative node |
| Sequence Source | Option A, "INPUT" if input sequence, "BLASTHIT" if identified in BLAST, "FAMILY" if from user-specified user-specified Pfam/InterPro family, "USER+BLASTHIT" if from BLAST and family |
| | Options B, C, and D, "USER" if from user-supplied file, "FAMILY" if from user-specified Pfam/InterPro family, "USER+FAMILY" if from both |
| Query IDs | Options C and D, Input Query ID(s) that identified a UniProt match in the idmapping file |
| Other IDs | Option C, headers for FASTA sequences that could not identify a UniProt match in the idmapping file |
| User IDs in Cluster | Options A, B, and C with UniRef family added and/or rep node SSNs, UniProt IDs for BLASTHITs or user-supplied sequences in metanode |
| Cluster Number | Number assigned to cluster, in order of decreasing number of sequences in the clusters |
| Cluster Sequence Count | Number of sequences in the cluster |
| Node.fillColor | Unique color assigned to cluster, in hexadecimal |
| Singleton Number | Number assigned to singleton |
| Present in ENA Database? | "true" if UniProt ID was found in an ENA file (see ENA Database Genome ID); otherwise "false" |
| Genome Neighbors in ENA Database? | "true" if ENA file has sequences for query plus neighbors; "false" if ENA file has no neighbors; "n/a" if not present in ENA database |
| ENA Database Genome ID | ENA file used to obtain genome neighbors |
| Neighbor Pfam Families | Pfam IDs of genome neighborhood proteins in the user-specified window and ≤0% query-neighbor co-occurrence |

| Neighbor InterPro Families | InterPro IDs of genome neighborhood proteins in the user-specified window and ≤0% query-neighbor co-occurrence |
|---|---|
| Organism | organism genus/genera and species, from UniProt taxonomy.xml |
| Taxonomy ID | NCBI taxonomy identifier(s), from UniProt |
| UniProt Annotation Stastus | SwissProt - manually annotated; TrEMBL - automatically annotated; from UniProt |
| Description | protein name(s)/annotation(s), from UniProtKB |
| SwissProt Description | protein name(s)/annotation(s), from UniProtKB for SwissProt reviewed entries |
| Sequence Length | number(s) of amino acid residues, from UniProt |
| Gene name | gene name(s) |
| NCBI IDs | RefSeq/GenBank IDs and GI numbers, from UniProt idmapping |
| Superkingdom | domain of life of the organism, from UniProt taxonomy.xml |
| Kingdom | kingdom of the organism, from UniProt taxonomy.xml |
| Phylum | Phylogenetic phylum of the organism, from UniProt taxonomy.xml |
| Class | Phylogenetic class of the organism, from UniProt taxonomy.xml |
| Order | Phylogenetic order of the organism, from UniProt taxonomy.xml |
| Family | Phylogenetic family of the organism, from UniProt taxonomy.xml |
| Genus | Phylogenetic genus of the organism, from UniProt taxonomy.xml |
| Species | Phylogenetic species of the organism, from UniProt taxonomy.xml |
| EC | EC number, from UniProt |
| PFAM | Pfam family, from UniProt |
| PDB | Protein Data Bank entry, from UniProt |
| InterPro (Domain) | InterPro domain(s), from InterPro |
| InterPro (Family) | InterPro family(ies), from InterPro |
| InterPro (Homologous Superfamily) | InterPro homologous superfamily(ies), from InterPro |
| InterPro (Other) | Other InterPro classes (repeat, site), from InterPro |
| BRENDA ID | BRENDA Database ID, from UniProt |
| CAZY Name | Carbohydrate-Active enZYmes (CAZy) family name(s), from UniProt |
| GO Term | Gene Ontology classification(s), from UniProt |
| KEGG ID | KEGG Database ID, from UniProt |
| PATRIC ID | PATRIC Database ID, from UniProt |
| STRING ID | STRING Database ID, from UniProt |
| HMP Body Site | location(s) of organism(s) in/on the body, if human microbiome organism, spreadsheet from HMP |
| HMP Oxygen | oxygen requirement(s), if human microbiome organism, spreadsheet from HMP |
| P01 gDNA | availability of gDNA(s) at EFI Protein Core, custom |

| Node Attribute | Option C without FASTA header reading |
|---|---|
| Name | zzznnn, where nnn = number of the sequence in FASTA file |
| Shared Name | zzznnn, where nnn = number of the sequence in FASTA file |
| Description | FASTA Header |
| Sequence Length | Length of sequence in FASTA entry |
| Present in ENA Database? | "false" |
| Genome Neighbors in ENA Database? | "n/a" |
| ENA Database Genome ID | none |

**Table S4. GNN Node Attributes for SSN Cluster Hub-Nodes and Pfam Family Spoke-Nodes**

| Node Attribute | SSN cluster hub-nodes |
|---|---|
| Shared name | Input SSN cluster number |
| Name | Input SSN cluster number |
| Cluster Number | Input SSN cluster number |
| # of Sequences in SSN Cluster | Total number of sequences in SSN cluster |
| # of Sequences in SSN Cluster with Neighbors | Number of sequences in SSN cluster with neighbors (queriable sequences) |
| Hub Queries with Pfam Neighbors | Summary of number of queriable sequences with a neighbor in the Pfam family |
| Hub Pfam Neighbors | Summary of the total # of Pfam neighbors found by the queriable sequences |
| Hub Average and Median Distances | Summary of average and median distances between the query and neighbors in each Pfam family |
| Hub Co-occurrence and Ratio | Summary of the query-neighbor co-occurrence (decimal value) and ratio (fraction) for each Pfam family |
| Node.fillColor | Hexadecimal color for the SSN cluster in the colored SSN, used by Cytoscape |
| Node.shape | "hexagon", used by Cytoscape |
| Node Size | "70.0", used by Cytoscape |

| Node Attribute | Pfam family spoke-nodes |
|---|---|
| Shared name | Pfam family short name |
| Name | Pfam family short name |
| SSN Cluster Number | SSN Cluster that found neighbors in the Pfam family |
| Pfam | Pfam family number (PFnnnnn) |
| Pfam description | Pfam family description |
| # of Queries with Pfam Neighbors | Number of queriable sequences with a neighbor in the Pfam family |
| # of Pfam Neighbors | Number of Pfam neighbors found by the queriable sequences |
| Query-Accessions | List of SSN cluster queries that found neighbors in the Pfam family |
| Query-Neighbor Accessions | Information about query-neighbor pairs in the Pfam family |
| Query-Neighbor Arrangement | Genome context information for the query-neighbor pairs in the Pfam family |
| Average Distance | Average distance (in ORFs) between the SSN cluster queries and Pfam neighbors |
| Median Distance | Median distance (in ORFs) between the SSN cluster queries and Pfam neighbors |
| Co-occurrence | Decimal value of ratio of queries that found neighbors to queriable sequences |
| Co-occurrence Ratio | Ratio of queries that found neighbors to queriable sequences |
| Node.fillColor | #EEEEEE, grey in hexadecimal, used by Cytoscape |

| | |
|---|---|
| Node.shape | "ellipse", "diamond", or "square"; explained in on-line tutorial, used by Cytoscape |
| Node.size | Co-occurrence * 100, used by Cytoscape |

**Table S5. GNN Node Attributes for Pfam Family Hub-Nodes and SSN Cluster Spoke-Nodes**

| Node Attribute | Pfam family hub-nodes |
|---|---|
| Shared name | Pfam family short name |
| Name | Pfam family short name |
| Pfam | Pfam family number (PFnnnnn) |
| Pfam description | Pfam family description |
| # of Sequences in SSN Cluster | Total number of sequences in SSN cluster |
| # of Sequences in SSN Cluster with Neighbors | Number of sequences in SSN cluster with neighbors (queriable sequences) |
| # of Queries with Pfam Neighbors | Number of queriable sequences with a neighbor in the Pfam family |
| # of Pfam Neighbors | Number of Pfam neighbors found by the queriable sequences |
| Query-Neighbor Accessions | Information about query-neighbor pairs in the Pfam family |
| Query-Neighbor Arrangement | Genome context information for the query-neighbor pairs in the Pfam family |
| Hub Average and Median Distances | Summary of average and median distances between the query and neighbors |
| Hub Co-occurrence and Ratio | Summary of the query-Pfam family co-occurrence (decimal value) and ratio (fraction) |
| Node.fillColor | "#FFFFFF", white in hexadecimal, used by Cytoscape |
| Node.shape | "hexagon", used by Cytoscape |
| Node.size | "70.0", used by Cytoscape |

| Node Attribute | Description - SSN cluster spoke-nodes |
|---|---|
| Shared name | Input SSN cluster number |
| Name | Input SSN cluster number |
| Cluster Number | Input SSN cluster number |
| # of Sequences in SSN Cluster | Total number of sequences in SSN cluster |
| # of Sequences in SSN Cluster with Neighbors | Number of sequences in SSN cluster with neighbors (queriable sequences) |
| # of Queries with Pfam Neighbors | Number of queriable sequences with a neighbor in the Pfam family |
| # of Pfam Neighbors | Number of Pfam neighbors found by the queriable sequences |
| Query-Accessions | List of queries in each SSN cluster that found neigbhors in the Pfam family |
| Query-Neighbor Accessions | Information about query-neighbor pairs in the Pfam family |
| Query-Neighbor Arrangement | Genome context information for the query-neighbor pairs in the Pfam family |
| Average Distance | Average distance (in ORFs) between the SSN cluster queries and Pfam neighbors |

| Median Distance | Median distance (in ORFs) between the SSN cluster queries and Pfam neighbors |
|---|---|
| Co-occurrence | Decimal value of ratio of queries that found neighbors to queriable sequences |
| Co-occurrence Ratio | Ratio of queries that found neighbors to queriable sequences |

**Table S6.  Additional Node Attributes for SSNs Generated by EFI-CGFP**

| Additional Node Attribute | SSNs with Marker Identification Results |
|---|---|
| Seed Sequences | ID for the (meta)node that contains the UniProt ID for the seed sequence for a ShortBRED family |
| Seed Sequence Cluster(s) | ID of ShortBRED family seed sequence to which the (meta)node contributes family members |
| Marker Types | "true", "quasi", or "junction" |
| Number of Markers | Number of markers identified for ShortBRED family seed sequence |

| Additional Node Attribute | SSNs with Metagenome Abundance Quantify Results |
|---|---|
| Metagenomes Identified by Markers | For ShortBRED family seed sequences, names of metagenome datasets identified by its markers |
| Metagenomes Identified by CD-HIT Family | For IDs that contribute to ShortBRED (CD-HIT) family seed sequences,  names of metagenome datasets identified with the seed sequence markers |