**SUPPLEMENTARY FIGURES**
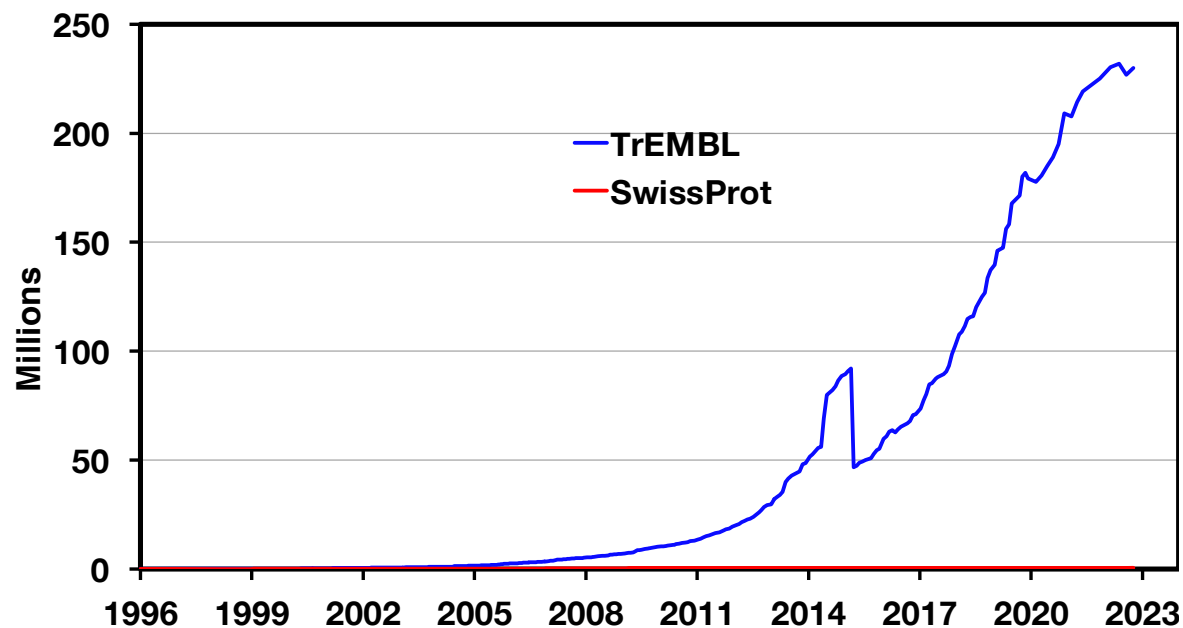
**EFI-EST, EFI-GNT, and EFI-CGFP:  Enzyme Function Initiative (EFI) Web Resource for Genomic Enzymology Tools**

Nils Oberg[1], Rémi Zallot[2,3], and John A. Gerlt[1,4,5*]

[1]Carl R. Woese Institute for Genomic Biology, [4]Department of Biochemistry, and [5]Department of Chemistry, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States.

[2]Department of Chemistry, [3]Manchester Institute of Biotechnology, The University of Manchester, 131 Princess Street, Manchester M1 7DN, UK

**Supplementary Figure S1**. **Growth of the UniProtKB database**. UniProtKB is the aggregate of the UniProtKB/TrEMBL database that contains computationally annotated entries [229,928,140 in Release 2022_04 (12-October-2022)] and the UniProtKB/SwissProt database that contains manually curated entries [(568,363 in Release 2022_04 (12-October-2022)]. The decrease in the number of entries in 2015 and the more modest increase in recent months is the explained by the inclusion of entries from reference proteomes to manage the growth of the database.

**Supplementary Figure S2. Panel A**, **Taxonomy Tool**. **Panel B**, **Filter by Taxonomy** option in the **Taxonomy Tool**, showing **Preselected conditions**. **Panel C**, **Filter by Taxonomy** option in the **Taxonomy Tool**, showing **Add Taxonomy Category**.

**Supplementary Figure S3. Taxonomy Sunburst.**

# A — All Categories



Dataset Summary | Taxonomy Sunburst | Length Histograms

The taxonomy distribution for the UniProt IDs identified as members of the input list of families is displayed.

The UniRef90 and UniRef50 clusters containing the UniProt IDs in the sunburst are identified using the lookup table provided by UniProt/UniRef. These UniRef90 and UniRef50 clusters may contain UniProt IDs from other families; in addition, the UniRef90 and UniRef50 clusters at a selected taxonomy category may contain UniProt IDs from other categories. This results from conflation of UniProt IDs in UniRef90 and UniRef50 clusters that share ≥90% and ≥50% sequence identity, respectively.

The numbers of UniProt IDs, UniRef90 cluster IDs, and UniRef50 cluster IDs for the selected category are displayed.

The sunburst is interactive, providing the ability to zoom to a selected taxonomy category by clicking on that category; clicking on the center circle will return the display to the next highest rank.

Root
Superkingdom
Kingdom
Phylum
Class
Order
Family
Genus
Species

21,636 UNIPROT, 6,938 UNIREF90, 1,199 UNIREF50 IDS

Lists of UniProt, UniRef90, and UniRef50 IDs and FASTA-formatted sequences can be downloaded.

The UniProt, UniRef90, or UniRef50 IDs can be transferred to the Accession IDs option of EFI-EST to generate an SSN. The Accession IDs option provides both Filter by Family and Filter by Taxonomy that should be used to remove internal UniProt IDs from UniRef90 or UniRef50 clusters that are not members of the selected families and/or taxonomy category.

The lists also can be transferred to the GND-Viewer to obtain GNDs.

ID type: ● UniProt ○ UniRef90 ○ UniRef50

Prepare ID Download | Prepare FASTA Download | Transfer to EFI-EST | Transfer to EFI-GND Viewer

# B — Superkingdom Bacteria



Dataset Summary | Taxonomy Sunburst | Length Histograms

The taxonomy distribution for the UniProt IDs identified as members of the input list of families is displayed.

The UniRef90 and UniRef50 clusters containing the UniProt IDs in the sunburst are identified using the lookup table provided by UniProt/UniRef. These UniRef90 and UniRef50 clusters may contain UniProt IDs from other families; in addition, the UniRef90 and UniRef50 clusters at a selected taxonomy category may contain UniProt IDs from other categories. This results from conflation of UniProt IDs in UniRef90 and UniRef50 clusters that share ≥90% and ≥50% sequence identity, respectively.

The numbers of UniProt IDs, UniRef90 cluster IDs, and UniRef50 cluster IDs for the selected category are displayed.

The sunburst is interactive, providing the ability to zoom to a selected taxonomy category by clicking on that category; clicking on the center circle will return the display to the next highest rank.

Superkingdom
Kingdom
Phylum
Class
Order
Family
Genus
Species

SUPERKINGDOM: BACTERIA
20,870 UNIPROT, 6,430 UNIREF90, 1,073 UNIREF50 IDS

Lists of UniProt, UniRef90, and UniRef50 IDs and FASTA-formatted sequences can be downloaded.

The UniProt, UniRef90, or UniRef50 IDs can be transferred to the Accession IDs option of EFI-EST to generate an SSN. The Accession IDs option provides both Filter by Family and Filter by Taxonomy that should be used to remove internal UniProt IDs from UniRef90 or UniRef50 clusters that are not members of the selected families and/or taxonomy category.

The lists also can be transferred to the GND-Viewer to obtain GNDs.

ID type: ● UniProt ○ UniRef90 ○ UniRef50

Prepare ID Download | Prepare FASTA Download | Transfer to EFI-EST | Transfer to EFI-GND Viewer

# C — Kingdom Terrabacteria Group



Dataset Summary | Taxonomy Sunburst | Length Histograms

The taxonomy distribution for the UniProt IDs identified as members of the input list of families is displayed.

The UniRef90 and UniRef50 clusters containing the UniProt IDs in the sunburst are identified using the lookup table provided by UniProt/UniRef. These UniRef90 and UniRef50 clusters may contain UniProt IDs from other families; in addition, the UniRef90 and UniRef50 clusters at a selected taxonomy category may contain UniProt IDs from other categories. This results from conflation of UniProt IDs in UniRef90 and UniRef50 clusters that share ≥90% and ≥50% sequence identity, respectively.

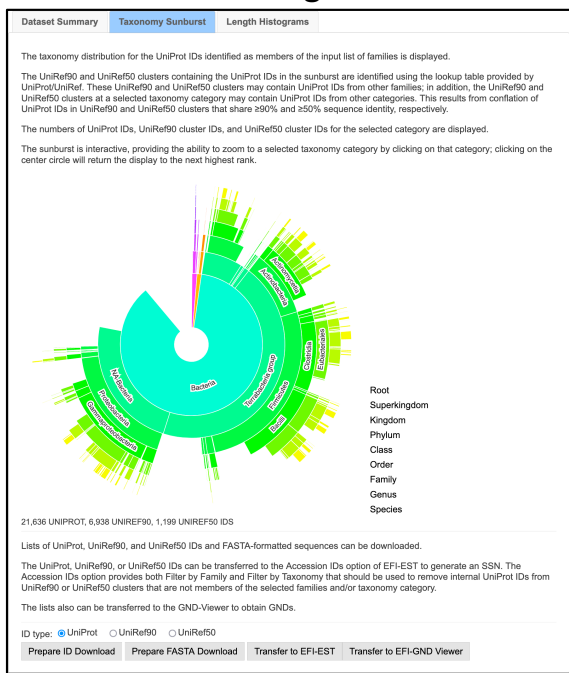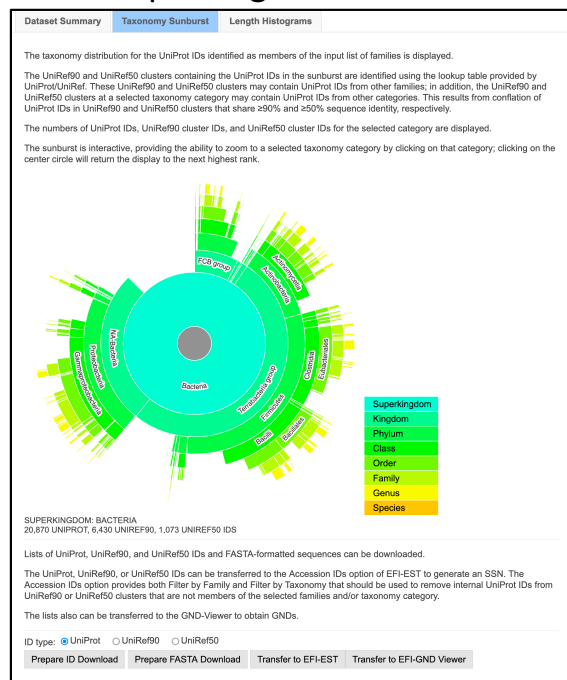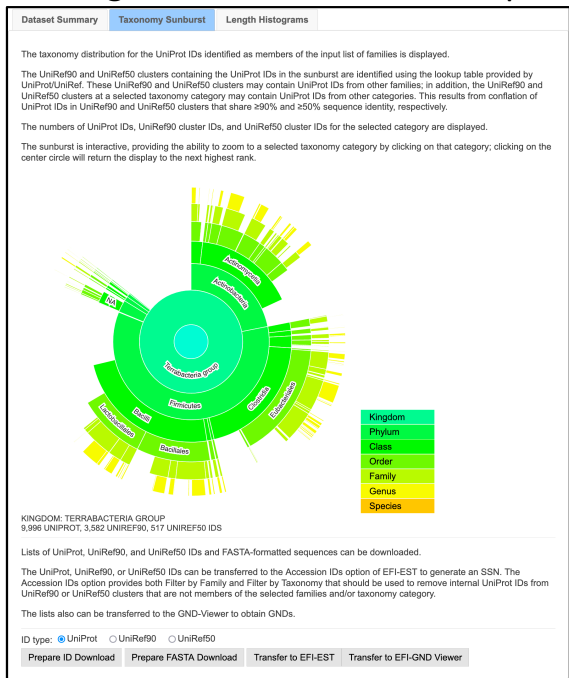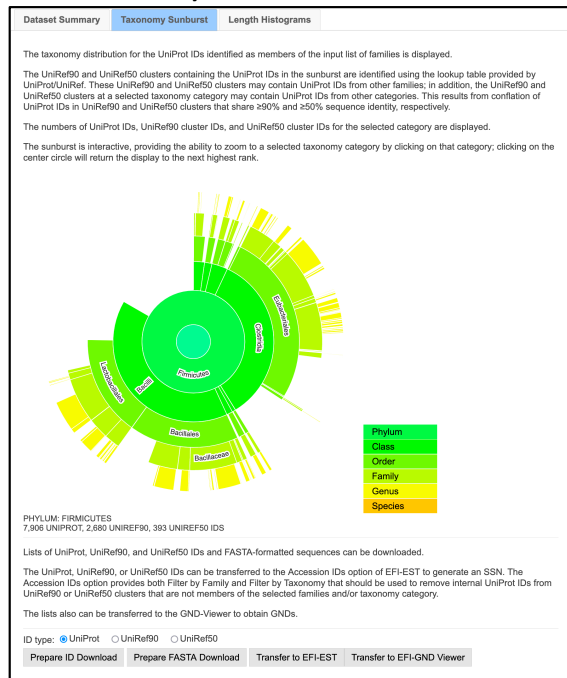The numbers of UniProt IDs, UniRef90 cluster IDs, and UniRef50 cluster IDs for the selected category are displayed.

The sunburst is interactive, providing the ability to zoom to a selected taxonomy category by clicking on that category; clicking on the center circle will return the display to the next highest rank.

Kingdom
Phylum
Class
Order
Family
Genus
Species

KINGDOM: TERRABACTERIA GROUP
9,996 UNIPROT, 3,582 UNIREF90, 517 UNIREF50 IDS

Lists of UniProt, UniRef90, and UniRef50 IDs and FASTA-formatted sequences can be downloaded.

The UniProt, UniRef90, or UniRef50 IDs can be transferred to the Accession IDs option of EFI-EST to generate an SSN. The Accession IDs option provides both Filter by Family and Filter by Taxonomy that should be used to remove internal UniProt IDs from UniRef90 or UniRef50 clusters that are not members of the selected families and/or taxonomy category.

The lists also can be transferred to the GND-Viewer to obtain GNDs.

ID type: ● UniProt ○ UniRef90 ○ UniRef50

Prepare ID Download | Prepare FASTA Download | Transfer to EFI-EST | Transfer to EFI-GND Viewer

# D — Phylum Firmicutes



Dataset Summary | Taxonomy Sunburst | Length Histograms

The taxonomy distribution for the UniProt IDs identified as members of the input list of families is displayed.

The UniRef90 and UniRef50 clusters containing the UniProt IDs in the sunburst are identified using the lookup table provided by UniProt/UniRef. These UniRef90 and UniRef50 clusters may contain UniProt IDs from other families; in addition, the UniRef90 and UniRef50 clusters at a selected taxonomy category may contain UniProt IDs from other categories. This results from conflation of UniProt IDs in UniRef90 and UniRef50 clusters that share ≥90% and ≥50% sequence identity, respectively.

The numbers of UniProt IDs, UniRef90 cluster IDs, and UniRef50 cluster IDs for the selected category are displayed.

The sunburst is interactive, providing the ability to zoom to a selected taxonomy category by clicking on that category; clicking on the center circle will return the display to the next highest rank.

Phylum
Class
Order
Family
Genus
Species

PHYLUM: FIRMICUTES
7,906 UNIPROT, 2,680 UNIREF90, 393 UNIREF50 IDS

Lists of UniProt, UniRef90, and UniRef50 IDs and FASTA-formatted sequences can be downloaded.

The UniProt, UniRef90, or UniRef50 IDs can be transferred to the Accession IDs option of EFI-EST to generate an SSN. The Accession IDs option provides both Filter by Family and Filter by Taxonomy that should be used to remove internal UniProt IDs from UniRef90 or UniRef50 clusters that are not members of the selected families and/or taxonomy category.

The lists also can be transferred to the GND-Viewer to obtain GNDs.

ID type: ● UniProt ○ UniRef90 ○ UniRef50

Prepare ID Download | Prepare FASTA Download | Transfer to EFI-EST | Transfer to EFI-GND Viewer
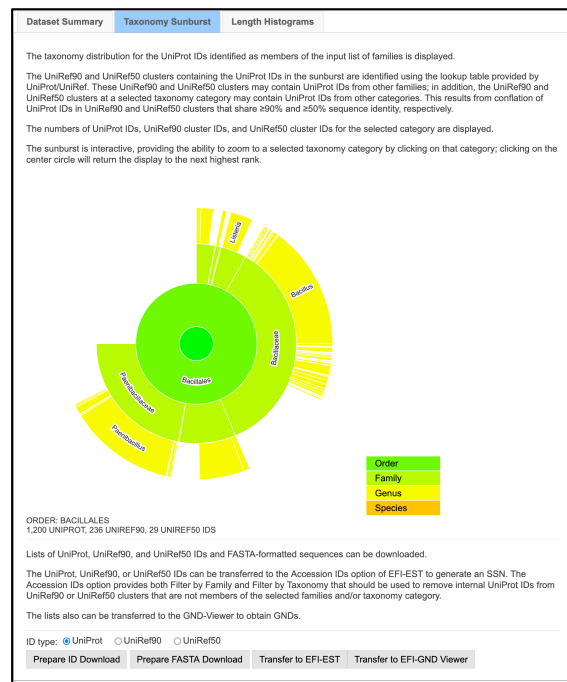
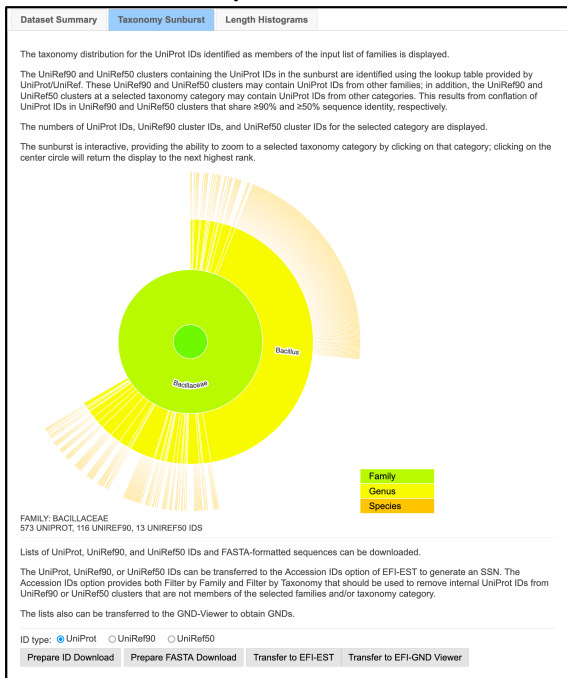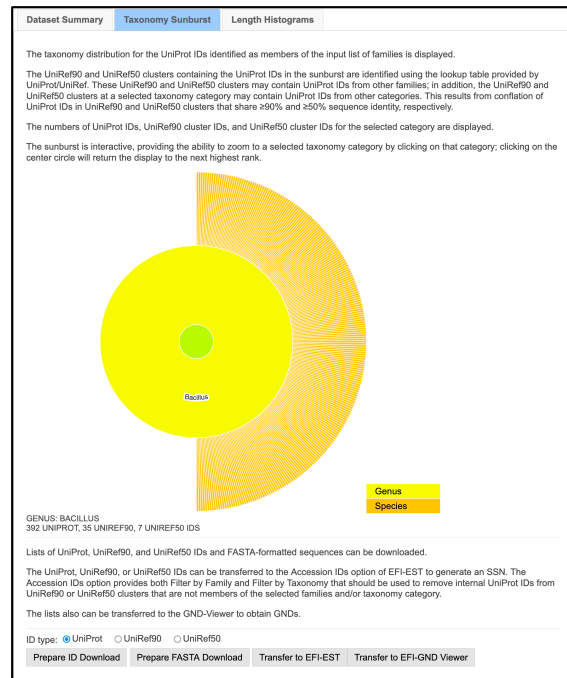**E**

# Class Bacilli



**F**

# Order Bacillales



**G**

# Family Bacillaceae



**H**

# Genuis Bacillus

**Supplementary Figure S4. Selection of Taxonomy Categories for GRE Superfamily in the Taxonomy Sunburst. Panel A**, All taxonomy categories (root). **Panel B**, Superkingdom Bacteria. **Panel C**, Kingdom Terrabacteria group. **Panel D**, Phylum Firmicutes. **Panel E**, Class Bacilli. **Panel F**, Order Bacillaes. **Panel G**, Family Bacillaceae. **Panel H**, Genus Bacillus.

| Sequence BLAST | Families | FASTA | Accession IDs | SSN Utilities |

**Generate a SSN from a list of UniProt, UniRef, NCBI, or Genbank IDs.**

An all-by-all BLAST ⑦ is performed to obtain the similarities between sequence pairs to calculate edge values to generate the SSN.

| Use UniProt IDs | Use UniRef50 or UniRef90 Cluster IDs |

Input a list of UniProt, NCBI, or Genbank (protein) accession IDs, or upload a text file.

**Accession IDs:**

**Accession ID File:** ⑦

⬆ Choose a file…

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

**Fragments:** ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**▾ Filter by Family**

**The input list of UniRef90 or UniRef50 cluster IDs should (must!) be filtered with the same list of Pfam families, InterPro families, and/or Pfam clans used to generate the IDs, if:**

The input list of UniRef90 or UniRef50 IDs is obtained from 1) the Color SSN or Cluster Analysis utility for a Families option (Option B) EFI-EST SSN, 2) the Families option of the Taxonomy Tool, or 3) the Accession IDs option of the Taxonomy Tool.

Input a list of Pfam families, InterPro families, and/or Pfam clans to restrict the UniProt and/or UniRef IDs in the SSN to these families.

**Family(s):**

The input format is a single family or comma/space separated list of families. Families should be specified as PFxxxxx (five digits), IPRxxxxxx (six digits) or CLxxxx (four digits) for Pfam clans.

For input lists of UniRef90 and UniRef50 clusters, the cluster ID (representative sequence) is used to identify those that match the list of families and are included in the SSN. The UniProt members in these clusters that do not match the input families are removed from the cluster and are not included in the SSN node attributes.

**▾ Filter by Taxonomy**

**The input list of UniRef90 or UniRef50 cluster IDs should (must!) be filtered with the same taxonomy categories used to generate the IDs, if:**

The input list of UniRef90 or UniRef50 IDs is obtained from 1) the Color SSN or Cluster Analysis utility for a Families option (Option B) EFI-EST SSN, 2) the Families option of the Taxonomy Tool, or 3) the Accession IDs option of the Taxonomy Tool.

From preselected conditions, the user can select "Bacteria, Archaea, Fungi", "Eukaryota, no Fungi", "Fungi", "Viruses", "Bacteria", "Eukaryota", or "Archaea" to restrict the UniProt IDs in the sunburst to these taxonomy groups.

"Bacteria, Archaea, Fungi", "Bacteria", "Archaea", and "Fungi" select organisms that may provide genome context (gene clusters/operons) useful for inferring functions.

The UniProt IDs also can be restricted to taxonomy categories within the Superkingdom, Kingdom, Phylum, Class, Order, Family, Genus, and Species ranks. Multiple conditions are combined to be a union of each other.

Preselected conditions: -- select a preset to auto populate --

Add Taxonomy category

**▸ Protein Family Addition Options**

**▸ Family Domain Boundary Options**

**▸ SSN Edge Calculation Option**

**Job name:** (required)

**E-mail address:** *Enter your e-mail address*

You will be notified by e-mail when your submission has been processed.

Submit Analysis

**Supplementary Figure S5.  EFI-EST, Accession IDs Option.**

**A**

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

Fragments:  ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**B**

**▾ Filter by Taxonomy**

From preselected conditions, the user can select "Bacteria, Archaea, Fungi", "Eukaryota, no Fungi", "Fungi", "Viruses", "Bacteria", "Eukaryota", or "Archaea" to restrict the UniProt IDs in the sunburst to these taxonomy groups.

"Bacteria, Archaea, Fungi", "Bacteria", "Archaea", and "Fungi" select organisms that may provide genome context (gene clusters/operons) useful for inferring functions.

The UniProt IDs also can be restricted to taxonomy categories within the Superkingdom, Kingdom, Phylum, Class, Order, Family, Genus, and Species ranks. Multiple conditions are combined to be a union of each other.

Preselected conditions: [ -- select a preset to auto populate -- ⌄ ]

[ Add Taxonomy category ]

**C**

**▾ Filter by Family**

Input a list of Pfam families, InterPro families, and/or Pfam clans to select sequences for inclusion in the Taxonomy Sunburst.

Family(s): [                                    ]

The input format is a single family or comma/space separated list of families. Families should be specified as PFxxxxx (five digits), IPRxxxxxx (six digits) or CLxxxx (four digits) for Pfam clans.

**Supplementary Figure S6. Filters. Panel A, Fragment Option**. **Panel B, Filter by Taxonomy**.

**Panel C, Filter by Family**.

**Families**   FASTA   Accession IDs

**Retrieve taxonomy for families.**

The UniProt IDs for family members are identified in UniProtKB with a list of Pfam families, InterPro families, and/or Pfam clans.

**Pfam and/or InterPro Families:**

The input format is a single family or comma/space separated list of families. Families should be specified as PFxxxxx (five digits), IPRxxxxxx (six digits) or CLxxxx (four digits) for Pfam clans.

Filter by Taxonomy can be used to remove UniProt IDs that do not match the specified taxonomy categories.

The remaining UniProt IDs are used to generate the sunburst.

UniRef90 and UniRef50 clusters that contain the UniProt IDs are retrieved from the UniRef90 andUniRef50 databases using the lookup table provided by UniProt/UniRef. Clusters for which the cluster ID (representative sequence) matches the list of families are retained.

The numbers of UniProt IDs and both UniRef90 cluster and UniRef50 cluster IDs are displayed on the sunburst; the UniProt IDs and both UniRef90 cluster and UniRef50 cluster IDs are available for download and/or transfer to the Accession ID option (Option D) of EFI-EST to generate SSNs.

**If the lists of UniRef90 or UniRef50 cluster IDs are used to generate SSNs with the Accession IDs option (Option D) of EFI-EST, the lists should (must!) be filtered with the same list of families (Filter by Family) and any specified taxonomy categories (Filter by Taxonomy) used to generate the lists.**

This filtering removes the UniRef90 and UniRef50 clusters with cluster IDs ("representative sequences") or internal UniProt IDs that are not members of the specified families or have the selected taxonomy categories.

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

**Fragments:**   ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**▸ Filter by Taxonomy**

**Job name:**  _____ (required)

**E-mail address:**  *Enter your e-mail address*

You will be notified by e-mail when your submission has been processed.

Submit Analysis

**Supplementary Figure S7.  Taxonomy Tool, Families Option**.

| Families | **FASTA** | Accession IDs |

**Retrieve taxonomy for FASTA files.**

The input is a list of FASTA-formatted sequences in which the headers contain a UniProt ID. The UniProt ID is required because it is used to retrieve the taxonomy from the UniProt database (FASTA header "reading").

The UniProt IDs for the family members are retrieved; these are used to calculate the sunburst.

**Sequences:**

**FASTA File:** ⑦

☁ Choose a file…

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

**Fragments:** ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**▸ Filter by Taxonomy**

**▸ Filter by Family**

**Job name:** _____ (required)

**E-mail address:** *Enter your e-mail address*

You will be notified by e-mail when your submission has been processed.

Submit Analysis

**Supplementary Figure S8.  Taxonomy Tool, FASTA Option**.

**Supplementary Figure S9. Taxonomy Tool, Accession IDs Option**.

**Supplementary Figure S10**. **SSN Utilities: Cluster Analysis Utility.**

**▾ Sequence Filter**

The MSA is generated with MUSCLE using the node IDs. Clusters containing less than the Minimum Node Count will be excluded from the analyses. Since MUSCLE can fail with a "large" number sequences (variable; anywhere from >750 to 1500), the Maximum Node Count parameter can be used to limit the number of sequences that MUSCLE uses.

**Minimum Node Count:** [          ] Minimum number of nodes in order to include a cluster in the computations [default: 5]

**Maximum Node Count:** [          ] Maximum number of nodes to include in the MSA [default: no maximum]

**▾ WebLogos**

A MSA for the (length-filtered) node IDs is generated using MUSCLE; the WebLogo is generated with the **http://weblogo.threeplusone.com** code.

**Make Weblogo:** ☑ Make Weblogos for each cluster [default: on]

**▾ Consensus Residues**

The positions and selected percent identities of the selected residues in the MSA are determined.

**Compute Consensus Residues:** ☑ Compute consensus residues [default: on]

[ C ] Residues to compute for (comma-separated list of amino acid codes)

[ 0.9,0.8,0.7,0.6,0.5,0.4,0 ] Percent identity threshold(s) for determining conservation (multiple comma-separated values allowed) [default: 0.9,0.8,0.7,0.6,0.5,0.4,0.3,0.2,0.1]

**▾ HMMs**

The MSA for the (length-filtered) node IDs is used to generate the HMM with hmmbuild from **HMMER3 (http://hmmer.org)**.

**Make HMMs:** ☑ Make HMMs for each cluster [default: on]

**▾ Length Histograms**

Length histograms for the node IDs (where applicable, UniProt, UniRef90, and UniRef50 IDs).

**Make Length Histograms:** ☑ Make length histograms for each cluster [default: on]

**Supplementary Figure S11**.  **SSN Utilities:  Cluster Analysis Utility Options.**

| Sequence BLAST | Families | FASTA | Accession IDs | **SSN Utilities** |

| Color SSNs | Cluster Analysis | Neighborhood Connectivity | **Convergence Ratio** |

**Convergence ratio is calculated per cluster.**

**SSN File:** ⑦

⬆ Choose a file…

A Color SSN (from either the Color SSN or Cluster Analysis utility) is the required input (cluster numbers are required).

**Alignment Score:** 5    The alignment score to calculate convergence ratio per cluster (should be the same as the original SSN alignment score).

The "convergence ratio" is the ratio of the actual number of edges in the cluster to the maximum possible number of edges (each node connected to every other node). For UniRef SSNs, two convergence ratios are calculated, one for the edges connecting the UniRef nodes in the input SSN and the second for the "hypothetical" edges that would connect the internal UniProt IDsin the cluster. The user specifies the value of the alignment score to be used (usually the same a alignment score used to generate the SSN).

The value of the convergence ratio ranges from 1.0 for sequences that are very similar ("identical") to 0.0 for sequences that are unrelated at the specified alignment score. The convergence ratio can be used as a criterion to infer whether an SSN cluster is isofunctional—the convergence ratio of a cluster containing orthologous sequences is expected to be close to 1.0 even at large alignment scores.

**E-mail address:**

You will be notified by e-mail when your submission has been processed.

Submit Analysis

**Supplementary Figure S12**. **SSN Utilities:  Convergence Ratio Utility.**

**Supplementary Figure S13**. SSN Utilities: Neighborhood Connectivity Utility.

**Supplementary Figure S14.  EFI-EST, Sequence BLAST Option.**

| Sequence BLAST | **Families** | FASTA | Accession IDs | SSN Utilities |

**Generate a SSN for a protein family.**

The members of the input Pfam families, InterPro families, and/or Pfam clans are selected from the UniProt, UniRef90, or UniRef50 database.

**Pfam and/or InterPro Families and/or Pfam clans:**

☐ Use [UniRef90 ⌄] cluster ID sequences instead of UniProt IDs (UniProt is default).

The input format is a single family or comma/space separated list of families. Families should be specified as PFxxxxx (five digits), IPRxxxxxx (six digits) or CLxxxx (four digits) for Pfam clans.

UniRef90 clusters contain UniProt IDs that share ≥90% sequence identity and have 80% overlap with the longest sequence in the cluster ("seed sequence"); as a result, the UniProt IDs in the cluster usually are functionally homogeneous, i.e., orthologues. UniRef50 clusters contain UniProt IDs that share ≥50% sequence identity and have 80% overlap with the seed sequence; as a result, the UniProt IDs in the cluster often are functionally heterogeneous, e.g., paralogues.

The sequences from the UniRef90 and UniRef90 databases are the UniRef90 and UniRef50 clusters for which the cluster ID ("representative sequence") matches the specified families. The UniProt members in these UniRef90 and Uni/Ref50 clusters that do not match the specified families are removed from the cluster.

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

**Fragments:** ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**▾ Filter by Taxonomy**

From preselected conditions, the user can select "Bacteria, Archaea, Fungi", "Eukaryota, no Fungi", "Fungi", "Viruses", "Bacteria", "Eukaryota", or "Archaea" to restrict the retrieved sequences to these taxonomy groups.

"Bacteria, Archaea, Fungi", "Bacteria", "Archaea", and "Fungi" select organisms that may provide genome context (gene clusters/operons) useful for inferring functions.

The retrieved sequences also can be restricted to taxonomy categories within the Superkingdom, Kingdom, Phylum, Class, Order, Family, Genus, and Species ranks. Multiple conditions are combined to be a union of each other.

The sequences from the UniRef90 and UniRef50 databases are the UniRef90 and UniRef50 clusters for which the cluster ID ("representative sequence") matches the specified taxonomy categories. The UniProt members in these clusters that do not match the specified taxonomy categories are removed from the cluster.

Preselected conditions: [ -- select a preset to auto populate -- ⌄ ]

[ Add Taxonomy category ]

**▸ Protein Family Size Options**

**▸ Family Domain Boundary Option**

**▸ SSN Edge Calculation Option**

**Job name:** [_____] (required)

**E-mail address:** [_____]

You will be notified by e-mail when your submission has been processed.

[ Submit Analysis ]

**Supplementary Figure S15**. **EFI-EST, Families Option.**

| Sequence BLAST | Families | FASTA | Accession IDs | SSN Utilities |

**Generate a SSN from FASTA-formatted UniProt sequences.**

An all-by-all BLAST ⑦ is performed to obtain the similarities between sequence pairs to calculate edge values to generate the SSN.

Input a list of sequences in the FASTA format or upload a FASTA-formatted sequence file.

Two options are available for generating the SSN:

1) The sequences are used "as is", with the node attributes including only the information in the header as the description and the number of residues in the sequence.

2) The ID in the header that immediately follows the ">" is used to retrieve node attribute information. Acceptable IDs include UniProt IDs, PDB IDs, and NCBI GenBank IDs that have equivalent entries in the UniProt database. ⑦ To use this option, check the "Read FASTA headers" box.

**Sequences:**

☐ **Read FASTA headers**

**FASTA File:** ⑦

⬆ Choose a file…

**▾ Fragment Option**

UniProt designates a Sequence Status for each member: Complete if the encoding DNA sequence has both start and stop codons; Fragment if the start and/or stop codon is missing. Approximately 10% of the entries in UniProt are fragments.

**Fragments:**   ☐ Check to exclude UniProt-defined fragments in the results. (default: off)

For the UniRef90 and UniRef50 databases, clusters are excluded if the cluster ID ("representative sequence") is a fragment.

UniProt IDs in UniRef90 and UniRef50 clusters with complete cluster IDs are removed from the clusters if they are fragments.

**▾ Filter by Family**

Input a list of Pfam families, InterPro families, and/or Pfam clans to restrict the UniProt and/or UniRef IDs in the SSN to these families.

**Family(s):**

The input format is a single family or comma/space separated list of families. Families should be specified as PFxxxxx (five digits), IPRxxxxxx (six digits) or CLxxxx (four digits) for Pfam clans.

**▾ Filter by Taxonomy**

From preselected conditions, the user can select "Bacteria, Archaea, Fungi", "Eukaryota, no Fungi", "Fungi", "Viruses", "Bacteria", "Eukaryota", or "Archaea" to restrict the input UniProt sequences to these taxonomy groups.

"Bacteria, Archaea, Fungi", "Bacteria", "Archaea", and "Fungi" select organisms that may provide genome context (gene clusters/operons) useful for inferring functions.

The input UniProt sequences also can be restricted to taxonomy categories within the Superkingdom, Kingdom, Phylum, Class, Order, Family, Genus, and Species ranks. Multiple conditions are combined to be a union of each other.

Preselected conditions: [ -- select a preset to auto populate -- ▾ ]

[ Add Taxonomy category ]

**▸ Protein Family Addition Options**

**▸ Family Domain Boundary Options**

**▸ SSN Edge Calculation Option**

**Job name:** [                    ] (required)

**E-mail address:** [                    ]

You will be notified by e-mail when your submission has been processed.

[ Submit Analysis ]

**Supplementary Figure S16**.  **EFI-EST, FASTA Option.**

**Supplementary Figure S17**. **SSN Utilities:   Color SSNs.**