



HHS Public Access

Author manuscript

ACS Bio Med Chem Au. Author manuscript; available in PMC 2023 February 16.

Published in final edited form as:

ACS Bio Med Chem Au. 2022 February 16; 2(1): 22–35. doi:10.1021/acsbioimedchemau.1c00048.

RadicalSAM.org: A Resource to Interpret Sequence-Function Space and Discover New Radical SAM Enzyme Chemistry

Nils Oberg¹, Timothy W. Precord^{1,2}, Douglas A. Mitchell^{1,2,3}, John A. Gerlt^{1,2,4}

¹Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States

²Department of Chemistry, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States

³Department of Microbiology, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States

⁴Department of Biochemistry, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana, Illinois 61801, United States

Abstract

The radical SAM superfamily (RSS), arguably the most functionally diverse enzyme superfamily, is also one of the largest with ~700K members currently in the UniProt database. The vast majority of the members have uncharacterized enzymatic activities and metabolic functions. In this Perspective, we describe [RadicalSAM.org](https://www.radicalsam.org), a new web-based resource that enables a user-friendly genomic enzymology strategy to explore sequence-function space in the RSS. The resource attempts to enable identification of isofunctional groups of radical SAM enzymes using sequence similarity networks (SSNs) and the genome context of the bacterial, archaeal, and fungal members provided by genome neighborhood diagrams (GNDs). Enzymatic activities and *in vivo* functions frequently can be inferred from genome context given the tendency for genes of related function to be clustered. We invite the scientific community to use [RadicalSAM.org](https://www.radicalsam.org) to (i) guide their experimental studies to discover new enzymatic activities and metabolic functions, (ii) contribute experimentally verified annotations to [RadicalSAM.org](https://www.radicalsam.org) to enhance the ability to predict novel activities and functions, and (iii) provide suggestions for improving this resource.

Keywords

Radical SAM superfamily; genomic enzymology; functional assignment; isofunctional families; sequence similarity networks; genome neighborhood diagrams; web resource

This issue celebrates the 20th anniversary of the discovery of the radical SAM superfamily (RSS)¹. The seminal bioinformatic study of 645 proteins by Sofia et al. in 2001 revealed a conserved CX₃CX₂C motif located near the N-terminus of a (β/α)₆-barrel domain that coordinates a [4Fe-4S] center that binds *S*-adenosylmethionine (SAM). As of mid-2021,

genome projects had identified ~700K additional members that include orthologs of characterized members as well as many uncharacterized members with potentially new enzyme activities and metabolic functions. However, the sheer size and accelerating growth of the RSS create a classic “big data” problem: exploration and interpretation of the sequence-function space has become untenable for non-bioinformaticians. In this Perspective, we describe [RadicalSAM.org](https://radicalsam.org/) (<https://radicalsam.org/>), an open-access “genomic enzymology” web resource designed for experimental biochemists that leverages the UniProt² (protein) and European Nucleotide Archive³ (ENA; nucleotide) databases (Figure 1). This Perspective provides a general textual and graphical overview of the resource; we also provide videos on the [RadicalSAM.org](https://radicalsam.org/tutorials.php) tutorial page (<https://radicalsam.org/tutorials.php>) that describe features of [RadicalSAM.org](https://radicalsam.org/).

RSS in the Structure-Function Linkage Database (SFLD)

The now archival Structure-Functional Linkage Database (SFLD) linked sequence-structure features to different chemical capabilities in several functionally diverse superfamilies⁸, including the RSS (<https://sfld.rbvi.ucsf.edu/archive/django/superfamily/29/index.html>). The SFLD segregated the sequence similarity network (SSN) of the RSS into 20 subgroups with functionally characterized members and 22 subgroups with uncharacterized members (designated by numbers and colors in Figure 2; see expanded image in the Supplementary Information); the 20 subgroups and their names are provided in Table 1. In the last update (2017)⁹, the SSN was generated using ~114K sequences then available in Pfam¹⁰ family PF04055 and InterPro¹¹ family IPR007197 collected at 50% sequence identity into 10,741 representative nodes (PF04055 and additional families defined by SFLD and PROSITE¹² are incorporated into IPR007197). The subgroups were identified by segregating the nodes in the SSN using a maximum e-value threshold of 1e-20 to draw edges.

As the protein databases continue to grow (doubling time ~2 years), the archival SFLD provides an increasingly outdated description of the RSS sequence-function space. Additionally, not all RSS subgroups are curated as Pfam and/or InterPro families, so current membership in these cases is unavailable. Furthermore, many of the hidden Markov models (HMMs) that were used to define the membership for the previously curated subgroups are no longer sufficient to reliably classify individual RSS proteins.

RadicalSAM.org: A Resource to Accelerate the Discovery of New Enzyme Chemistry

We anticipate that many readers want to assign *in vitro* enzymatic activities and *in vivo* metabolic functions to uncharacterized members of the RSS. We propose that this can be facilitated using a “genomic enzymology” strategy¹³. With this strategy, the members of a functionally diverse superfamily are segregated into potential isofunctional families using SSNs (separate SSN clusters). Then, the genomic contexts of the bacterial, archaeal, and fungal members of the clusters are retrieved and profiled. When the genomic neighborhood is reminiscent of a known pathway, the local context can be leveraged to rapidly formulate high quality hypotheses and aid in the design of experiments to confirm or refute the predicted enzymatic activity. When the genomic neighborhood is not similar to that for a

known pathway, the user can use contextual clues and other information to inform decisions on further experimental characterization.

The identification of isofunctional families cannot be accomplished using sequence identity alone—the sequence boundaries between homologs (proteins derived from a common ancestor) and orthologs (homologs separated by speciation that likely exhibit the same enzymatic activities and metabolic functions) are not easily determined. As a further complication, a single pairwise minimum sequence identity threshold (or minimum edge alignment score threshold in SSNs) likely will not separate orthologous groups (clusters) across any superfamily, especially those that are functionally diverse like the RSS.

Genome context is a powerful approach for inferring isofunctionality for bacterial, archaeal, and fungal enzymes: shared genome context (encoding the same metabolic pathway) can be used as evidence for shared function. However, genome context is not necessarily preserved across diverse taxa, so the ability to readily survey the genome contexts for *all* members of isofunctional families is essential for inferring metabolic functions and enzymatic activities.

To “democratize” the genome enzymology strategy, we developed and provide a web-based resource^{14–17} (<https://efi.igb.illinois.edu/>) with tools for (i) generating SSNs for protein families and separating these into clusters based on pairwise sequence identity thresholds (EFI-EST; <https://efi.igb.illinois.edu/efi-est/>); and (ii) collecting, visualizing, and analyzing genome context of proteins in the SSN clusters using genome neighborhood networks (GNNs) and genome neighborhood diagrams (GNDs; EFI-GNT; <https://efi.igb.illinois.edu/efi-gnt/>). The resource also provides a tool for prioritizing uncharacterized isofunctional SSN clusters for functional discovery based on metagenome abundance using chemically guided functional profiling (CGFP; EFI-CGFP; <https://efi.igb.illinois.edu/efi-cgfp/>).

Although the tools have accelerated the biochemical characterization of many proteins^{15, 18}, the size of the RSS prevents the experimental community from fully utilizing our resource, *i.e.*, the necessary computational resources for visualizing the SSNs far exceed what is common for personal computers. Therefore, we developed [RadicalSAM.org](https://radicalsam.org/) (<https://radicalsam.org/>), an open-access, web-based resource, for exploring sequence-function space for radical SAM enzymes that share the conserved CX₃CX₂C motif identified by Sofia et al. (Figure 3).

Following the approach used by the SFLD⁹, the SSN for the RSS is segregated into functionally curated as well as uncharacterized subgroups (*vide infra*). [RadicalSAM.org](https://radicalsam.org/) provides lists of the UniProt IDs for the sequences in the various subgroups as well as precalculated SSNs so that users can explore their regions of interest within the larger sequence-function space. [RadicalSAM.org](https://radicalsam.org/) also enables straightforward access to the genome context (GNDs) for members of the subgroups as well as for individual proteins in UniProt.

This remainder of this Perspective describes general features and salient details of [RadicalSAM.org](https://radicalsam.org/). Readers interested in using [RadicalSAM.org](https://radicalsam.org/) to guide their experimental work are encouraged to view the videos on the [RadicalSAM.org](https://radicalsam.org/tutorials.php) tutorial page (<https://radicalsam.org/tutorials.php>) that provide an overview of [RadicalSAM.org](https://radicalsam.org/) as well as

descriptions of some of its tools. The first video (Index of Tutorial Videos at the end of the manuscript) provides an overview of [RadicalSAM.org](https://www.RadicalSAM.org).

Sequences in [RadicalSAM.org](https://www.RadicalSAM.org)

[RadicalSAM.org](https://www.RadicalSAM.org) includes sequences for radical SAM enzymes that share the conserved CX₃CX₂C motif identified by Sofia et al.¹. Three known radical SAM families that do not share this motif are not included: (i) diphthamide synthase¹⁹ (PF01866), (ii) phosphomethylpyrimidine synthase²⁰ (ThiC; PF01964), and (iii) α-D-ribose 1-methylphosphonate C-P lyase²¹ (PhnJ; PF06007).

The UniProt 2020_05 (October 7, 2020) and InterPro 82 (October 8, 2020) databases were used to develop the release of [RadicalSAM.org](https://www.RadicalSAM.org) described in this Perspective. We used Option B of EFI-EST to specify one Pfam family and 172 InterPro families/domains (including PF04055 and IPR007197 used by the SFLD) for collecting sequences to include in the SSN, with the goal of providing a more comprehensive inventory of the membership than was provided by the SFLD. These families/domains are listed on the “Sequence Families” subtab of the “Current Release” tab on Home page of [RadicalSAM.org](https://www.RadicalSAM.org).

We identified 664,196 UniProt IDs representing 579,102 unique sequences in 66,428 UniRef50 clusters. [RadicalSAM.org](https://www.RadicalSAM.org) uses UniRef50 and UniRef90 sequence clusters (sequences that share 50% and 90% sequence identity, respectively; <https://www.uniprot.org/help/uniref>) to decrease the computational requirements for generating SSNs and enable visualization of the SSNs with Cytoscape²². A representative UniProt accession within each UniRef cluster is assigned by UniProt as its identifier. The sequence of the identifier is used to generate SSNs, multiple sequence alignments (MSAs), HMMs, and length histograms.

Not every retrieved sequence is “full length” (*vide infra*). We removed truncated sequences from [RadicalSAM.org](https://www.RadicalSAM.org) to 1) improve the quality and reliability of the multiple sequence alignments used to generate WebLogos²³ and HMMs and 2) reduce the number of isolated nodes (singletons) in SSNs generated with alignment score edge thresholds that collect sequences into putative isofunctional clusters.

We used two procedures to remove truncated sequences:

1) UniProt designates a “Sequence Status” for each accession: “Complete” if the encoding DNA includes both a start and stop codon; “Fragment” if either a start or stop codon is absent. After excluding fragments with the “Fragment Option” of Option B of EFI-EST, the sequence set contained 620,386 UniProt IDs represented by 535,892 unique sequences in 52,886 UniRef50 clusters.

2) A “Complete” sequence may be truncated because of sequencing errors. The shortest “Complete” sequence in PF04055 (UniProt ID A0A351TBI7) contains 39 residues and includes the CX₃CX₂C motif; another “Complete” sequence (UniProt ID A0A376TI31) contains 58 residues without the CX₃CX₂C motif. We inspected the UniProt ID length histograms for all clusters in a prototype version of [RadicalSAM.org](https://www.RadicalSAM.org) and identified

anaerobic ribonucleotide-triphosphate reductase activating enzymes as the “shortest” family with 140 residues. Therefore, we used 140 residues for the minimum length filter for generating the SSN used by [RadicalSAM.org](https://radicalsam.org). The final sequence set contained 616,009 UniProt IDs represented by 531,705 unique sequences (in 50,232 UniRef50 clusters).

In unpublished SSNs generated with more recent Pfam/InterPro releases, we used three additional Pfam families and 28 additional InterPro families/domains to collect sequences missing in the sequence set used for the current release of [RadicalSAM.org](https://radicalsam.org) (Supplementary Table S1). As an example, PoyD^{24, 25}, a radical SAM enzyme involved in polytheonamide biosynthesis that epimerizes L-amino acids within a peptide substrate to the D-configuration (UniProt ID J9ZW29), is missing in [RadicalSAM.org](https://radicalsam.org); it is included in IPR030950 that was added. If/when users recognize the absence of *bona fide* members of the RSS, they should contact us so that these can be included in future updates of [RadicalSAM.org](https://radicalsam.org). These omissions could occur if additional Pfam/InterPro families and domains are needed to identify sequences; alternatively, sequences will be missing if they had not been deposited in the UniProt database used to identify sequences.

The NCBI nr protein database is much larger than the UniProt database (~424M sequences on October 8, 2021 vs. ~220M sequences in UniProt Release 2021_03), so it includes more members of the RSS. As a result, radical SAM proteins that have been described/characterized in the literature may not be deposited in the UniProt database and, therefore, are not in [RadicalSAM.org](https://radicalsam.org). As an example, NxxcB, a radical SAM enzyme from *Streptococcus orisratti* that installs a β -thioether bond in a ribosomally synthesized and post-translationally modified peptide (RiPP)²⁶, is present in the NCBI database (accession identifier WP_018375754.1) but not UniProt or [RadicalSAM.org](https://radicalsam.org). Orthologs from other *Streptococcus* sp. can be located in [RadicalSAM.org](https://radicalsam.org) with the “Search by Sequence” feature (*vide infra*) and are present in Megacluster-1-1, SFLD subgroup 17, SPASM/twitch domain-containing.

When this manuscript was in preparation (October 2021), the UniProt 2021_03 and InterPro 86 databases were the most current. Using the current list of Pfam families and InterPro families/domains, we identified 754,719 UniProt IDs representing 664,851 unique sequences in 78,535 UniRef50 clusters. After excluding fragments, the sequence set contained 700,346 UniProt IDs represented by 611,187 unique sequences in 62,118 UniRef50 clusters. After applying the 140 residue minimum length filter, the sequence set contained 694,831 UniProt IDs represented by 605,897 unique sequences in 58,733 UniRef50 clusters.

RSS SSN

The SSN generated with the full length UniRef50 cluster identifiers can be visualized and edited with Cytoscape 3.8.2²² installed on a Mac Pro computer with 768 GB RAM. By visual inspection of SSNs with nodes colored by SFLD subgroup (with curated InterPro families/domains), a minimum edge alignment score threshold of 11 groups UniRef50 clusters into SFLD subgroups and, also, allows segregation of the SFLD subgroups (Figure 4). The resulting large cluster contained 615,705 UniProt IDs represented by 531,425 unique

sequences in 50,084 UniRef50 clusters. All further analysis and subgroup identification used these sequences.

Identification of Subgroups

In contrast to the SFLD's SSN, none of the functionally characterized SFLD subgroups are separated into distinct clusters in this SSN, the result of the larger number of UniProt IDs/UniRef50 clusters and the choice of a smaller minimum edge alignment score threshold to prevent separation of the SFLD subgroups into multiple clusters (the SFLD's SSN used a maximum edge e-value threshold of $1e-20$).

The SSN clusters containing the SFLD subgroups were segregated by manual deletion of edges using Cytoscape as described in the SUBGROUPS/SUBGROUP IDENTIFICATION tab on the [RadicalSAM.org](https://www.RadicalSAM.org) home page. The resulting SSN contained 10 clusters: five have been designated "megaclusters" because they contain multiple SFLD subgroups (and uncharacterized subgroups) and five have been designated as "clusters" that contain only a single SFLD subgroup (Clusters 6-10). The (mega)clusters are numbered in order of decreasing number of UniRef50 IDs/nodes; Megacluster-1 through Megacluster-5 then Cluster-6 through Cluster-10 (Figure 5). The megaclusters were also segregated by manual deletion of edges into SFLD-curated and uncharacterized subgroups as described in the SUBGROUPS tab on the [RadicalSAM.org](https://www.RadicalSAM.org) home page.

The segregation resulted in the 56 clusters/subgroups included in [RadicalSAM.org](https://www.RadicalSAM.org) (Figure 6). As described on their Explore pages (next section), some of the subclusters in Megaclusters-3 and -4 were further manually segregated (using increased minimum edge alignment score thresholds) to separate UniProtKB/SwissProt functions.

Explore Pages

[RadicalSAM.org](https://www.RadicalSAM.org) provides bioinformatic and genome context information (GNDs) for each of the clusters/subgroups on cluster-specific Explore pages. A representative Explore page, for Megacluster-3-1, 7-carboxy-7-deazaguanine synthase-like, SFLD subgroup 1, is shown in Figure 7 (see expanded image in the Supplementary Information).

This section highlights several types useful of information provided by Explore pages. The second video (Index of Tutorial Videos at the end of the manuscript) provides a description of the contents of an Explore page.

1. The convergence ratios (CRs) calculated using the UniRef 50 or UniRef90 cluster identifiers and the UniProt IDs in the SSN cluster using the minimum edge alignment score threshold (e-value) used to generate the cluster. The CR is a measure of sequence similarity and is the ratio of the number of edges at the specified alignment score relative to the total number of sequence pairs (maximum number of edges). The value decreases from 1.0 for sequences that are highly similar (not necessarily identical, depending on the value of the alignment score) to a value approaching 0 for very divergent sequences.

The CR is particularly useful for analyzing the “diced” clusters (next section) for which [RadicalSAM.org](#) provides a series of SSNs generated as a function of increasing minimum edge alignment score thresholds (increasing pairwise percent identity); as the clusters become isofunctional and the sequences become orthologous, the value of CR approaches 1.0. Exceptions occur when an isofunctional cluster contains orthologs from diverse taxa, *i.e.*, the value of CR in an isofunctional cluster may decrease as the minimum edge alignment score threshold increases as taxonomic divergence in sequence dominates the CR.

2. The numbers of conserved Cys residues in the multiple sequence alignment (MSA) calculated from 90 to 10% conservation in steps of 10%. This number is influenced by both sequence and length heterogeneity within the cluster (*vide infra*). By definition, all members of the RSS contain three conserved Cys residues in the Cx₃Cx₂C motif that participates in the [4Fe-4S] center than binds SAM; however, many subgroups contain additional [Fe-S] centers coordinated to additional conserved Cys residues. Examples are provided in a later section.

3. A list of community-provided annotations (ANNO button). The UniProtKB/SwissProt database is not a comprehensive list of experimentally verified functions; annotations provided by the community using the form provided on the SUBMIT tab at the top of each page inform inference of possible uncharacterized functions. We compiled the annotations currently provided; moving forward, we ask users to contribute annotations to add to the resource (*vide infra*).

4. The taxonomy sunburst (Figure 8) provides a graphical display of the taxonomic distribution of sequences in the cluster (TAXONOMY button; <https://github.com/vasturiano/sunburst-chart>). Clicking on a wedge expands the view to include the sequences represented in that wedge. Lists of accession IDs and FASTA files are available for download at any selected taxonomic level.

5. GNDs (GENOME NEIGHBORHOOD DIAGRAM button; Figure 9; see expanded image in the Supplementary Information) can be viewed for the UniRef50 cluster identifiers, UniRef90 cluster identifiers, and UniProt IDs for UniRef50 clusters and for the UniRef90 cluster identifiers and UniProt IDs for UniRef90 SSN clusters. The GNDs are used in the identification of isofunctional clusters as well as to provide genome context (metabolic pathway context) for discovering novel enzymatic activities and metabolic functions. The third video (Index of Tutorial Videos at the end of the manuscript) explains the use of the GND viewer.

6. The multiple sequence alignment (MSA) for the UniRef cluster identifiers in an SSN cluster is generated with MUSCLE^{27, 28} and can be opened with Jalview^{29, 30} (<https://www.jalview.org/>). The MSA is used to assess function/sequence heterogeneity.

7. The WebLogo²³ (<http://weblogo.threeplusone.com>) for the SSN cluster generated from the MSA. Given their importance in RSS structure and function, Cys residues are highlighted in red to allow their easy identification.

8. The HMM for the SSN cluster displayed using Skyalign³¹ (<https://skylign.org/>; Figure 10). This HMM viewer allows quick visualization of the consensus sequence for the cluster. HMMs provide the probability of a residue at any position in the sequence, not the percent conservation.

“Dicing” of Functionally Diverse SFLD Subgroups

SFLD subgroups 17, 5, 2, and 16, are particularly large and functionally diverse. In RadicalSAM.org, these are Megacluster-1-1 (SPASM/twitch domain-containing), Megacluster-2-1 (B12-binding domain-containing), Megacluster-2-2 (anaerobic coproporphyrinogen III oxidase-like), and Cluster-7 (PLP-dependent), respectively. As noted previously, the SSNs for functionally diverse superfamilies do not segregate into isofunctional families/SSN clusters using a single minimum edge alignment score threshold, making identification of orthologs a challenge. To aid in determining isofunctionality, [RadicalSAM.org](https://radicalsam.org) provides a series of SSNs for these subgroups with an increasing minimum edge alignment score threshold that we designate as SSN “dicing”. [RadicalSAM.org](https://radicalsam.org) also provides GNDs for each cluster in the diced SSNs to permit genome context to be used to infer isofunctionality and possible enzymatic activities and metabolic pathways.

The diced SSNs are generated with UniRef90 cluster identifiers instead of the UniRef50 cluster identifiers used to generate the SSN for the RSS so the sequences within each node are more likely to be isofunctional (sequences with 90% sequence identity). The sequences in the lower resolution UniRef50 clusters may be heterofunctional so the GNDs for the cluster identifiers as well as the UniProt IDs in the cluster may provide misleading information about metabolic context.

Using Megacluster-1-1 (SPASM³²/twitch³³ domain-containing) as an example, a series of 33 SSNs was generated as a function of minimum edge alignment score threshold, ranging from 25 in which the SSN is dominated by a large, complex cluster to 300 in which the SSN contains only a few small clusters with large CR values and likely isofunctional nodes (Figure 11). In the diced SSNs, only clusters with 3 nodes are shown (an Explore page is provided for each cluster).

As the minimum edge alignment score threshold increases, the SSN clusters that emerge are more similar in sequence and, therefore, are more likely to carry out similar reactions. However, since a single alignment score edge threshold for generating isofunctional clusters does not apply across the entire SSN, each diced SSN will be a mixture of heterofunctional and isofunctional clusters. The GND that can be accessed for each cluster in each SSN can be used to assess functional homogeneity in the cluster, realizing that genome context for orthologs often is not conserved across taxonomic divisions, as well as potentially provide metabolic pathway context for inferring functions of uncharacterized enzymes.

The Explore page for each cluster in each diced SSN provides the ability (Figure 12; see expanded image in the Supplementary Information) to (i) step forward/backward through the clusters in each SSN (left panel), (ii) select any cluster in the current SSN (center panel), or (iii) select any SSN in the diced series (right panel).

The alignment score (AS) “AS Walk-Through” button is located above the image for each cluster in the “diced” SSNs (Figure 13; see expanded image in the Supplementary Information). This button opens a pop-up window that identifies the progeny clusters at the next alignment score or the progenitor cluster at the previous alignment score; these can be accessed by clicking the name of the cluster, thereby allowing the identification of functionally diverse homologs, i.e., homologs that share a common mechanistic attributes such regiochemistry/stereochemistry of hydrogen abstraction. The walk-through window provides the UniProtKB/SwissProt function(s), if any, and user-supplied annotation(s), if any, in each progenitor and progeny cluster. The fourth video (Index of Tutorial Videos at the end of the manuscript) describes navigation through a series of diced SSNs.

Alternatively, using the Search tab (Figure 14; see expanded image in the Supplementary Information), the diced SSNs can be searched with either a UniProt ID (“Find by UniProt ID” in the Search tab in the menu bar at the top of the page) or a sequence (“Find by Sequence” in the Search tab in the menu bar at the top of the page). “Find by UniProt ID” identifies the cluster containing the specified UniProt ID in each diced SSN and provides its CR; clicking on the cluster number opens the Explore page for the cluster. “Find by Sequence” uses hmmscan³⁴ to scan the HMMs for the clusters in each SSN in the series and provides a list of matching clusters, their e-values, and CRs; clicking on the cluster number opens the Explore window for the cluster. The fifth video (Index of Tutorial Videos at the end of the manuscript) provides a description of the Search tab.

If the diced SSNs are searched with a UniProt ID, the AS Walk-Through window identifies the progeny cluster that contains the UniProt ID. If the diced SSNs are searched with a sequence, the walk-through window identifies the progeny cluster with the lowest E-value to the HMM for the cluster (best match to the cluster HMM). These features facilitate the selection of progeny clusters for more detailed investigation. The sixth video (Index of Tutorial Videos at the end of the manuscript) illustrates the use of dicing to aid in identification of isofunctional families/clusters using PqqE, the PqqA peptide cyclase in pyrroloquinoline quinone (PQQ) biosynthesis, as a case study.

Conserved Cys Residues

The chemistry of [Fe-S] centers is central to understanding reaction mechanisms in the RSS. RadicalSAM.org provides easy access to the number and positions of conserved Cys residues in isofunctional clusters, e.g., allowing identification of Cys motifs that form [Fe-S] centers in addition to the C_x₃C_x₂C motif that binds SAM as well other conserved Cys residues not involved in forming [Fe-S] centers but may be involved in the reaction mechanism.

Each Explore page provides a list of number of Conserved Cys Residues at 90, 80, 70, 60, 50, 40, 30, 20, and 10% sequence conservation (from the downloadable Consensus Residue table). The number almost always is a function of percent conservation, with the number increasing as percent conservation decreases. This occurs for two reasons: 1) length heterogeneity resulting from the presence of “truncated” sequences, although we have tried

to remove as many of these as possible; and 2) sequence heterogeneity resulting from functional heterogeneity (the cluster is not isofunctional).

Inspection of the WebLogo, MSA, and/or Skylign display of the HMM (“View HMM in Skylign”) for a cluster allows visual identification of conserved Cys motifs and their locations in the sequence. These motifs are most quickly identified by inspection of the Skylign HMM display (“View HMM in Skylign”) since the HMM display “removes” insertions/deletions.

By definition, a member of the RSS contains 3 conserved Cys residues in the conserved Cx_3Cx_2C motif that form the SAM-binding [4Fe-4S] center. However, variations of this motif can be identified using [RadicalSAM.org](https://radicalsam.org), e.g., Megacluster-2-4-1 with a $Cx_{11}Cx_2C$ motif, Megacluster-2-5 with a Cx_7Cx_2C motif, Megaclusters-4-6-1, -4-6-3, -4-6-5, and 4-6-6 with Cx_8Cx_2C motifs, Megacluster-4-6-4 with a Cx_9Cx_2C motif, Megacluster-4-10 with a Cx_5Cx_2C motif, Megacluster-4-11 with a Cx_4Cx_2C motif, and Megacluster-5-2 with a Cx_6Cx_2C motif.

Many subgroups/clusters contain additional conserved Cys motifs. Although widely distributed across the RSS, many are found in Megacluster-1-1, which contain SPASM/twitch domains. SPASM domains located C-terminal to the $(\beta/\alpha)_6$ -barrel domain contain 7 or 8 conserved Cys residues that bind two additional [4Fe-4S] centers³²; similarly positioned twitch domains contain 3 or 4 conserved Cys residues that bind one additional [4Fe-4S] center³³. Many isofunctional clusters that contain SPASM or twitch domains can be identified in the diced SSNs for Megacluster-1-1 generated with alignments scores ≥ 60 (as inferred from conserved genome context in the GNDs); interestingly, the conserved Cys residues do not share common sequence motifs, *i.e.*, they occur with different inter-Cys residue spacings.

The diced SSNs for Megacluster-1-1 also contains many clusters with other numbers of conserved Cys residues, ranging from 1 to >28 (e.g., Megacluster-1-1-11 at alignment score 300). Although most of the additional conserved Cys residues are located C-terminal to the $(\beta/\alpha)_6$ -barrel domain, some are located N-terminal to the $(\beta/\alpha)_6$ -barrel domain; in some members of Megacluster-1-1, additional conserved Cys residues are found both N- and C-terminal to the $(\beta/\alpha)_6$ -barrel domain. Although structures are not available for these proteins, perusal of models generated by trRosetta³⁵/RoseTTAfold³⁶/AlphaFold³⁷ suggest that these Cys residues usually are located in additional domains. Users of [RadicalSAM.org](https://yanglab.nankai.edu.cn/trRosetta/) can submit sequences to trRosetta (<https://yanglab.nankai.edu.cn/trRosetta/>) or RoseTTAfold (<https://rosetta.bakerlab.org/>) and then highlight the positions of Cys residues when viewing the predicted structures to determine whether the additional conserved Cys residues are proximal in space and, therefore, may participate in binding of additional [Fe-S] centers.

Additional conserved Cys motifs are also found in other subgroups/clusters. Notable examples include SFLD subgroup 14, methyltransferase Class D (Megacluster-1-3), with 6 conserved Cys residues N-terminal to the $(\beta/\alpha)_6$ -barrel domain; SFLD subgroup 12, methylthiotransferase (Megacluster-2-3), with conserved 3 Cys residues N-terminal to the $(\beta/\alpha)_6$ -barrel domain; SFLD subgroup 15, organic radical activating enzymes

(Megacluster-3-2-1-1), with 9 additional conserved Cys residues in the $(\beta/\alpha)_6$ -barrel domain; SFLD subgroup 15, organic radical activating enzymes (Megacluster-3-2-1-3) with 2 conserved Cys residues N-terminal to the $(\beta/\alpha)_6$ -barrel domain and 6 additional conserved Cys residues in the RSS $(\beta/\alpha)_6$ -barrel domain.

A Community Resource

We invite members of the community to enhance the capabilities of this resource by contributing their experimentally determined enzymatic activities and metabolic functions for previously uncharacterized members of the RSS. The SUBMIT tab on the Home page provides a form for providing this information (Figure 15; see expanded image in the Supplementary Information).

Informed inference of function in sequence-function space requires as many experimentally confirmed landmarks as possible. Although the SSNs generated by EFI-EST include a SwissProt Description node attribute (from UniProt), the SwissProt database is not a comprehensive (and accurate/reliable) list of characterized functions. Mining the literature for experimentally verified functions is tedious, and, unfortunately, many publications fail to include an accession identifier (UniProt or NCBI) for experimentally investigated proteins. Therefore, it is challenging (sometimes impossible) to associate a published experimentally established function with a specific protein. Nonetheless, we have provided a large number of literature citations for experimentally characterized members that are accessible using the ANNO(TATION) buttons on the Explore pages and AS Walk-Through pop-up windows in the diced SSNs. The “Submit” tab provides users with the ability to submit annotations and publication DOI’s for specific accession IDs. With each update of [RadicalSAM.org](https://www.radicalsam.org), we will make these functions and publications available.

We ask that members of the community include the UniProt and/or NCBI accession IDs for proteins that are functionally characterized in their publications. The IDs not only make it easier for UniProtKB/SwissProt to target new annotations for curation but also for members of the community to associate enzymatic activities and metabolic functions with specific proteins in the databases. BIOCHEMISTRY requires that authors include accession identifiers³⁸; we hope that additional journals will do the same in the future. In the meantime, please help future scientific data curation efforts by including accession identifiers in all of your publications!

Finally, we encourage suggestions for improving [RadicalSAM.org](https://www.radicalsam.org). Our guiding principle has been to provide information in a format that is friendly to experimentalists. Indeed, as [RadicalSAM.org](https://www.radicalsam.org) has been developed, we realized the need for new useful features that have been incorporated, *e.g.*, “dicing” and the AS Walk-Through pop-up windows. The CONTACT tab on the Home page can be used to submit suggestions (Figure 16).

Future/Planned Enhancements

[RadicalSAM.org](https://www.radicalsam.org) will be most useful when it is updated regularly, ideally with each update of the UniProt database (every eight weeks). We plan to implement a pipeline for updates in which the subgroups and their HMMs are defined annually (using the first annual release of

UniProt) by manual dissection of the SSN; then, for subsequent releases in the calendar year, the members of the subgroups will be retrieved using the HMMs.

The current release of [RadicalSAM.org](https://www.radicalsam.org) uses dicing to facilitate identification of isofunctional clusters/families in four of the largest, functionally diverse SFLD subgroups, i.e., the ability to provide genome context/GNDs for inference of function is the central concept of genomic enzymology. We are pleased with this feature and would like to extend it to all other subgroups. In principle, this can be automated, *i.e.*, generating SSNs as a function of alignment score followed by using the EFI-EST Cluster Analysis and Convergence Ratio utilities to obtain the information provided on the Explore pages. In practice, the automated generation of images of individual clusters for diced SSNs with thousands of clusters is computationally demanding but required for the user to assess divergence of sequence (and function) as the clusters segregate into smaller clusters as the minimum edge alignment score threshold increases.

We plan to provide predicted three-dimensional structures for the UniRef90 cluster identifiers in the various SSN clusters/subgroups generated by AlphaFold/DeepMind³⁷. We will implement this feature when UniProt/InterPro makes this comprehensive set of AlphaFold predictions available; currently predicted structures are available for the proteins encoded by 22 prototype organisms (<https://www.alphafold.ebi.ac.uk/>). These structures will enable visualization of domain organization/topology in multidomain RSS members as well as the locations of conserved Cys motifs that bind [Fe-S] centers. Perhaps some of these structures will be useful for prediction of substrates, products, and mechanisms using virtual ligand/intermediate docking^{39, 40}.

We would like to adapt our tools to use the NCBI and JGI-IMG databases, both of which are larger than UniProt; we have received requests from the community to do this. The sequences in the NCBI and JGI-IMG databases are not as well curated as those in UniProt, so any SSNs using these databases would contain fewer node attributes. Also, the NCBI database does not always assign membership in Pfam families and InterPro families/domains to its sequences so identification of the members of the RSS is associated with certain challenges. As time and resources permit, we will investigate whether future versions of [RadicalSAM.org](https://www.radicalsam.org) can be developed to be compatible with NCBI and JGI-IMG.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

The conceptualization and computational infrastructure that enabled the development of [RadicalSAM.org](https://www.radicalsam.org) was supported by two grants from NIH/NIGMS (U54GM093342 and P01GM118303). We thank Daniel Davidson and David Slater for their assistance with programming and cluster support. We also thank members of the Mitchell group for locating many UniProt IDs for previously characterized radical SAM enzymes by literature and database mining.

References

- [1]. Sofia HJ, Chen G, Hetzler BG, Reyes-Spindola JF, and Miller NE (2001) Radical SAM, a novel protein superfamily linking unresolved steps in familiar biosynthetic pathways with radical mechanisms: functional characterization using new analysis and information visualization methods, *Nucleic Acids Res* 29, 1097–1106. [PubMed: 11222759]
- [2]. UniProt C (2021) UniProt: the universal protein knowledgebase in 2021, *Nucleic Acids Res* 49, D480–D489. [PubMed: 33237286]
- [3]. Harrison PW, Ahamed A, Aslam R, Alako BTF, Burgin J, Buso N, Courtot M, Fan J, Gupta D, Haseeb M, Holt S, Ibrahim T, Ivanov E, Jayathilaka S, Balavenkataraman Kadhirvelu V, Kumar M, Lopez R, Kay S, Leinonen R, Liu X, O’Cathail C, Pakseresht A, Park Y, Pesant S, Rahman N, Rajan J, Sokolov A, Vijayaraja S, Waheed Z, Zyoud A, Burdett T, and Cochrane G (2021) The European Nucleotide Archive in 2020, *Nucleic Acids Res* 49, D82–D85. [PubMed: 33175160]
- [4]. Frey PA, Hegeman AD, and Ruzicka FJ (2008) The Radical SAM Superfamily, *Crit Rev Biochem Mol Biol* 43, 63–88. [PubMed: 18307109]
- [5]. Vey JL, and Drennan CL (2011) Structural insights into radical generation by the radical SAM superfamily, *Chem Rev* 111, 2487–2506. [PubMed: 21370834]
- [6]. Booker SJ (2012) Radical SAM enzymes and radical enzymology, *Biochim Biophys Acta* 1824, 1151–1153. [PubMed: 22850428]
- [7]. Broderick WE, Hoffman BM, and Broderick JB (2018) Mechanism of Radical Initiation in the Radical S-Adenosyl-l-methionine Superfamily, *Acc Chem Res* 51, 2611–2619. [PubMed: 30346729]
- [8]. Akiva E, Brown S, Almonacid DE, Barber AE 2nd, Custer AF., Hicks MA., Huang CC., Lauck F, Mashiyama ST., Meng EC., Mischel D., Morris JH., Ojha S., Schnoes AM., Stryke D., Yunes JM., Ferrin TE., Holliday GL., and Babbitt PC. (2014) The Structure-Function Linkage Database, *Nucleic Acids Res* 42, D521–530. [PubMed: 24271399]
- [9]. Holliday GL, Akiva E, Meng EC, Brown SD, Calhoun S, Pieper U, Sali A, Booker SJ, and Babbitt PC (2018) Atlas of the Radical SAM Superfamily: Divergent Evolution of Function Using a “Plug and Play” Domain, *Methods Enzymol* 606, 1–71. [PubMed: 30097089]
- [10]. Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, Tosatto SCE, Paladin L, Raj S, Richardson LJ, Finn RD, and Bateman A (2021) Pfam: The protein families database in 2021, *Nucleic Acids Res* 49, D412–D419. [PubMed: 33125078]
- [11]. Blum M, Chang HY, Chuguransky S, Grego T, Kandasaamy S, Mitchell A, Nuka G, Paysan-Lafosse T, Qureshi M, Raj S, Richardson L, Salazar GA, Williams L, Bork P, Bridge A, Gough J, Haft DH, Letunic I, Marchler-Bauer A, Mi H, Natale DA, Necci M, Orengo CA, Pandurangan AP, Rivoire C, Sigrist CJA, Sillitoe I, Thanki N, Thomas PD, Tosatto SCE, Wu CH, Bateman A, and Finn RD (2021) The InterPro protein families and domains database: 20 years on, *Nucleic Acids Res* 49, D344–D354. [PubMed: 33156333]
- [12]. Sigrist CJ, de Castro E, Cerutti L, Cuče BA, Hulo N, Bridge A, Bougueleret L, and Xenarios I (2013) New and continuing developments at PROSITE, *Nucleic Acids Res* 41, D344–347. [PubMed: 23161676]
- [13]. Gerlt JA, and Babbitt PC (2001) Divergent evolution of enzymatic function: mechanistically diverse superfamilies and functionally distinct suprafamilies, *Annu Rev Biochem* 70, 209–246. [PubMed: 11395407]
- [14]. Gerlt JA, Bouvier JT, Davidson DB, Imker HJ, Sadkhin B, Slater DR, and Whalen KL (2015) Enzyme Function Initiative-Enzyme Similarity Tool (EFI-EST): A web tool for generating protein sequence similarity networks, *Biochim Biophys Acta* 1854, 1019–1037. [PubMed: 25900361]
- [15]. Gerlt JA (2017) Genomic Enzymology: Web Tools for Leveraging Protein Family Sequence-Function Space and Genome Context to Discover Novel Functions, *Biochemistry* 56, 4293–4308. [PubMed: 28826221]
- [16]. Zallot R, Oberg NO, and Gerlt JA (2018) ‘Democratized’ genomic enzymology web tools for functional assignment, *Curr Opin Chem Biol* 47, 77–85. [PubMed: 30268904]

- [17]. Zallot R, Oberg N, and Gerlt JA (2019) The EFI Web Resource for Genomic Enzymology Tools: Leveraging Protein, Genome, and Metagenome Databases to Discover Novel Enzymes and Metabolic Pathways, *Biochemistry* 58, 4169–4182. [PubMed: 31553576]
- [18]. Zallot R, Oberg N, and Gerlt JA (2021) Discovery of new enzymatic functions and metabolic pathways using genomic enzymology web tools, *Curr Opin Biotechnol* 69, 77–90. [PubMed: 33418450]
- [19]. Zhang Y, Zhu X, Torelli AT, Lee M, Dzikovski B, Koralewski RM, Wang E, Freed J, Krebs C, Ealick SE, and Lin H (2010) Diphthamide biosynthesis requires an organic radical generated by an iron-sulphur enzyme, *Nature* 465, 891–896. [PubMed: 20559380]
- [20]. Chatterjee A, Li Y, Zhang Y, Grove TL, Lee M, Krebs C, Booker SJ, Begley TP, and Ealick SE (2008) Reconstitution of ThiC in thiamine pyrimidine biosynthesis expands the radical SAM superfamily, *Nat Chem Biol* 4, 758–765. [PubMed: 18953358]
- [21]. Kamat SS, Williams HJ, Dangott LJ, Chakrabarti M, and Raushel FM (2013) The catalytic mechanism for aerobic formation of methane by bacteria, *Nature* 497, 132–136. [PubMed: 23615610]
- [22]. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, and Ideker T (2003) Cytoscape: a software environment for integrated models of biomolecular interaction networks, *Genome Res* 13, 2498–2504. [PubMed: 14597658]
- [23]. Crooks GE, Hon G, Chandonia JM, and Brenner SE (2004) WebLogo: a sequence logo generator, *Genome Res* 14, 1188–1190. [PubMed: 15173120]
- [24]. Parent A, Benjdia A, Guillot A, Kubiak X, Balty C, Lefranc B, Leprince J, and Berteau O (2018) Mechanistic Investigations of PoyD, a Radical S-Adenosyl-l-methionine Enzyme Catalyzing Iterative and Directional Epimerizations in Polytheonamide A Biosynthesis, *J Am Chem Soc* 140, 2469–2477. [PubMed: 29253341]
- [25]. Morinaka BI, Vagstad AL, Helf MJ, Gugger M, Kegler C, Freeman MF, Bode HB, and Piel J (2014) Radical S-adenosyl methionine epimerases: regioselective introduction of diverse D-amino acid patterns into peptide natural products, *Angew Chem Int Ed Engl* 53, 8503–8507. [PubMed: 24943072]
- [26]. Caruso A, Bushin LB, Clark KA, Martinie RJ, and Seyedsayamdost MR (2019) Radical Approach to Enzymatic beta-Thioether Bond Formation, *J Am Chem Soc* 141, 990–997. [PubMed: 30521328]
- [27]. Edgar RC (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput, *Nucleic Acids Res* 32, 1792–1797. [PubMed: 15034147]
- [28]. Edgar RC (2004) MUSCLE: a multiple sequence alignment method with reduced time and space complexity, *BMC Bioinformatics* 5, 113. [PubMed: 15318951]
- [29]. Waterhouse AM, Procter JB, Martin DM, Clamp M, and Barton GJ (2009) Jalview Version 2--a multiple sequence alignment editor and analysis workbench, *Bioinformatics* 25, 1189–1191. [PubMed: 19151095]
- [30]. Procter JB, Carstairs GM, Soares B, Mourao K, Ofoegbu TC, Barton D, Lui L, Menard A, Sherstnev N, Roldan-Martinez D, Duce S, Martin DMA, and Barton GJ (2021) Alignment of Biological Sequences with Jalview, *Methods Mol Biol* 2231, 203–224. [PubMed: 33289895]
- [31]. Wheeler TJ, Clements J, and Finn RD (2014) Skylign: a tool for creating informative, interactive logos representing sequence alignments and profile hidden Markov models, *BMC Bioinformatics* 15, 7. [PubMed: 24410852]
- [32]. Haft DH, and Basu MK (2011) Biological systems discovery in silico: radical S-adenosylmethionine protein families and their target peptides for posttranslational modification, *J Bacteriol* 193, 2745–2755. [PubMed: 21478363]
- [33]. Grell TA, Goldman PJ, and Drennan CL (2015) SPASM and twitch domains in S-adenosylmethionine (SAM) radical enzymes, *J Biol Chem* 290, 3964–3971. [PubMed: 25477505]
- [34]. Eddy SR (2011) Accelerated Profile HMM Searches, *PLoS Comput Biol* 7, e1002195. [PubMed: 22039361]

- [35]. Yang J, Anishchenko I, Park H, Peng Z, Ovchinnikov S, and Baker D (2020) Improved protein structure prediction using predicted interresidue orientations, *Proc Natl Acad Sci U S A* 117, 1496–1503. [PubMed: 31896580]
- [36]. Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, Lee GR, Wang J, Cong Q, Kinch LN, Schaeffer RD, Millan C, Park H, Adams C, Glassman CR, DeGiovanni A, Pereira JH, Rodrigues AV, van Dijk AA, Ebrecht AC, Opperman DJ, Sagmeister T, Buhlheller C, Pavkov-Keller T, Rathinaswamy MK, Dalwadi U, Yip CK, Burke JE, Garcia KC, Grishin NV, Adams PD, Read RJ, and Baker D (2021) Accurate prediction of protein structures and interactions using a three-track neural network, *Science* 373, 871–876. [PubMed: 34282049]
- [37]. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates R, Zidek A, Potapenko A, Bridgland A, Meyer C, Kohl SAA, Ballard AJ, Cowie A, Romera-Paredes B, Nikolov S, Jain R, Adler J, Back T, Petersen S, Reiman D, Clancy E, Zielinski M, Steinegger M, Pacholska M, Berghammer T, Bodenstein S, Silver D, Vinyals O, Senior AW, Kavukcuoglu K, Kohli P, and Hassabis D (2021) Highly accurate protein structure prediction with AlphaFold, *Nature* 596, 583–589. [PubMed: 34265844]
- [38]. Gerlt JA (2018) The Need for Manuscripts to Include Database Identifiers for Proteins, *Biochemistry* 57, 4239–4230. [PubMed: 30037234]
- [39]. Hermann JC, Marti-Arbona R, Fedorov AA, Fedorov E, Almo SC, Shoichet BK, and Raushel FM (2007) Structure-based activity prediction for an enzyme of unknown function, *Nature* 448, 775–779. [PubMed: 17603473]
- [40]. Song L, Kalyanaraman C, Fedorov AA, Fedorov EV, Glasner ME, Brown S, Imker HJ, Babbitt PC, Almo SC, Jacobson MP, and Gerlt JA (2007) Prediction and assignment of function for a divergent N-succinyl amino acid racemase, *Nat Chem Biol* 3, 486–491. [PubMed: 17603539]

Index of Tutorial Videos

1. [RadicalSAM.org](#) empowers researchers in the field of radical SAM enzymology to discover novel functions of uncharacterized members of the RSS by democratizing the genomic enzymology tools. This video provides a tour of the site and encourages the user to explore the tools using their own radical SAM enzymes.
2. The Explore tab gives the user the ability to maneuver through the enormous protein dataset in [RadicalSAM.org](#). This video provides a tour of the Explore tab and details the various tools in the tab.
3. GNDs are powerful tools in determining the potential function of a radical SAM enzyme. This video describes their use and how they have been incorporated into [RadicalSAM.org](#)
4. [RadicalSAM.org](#) uses SSNs generated at a series of increasing alignment scores, through a process known as dicing, to assist the user in determining points of isofunctionality in the radical SAM superfamily. This video provides a brief description of that process.
5. The Search tab enables the user to search for a specific radical SAM in [RadicalSAM.org](#). This video provides a tour of the Search tab and describes how a user can find useful information about their radical SAM enzyme.
6. This video provides a tutorial on determining an isofunctional alignment score for a radical SAM of interest using PqqE, the PqqA peptide cyclase in pyrroloquinoline quinone (PQQ) biosynthesis, as a case study.

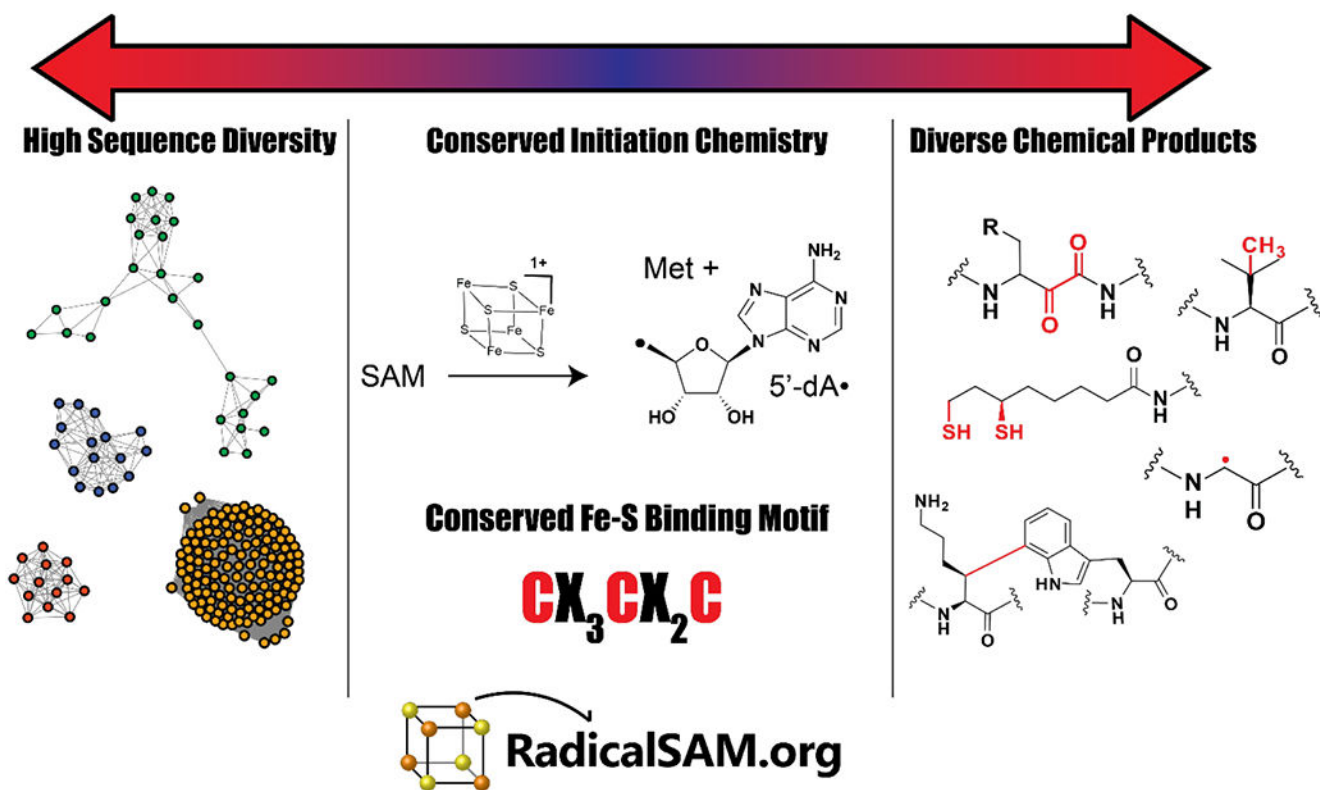


Figure 1.

Despite high diversity in sequence and reaction outcome, all members of the RSS generate 5'-deoxyadenosyl (5'-dA) radical using a conserved [4Fe-4S]-forming motif that binds and reductively liberates Met from SAM⁴⁻⁷. RadicalSAM.org provides easy access to a genomic enzymology strategy to catalog known enzymatic activities and discover new ones within the RSS.

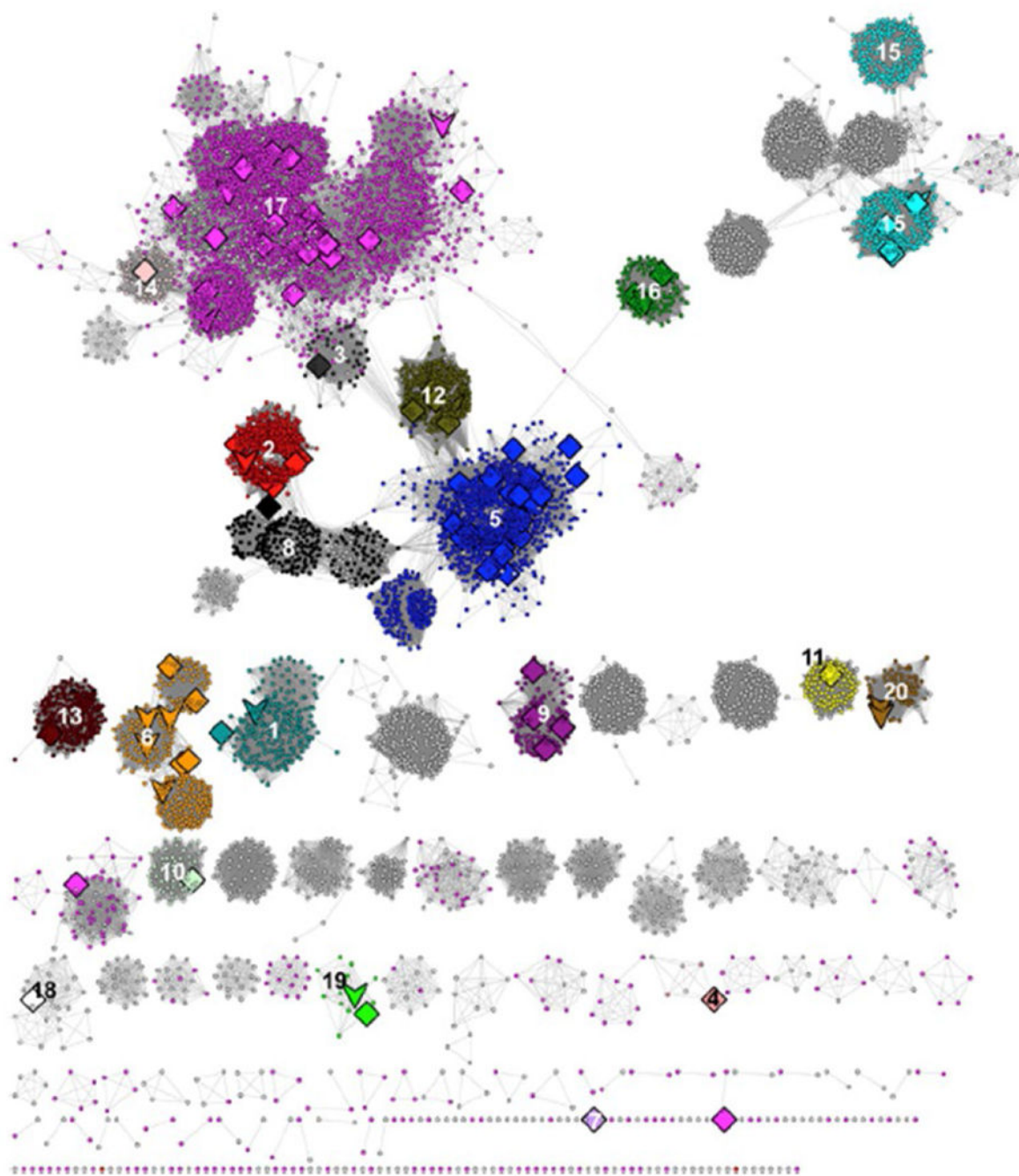


Figure 2. The SSN generated with a maximum e-value edge threshold of $1e-20$ used by the SFLD to identify its 20 functionally characterized subgroups (colored/numbered clusters) and 22 uncharacterized subgroups. Large nodes represent experimentally characterized proteins; downward arrows indicate a structurally characterized protein; diamonds indicate no structural characterization. Reproduced with permission from reference 9; <http://www.elsevier.com>

Home

The radical SAM superfamily (RSS) is arguably the largest and most functionally diverse enzyme superfamily. Many functions (and intriguing reaction mechanisms) have been discovered; many more remain to be discovered!

RadicalSAM.org is designed to leverage "top-down" discovery of function using the EFI's genomic enzymology web tools. The sequence similarity network (SSN) for the RSS is too large to be analyzed with Cytoscape and the RAM available on most computers so has been inaccessible to RSS community.

We generated the SSN for the entire RSS using a computer with 768GB RAM and segregated it into clusters for 1) the 20 subgroups curated by the Structure-Function Linkage Database (SFLD) and 2) many additional subgroups not curated by the SFLD.

For each subgroup, RadicalSAM.org provides:

1. The SSN, multiple sequence alignment (MSA), WebLogo, hidden Markov model (HMM), length histogram, phylogenetic distribution, SwissProt annotations, and number and locations of conserved Cys residues.
2. Genome neighborhood diagrams (GNDs) that provide metabolic pathway context for inference of functions.
3. UniProt accession IDs and FASTA sequences that can be used with EFI-EST, EFI-GNT, and EFI-CGFP for user-specific applications.
4. For four large and functionally diverse subgroups, the ability to "walk" through a series of SSNs generated at increasing alignment scores. The progeny (walking forward) and progenitors of a cluster (walking backward) can be identified, allowing the discovery of related functions and/or substrate specificities.

We encourage users to submit experimentally characterized functional annotations for sequences that have not yet been curated by SwissProt so that these can be made available to the RSS community.

A Perspective was recently published in ACS Bio & Med Chem Au describing RadicalSAM.org. If you use RadicalSAM.org, please cite us:

Nils Oberg, Timothy W. Precord, Douglas A. Mitchell, and John A. Gerit, **RadicalSAM.org: A Resource to Interpret Sequence-Function Space and Discover New Radical SAM Enzyme Chemistry**, ACS Bio & Med Chem Au 2021 ??????. <https://doi.org/10.1021/acsbiochemau.??????>

BACKGROUND

CURRENT RELEASE

SUBGROUPS

FUNCTIONALLY DIVERSE SUBGROUPS

EXPLORE PAGES

SEARCH FUNCTIONS

COMMUNITY ANNOTATION

Figure 3.
Home page of [RadicalSAM.org](https://radicalsam.org).

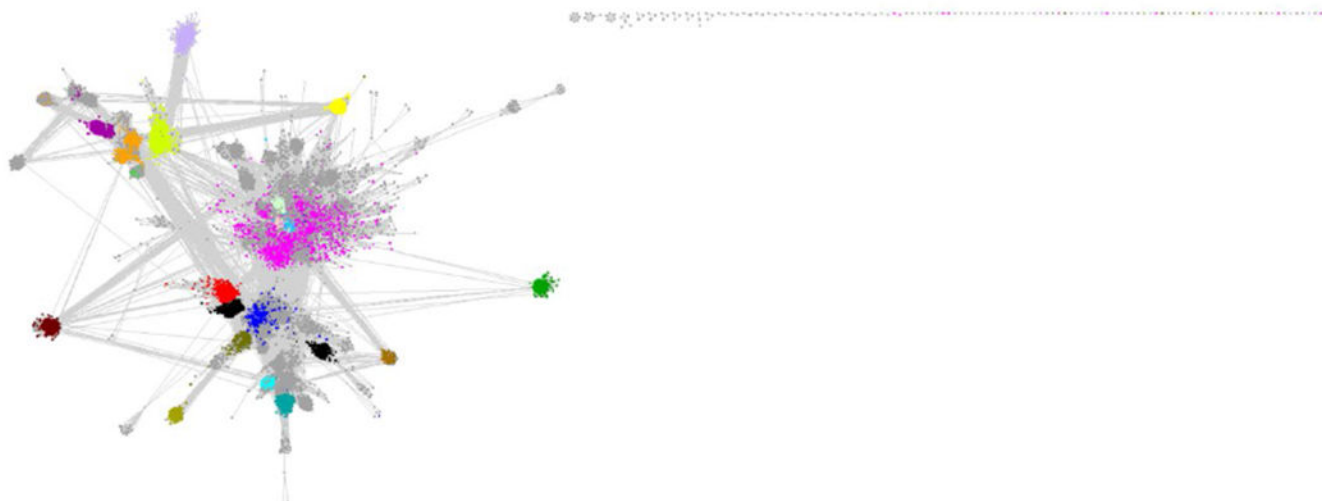


Figure 4. SSN of the starting point for RSS subgroup identification. The network is visualized with 11 as the minimum edge alignment score threshold and colored based on previously assigned SFLD subgroups (Figure 2 and Table 1).

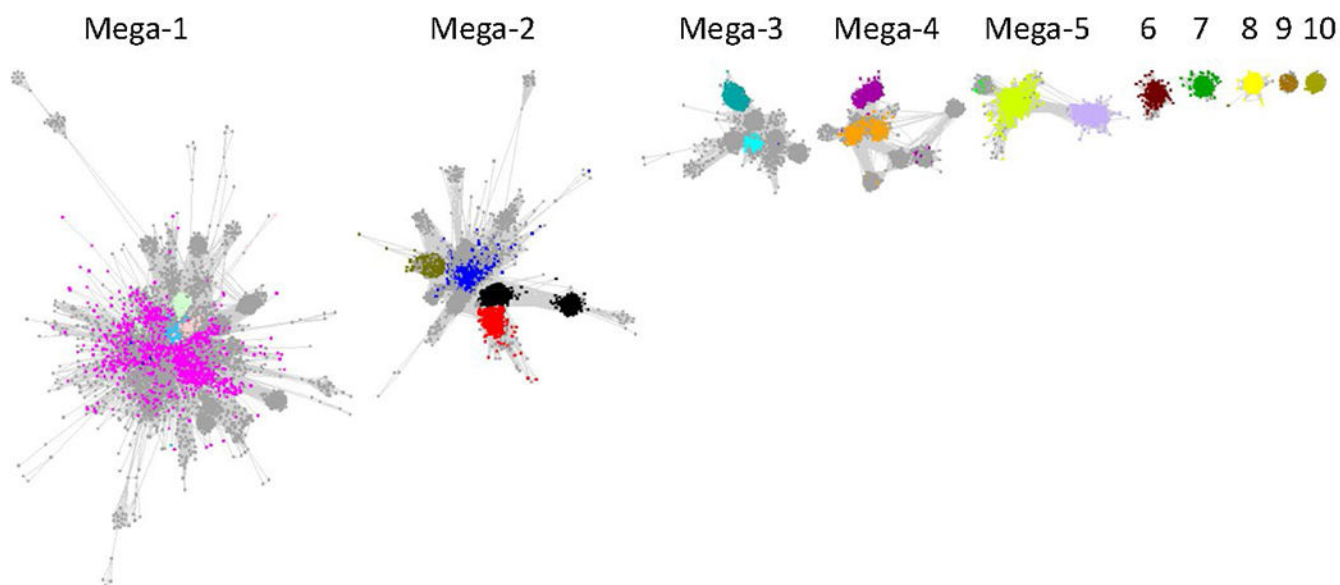


Figure 5.

The five megaclusters (containing multiple SFLD subgroups) and five standard clusters (containing single SFLD subgroups) after manual deletion of “long” edges that correspond to larger alignment scores and connect functionally divergent nodes/subgroups in the large cluster in Figure 4. The nodes are colored based on previously assigned SFLD subgroups (Figure 2 and Table 1).

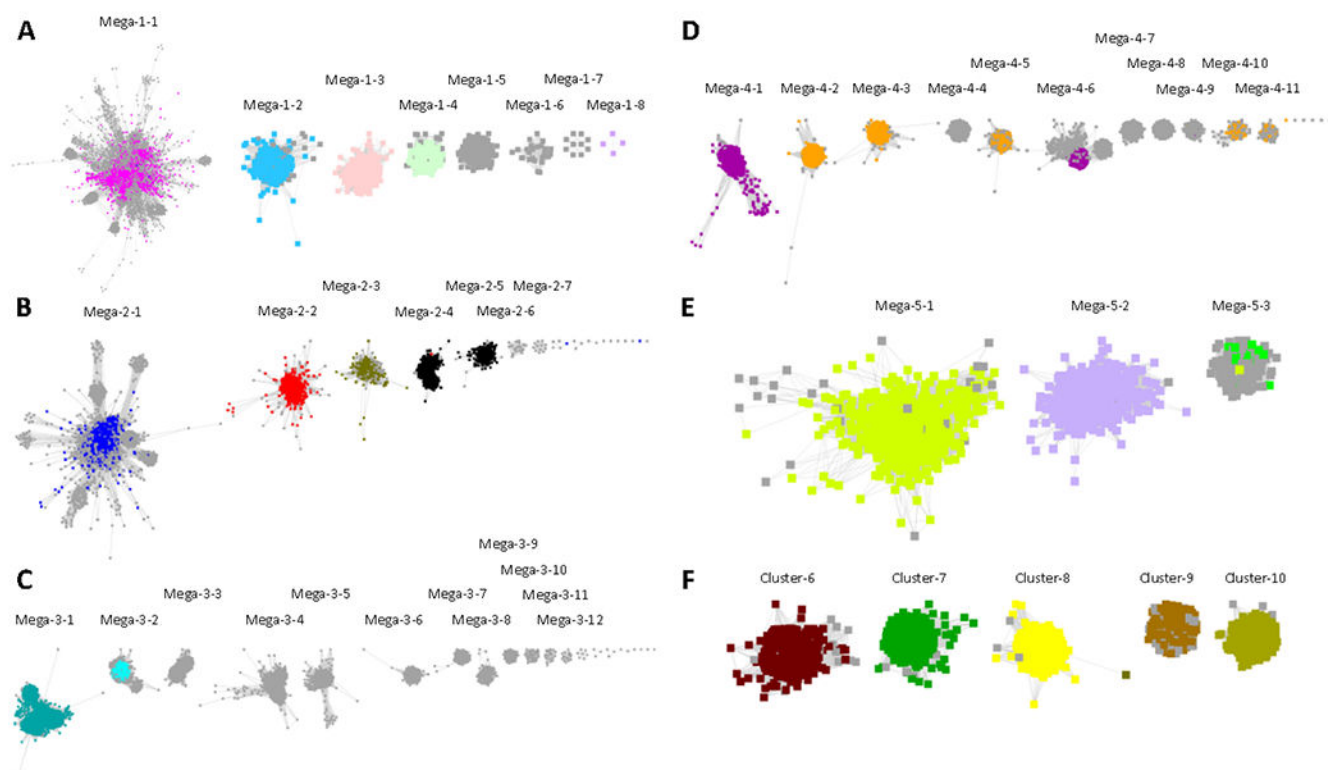


Figure 6.

The segregated subgroups in [RadicalSAM.org](https://radicalsam.org). The nodes are colored based on previously assigned SFLD subgroups (Figure 2 and Table 1). **Panel A**, UniRef50 SSNs for Megaclusters-1-2 through -1-8 were removed from Megacluster-1 using an alignment score edge threshold of 30 and displayed with 16 as the minimum edge alignment score threshold. The remaining Megacluster-1-1 is displayed using 11 as the minimum edge alignment score threshold. **Panel B**, UniRef50 SSNs for Megaclusters-2-1 through -2-7 were segregated using 12 as the minimum edge alignment score threshold. **Panel C**, UniRef90 SSNs for Megaclusters-3-1 through -3-12 segregated using 18 as the minimum edge alignment score threshold. **Panel D**, UniRef90 SSNs for Megaclusters-4-1 through -4-11 segregated using 22 as the minimum edge alignment score threshold. **Panel E**, UniRef50 SSNs for Megaclusters-5-1 through -5-3 were segregated using 12 as the minimum edge alignment score threshold. **Panel F**, UniRef50 SSNs for Clusters 6-10 from the segregated SSN in Figure 5.

Megacluster-3-1: 7-carboxy-7-deazaguanine synthase-like [1]

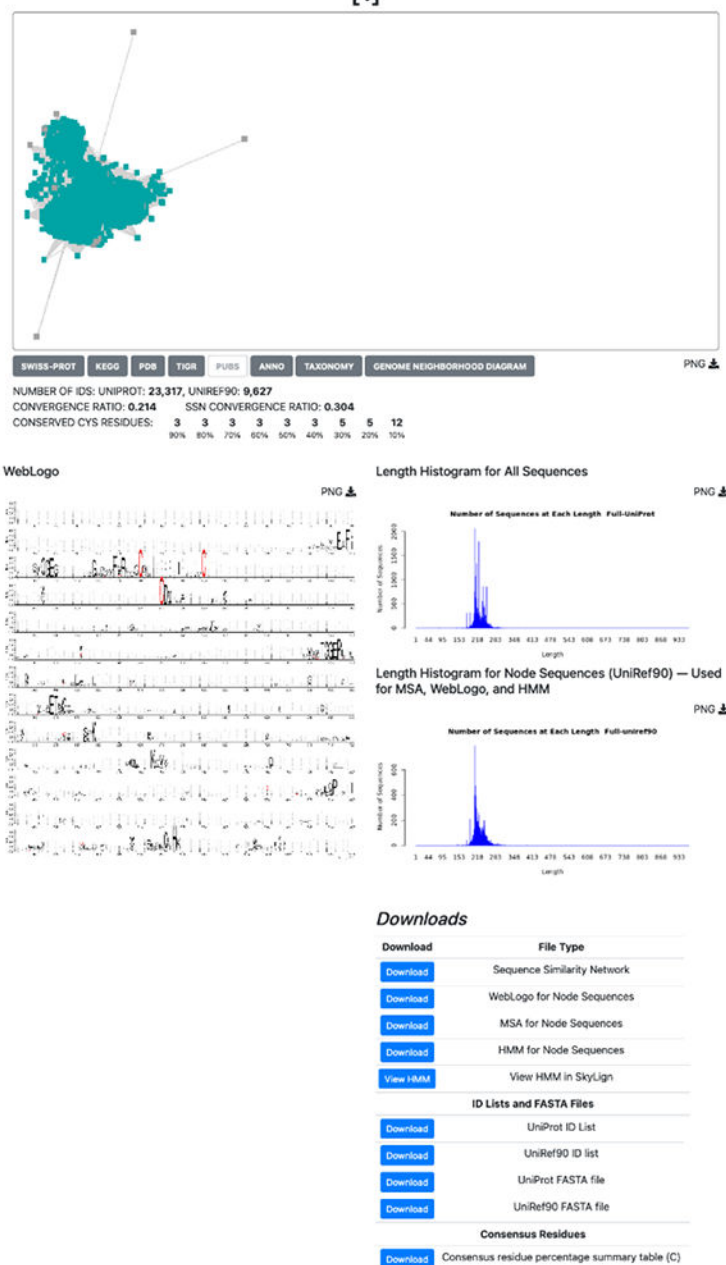
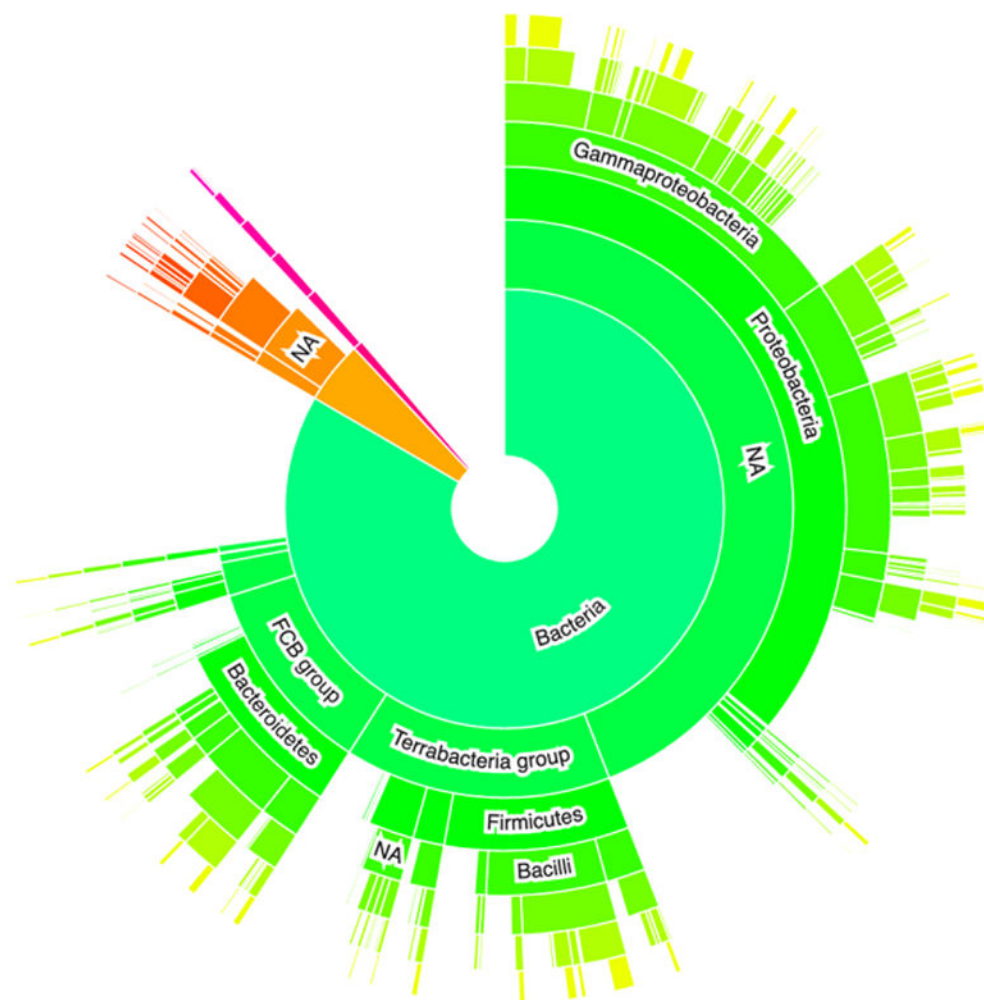


Figure 7. An example of an Explore page using Megacluster-3-1, the 7-carboxy-7-deazaguanine synthase-like radical SAM proteins (SFLD subgroup 1).

Species in Cluster



23,247 UNIPROT IDS VISIBLE

Click on a region to zoom into that part of the taxonomic hierarchy. Clicking on the center circle will zoom out to the next highest level.

Figure 8. Example of a taxonomy sunburst display. Shown is Megacluster-3-1: 7-carboxy-7-deazaguanine synthase-like (SFLD subgroup 1). Green, bacteria; orange, archaea; magenta, eukaryota.



Figure 9.
A representative GND Viewer page: Cluster-6 in the “diced” SSN for Megacluster-1-1 generated with a minimum edge alignment score threshold of 60.

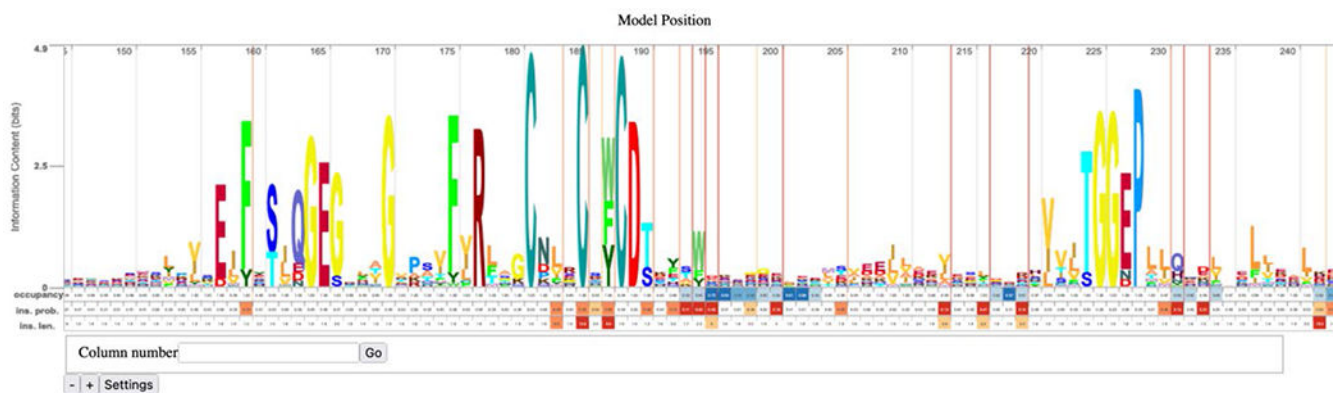


Figure 10.

A portion of the Skyline display of the HMM for Megacluster-3-1, 7-carboxy-7-deazaguanine synthase-like, SFLD subgroup 1, showing the conserved CX₃CX₂C motif that binds SAM.



Figure 11.
Representative diced SSNs for Megacluster-1-1 (SPASM/twitch domain-containing).

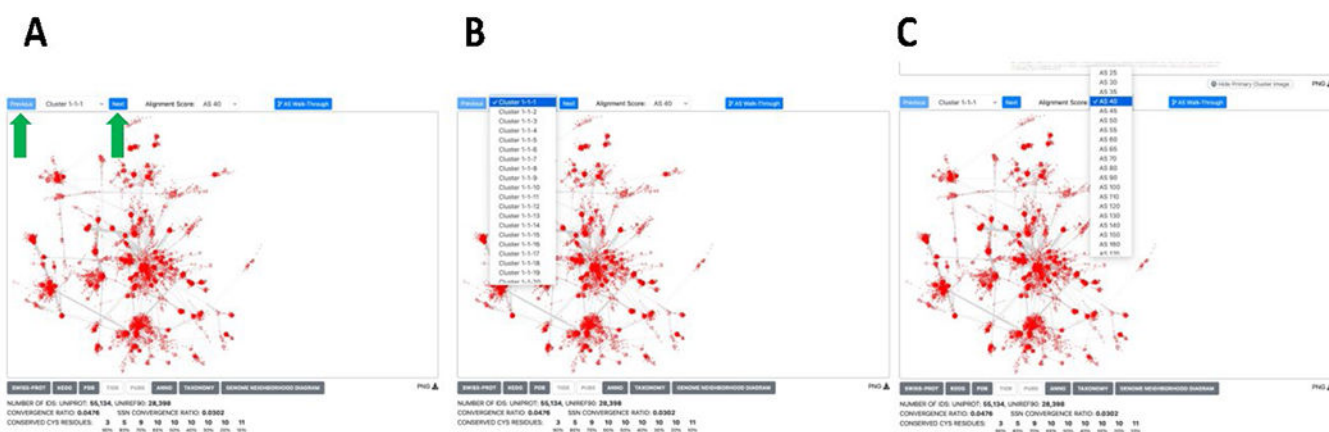


Figure 12.

Navigation through diced megaclusters. Panel A, “Previous” and “Next” navigation buttons (green arrows) direct the cluster selection backward and forward, respectively, in the selected diced SSN. Panel B, “Cluster” drop down menu allows the user to select any cluster in the currently viewed diced SSN. Panel C, Alignment score drop down menu allows the user to view the cluster in the selected diced SSN.

Alignment Score Walk-Through

Cluster ID	Num Nodes	Conv. Ratio	SwissProt	Annotation
Previous Cluster (AS35)				
Cluster-1-1-1	34396	0.0383	<ul style="list-style-type: none"> + AdaMet-dependent heme synthase + Anaerobic sulfatase-maturing enzyme ShortAnSME + Anaerobic sulfatase-maturing enzyme homolog AaIB ShortAnSME homolog + Anaerobic sulfatase-maturing enzyme homolog 1d6M ShortAnSME homolog + Antilisterial bacteriocin subtilisin biosynthesis protein ABA + Fe-coproporphyrin III synthase + PqgA peptide cyclase + Putative mycofactacin radical SAM maturase MFC + S-adenosyl-L-methionine-dependent 2-deoxy-scylo-inosamine dehydrogenase + Sporulation killing factor maturation protein SsFB + Tungsten-containing aldehyde ferredoxin oxidoreductase cofactor-modifying protein + Uncharacterized protein AF_2204 + Uncharacterized protein MJ0907 + Uncharacterized protein MTH_114 + Uncharacterized protein st1766 	<ul style="list-style-type: none"> + ADA0E3RYH2 + ADA1H7406 + ADA1I5E23 + ADA378YSA1 + ADP149 + A1B2Q7 + A3D0W1 + A4VXX3 + B8J367 + C1TQ82 + C3HC15 + C3HC16 + D0QZJ5 + F5A108 + F5A109 + GOL012 + GOL027 + O31423 + P27507 + P71011 + P71517 + P9WJ78 + P9WJ79 + Q07TH1 + Q46CH7 + Q51741 + Q841K9 + Q8D169 + Q8G907 + Q8R6P9 + Q8RAM6 + Q9X758
Next Clusters (AS45)				
Cluster-1-1-1	10995	0.0918	<ul style="list-style-type: none"> + AdaMet-dependent heme synthase + Fe-coproporphyrin III synthase + Putative mycofactacin radical SAM maturase MFC 	<ul style="list-style-type: none"> + ADA1I5E23 + ADP149 + B8J367 + P9WJ78 + P9WJ79 + Q46CH7 + Q8D169
Cluster-1-1-2	3728	0.464	<ul style="list-style-type: none"> + Anaerobic sulfatase-maturing enzyme ShortAnSME + Anaerobic sulfatase-maturing enzyme homolog AaIB ShortAnSME homolog + Anaerobic sulfatase-maturing enzyme homolog 1d6M ShortAnSME homolog + Uncharacterized protein MTH_114 + Uncharacterized protein st1766 	<ul style="list-style-type: none"> + ADA0E3RYH2 + ADA1H7406 + Q07TH1 + Q9X758

Cluster ID	Num Nodes	Conv. Ratio	SwissProt	Annotation
Cluster-1-1-3	2674	0.0638		
Cluster-1-1-4	2429	0.326		<ul style="list-style-type: none"> + ADA378YSA1 + A1B2Q7 + A3D0W1 + C3HC15 + F5A108 + GOL012 + GOL027 + Q8RAM6
Cluster-1-1-6	1610	0.993	+ PqgA peptide cyclase	<ul style="list-style-type: none"> + P27507 + P71517
Cluster-1-1-8	1323	0.101		+ Q8R6P9
Cluster-1-1-9	1136	0.139		
Cluster-1-1-15	625	0.663		
Cluster-1-1-17	474	0.702		
Cluster-1-1-20	429	0.352	+ Tungsten-containing aldehyde ferredoxin oxidoreductase cofactor-modifying protein	+ Q51741
Cluster-1-1-21	419	0.152		
Cluster-1-1-26	267	0.206		
Cluster-1-1-32	176	0.406		
Cluster-1-1-36	129	0.27	+ Uncharacterized protein MJ0907	
Cluster-1-1-38	128	0.192		
Cluster-1-1-46	107	0.71		
Cluster-1-1-50	103	0.994		
Cluster-1-1-61	88	0.31	+ Antilisterial bacteriocin subtilisin biosynthesis protein ABA	<ul style="list-style-type: none"> + A3VXX3 + D0QZJ5 + P71011
Cluster-1-1-64	77	0.524		
Cluster-1-1-68	73	0.493		
Cluster-1-1-81	62	0.401		
Cluster-1-1-105	47	0.229		
Cluster-1-1-113	42	0.321		
Cluster-1-1-117	41	0.0734		

Figure 13. AS Walk-Through pop-up window showing identity of the progenitor cluster (“Previous Cluster”) and progeny clusters (“Next Clusters”).

RADICALSAM.ORG EXPLORE SEARCH SUBMIT TUTORIALS ABOUT CONTACT

Search

Find by UniProt ID

Input a UniProt ID to identify its cluster.

For all but Megacluster-1-1 (SFLD Subgroup-17, SPASM/Twitch domain), Megacluster-2-1 (SFLD Subgroup 5, B12-binding domain), Megacluster-2-2 (SFLD Subgroup 2, anaerobic coproporphyrinogen-III oxidase-like), or Cluster-7 ((SFLD Subgroup 16, PLP-dependent), the search opens the Explore page for the cluster that contains the user-specified UniProt ID.

For Megacluster-1-1 (SFLD Subgroup-17, SPASM/Twitch domain), Megacluster-2-1 (SFLD Subgroup 5, B12-binding domain), Megacluster-2-2 (SFLD Subgroup 2, anaerobic coproporphyrinogen-III oxidase-like), or Cluster-7 ((SFLD Subgroup 16, PLP-dependent), the search identifies the cluster (if ≥ 3 nodes/UniRef50 IDs) in each "diced" SSN that contains the ID. The number of UniProt IDs, number of cluster nodes, and UniProt ID convergence ratio (CR; described on the Subgroups tab) are provided for each identified cluster.

In the generation of the megaclusters and clusters, some UniProt IDs for a singleton may be deleted. The Search will report: "ID not found".

[Find Cluster](#)

Find by Sequence

Input a sequence to find clusters that contain homologues. The sequence is used to query the HMMs for the clusters (≥ 3 UniRef IDs/nodes). The results reports matches for the "top" three clusters if the e-value is $\leq 1e-10$. The cluster is a link to the Explore page for the cluster.

For Megacluster-1-1 (SFLD Subgroup-17, SPASM/Twitch domain), Megacluster-2-1 (SFLD Subgroup 5, B12-binding domain), Megacluster-2-2 (SFLD Subgroup 2, anaerobic coproporphyrinogen-III oxidase-like), or Cluster-7 ((SFLD Subgroup 16, PLP-dependent), the second section reports matches for clusters in the "diced" SSNs; the clusters with the three smallest e-values are listed. The number of UniProt IDs, number of cluster nodes, and UniProt ID convergence ratio (CR; described on the Subgroups tab) are provided for each identified cluster.

The "Exploring Subgroups" subtab under the "Functionally Diverse Subgroups" tab provides advice about interpreting the search results for these subgroups.

[Find Clusters](#)

GND Lookup

The [EFI-GNT web tools](#) allow users to lookup genome neighborhood diagrams (GNDs) for lists of UniProt IDs. Users may find it convenient to be able to access the GNDs for members of the RSS within RadicalSAM.org.

The GND Viewer can be accessed with the button below. The input is a list of UniProt IDs. The GNDs will be displayed.

[GND Viewer](#)

Find by Taxonomy

Input the genus/species/strain for an organism.

If only the genus is entered, a pop-up list of matching genus-species-strains is provided for selection of the desired genus/species/strain. If the genus and species are entered, a pop-up list of matching genus-species-strains is provided for selection of the desired genus/species/strain.

The search provides a list of sequences in the RSS. The list provides the UniProt ID (link to the UniProt page for sequence), UniProt description, organism name, UniProt annotation status (SwissProt or TrEMBL), and link to its Explore page.

[Find Sequences](#)

[CLICK HERE TO CONTACT US FOR HELP, REPORTING ISSUES, OR SUGGESTIONS](#)
THIS WEBSITE USES DATA DERIVED FROM THE INTERPRO AND UNIPROT DATABASES.

Figure 14.

Search functions: "Find by UniProt ID" identifies the cluster(s) containing the user-specified accession ID; "Find by Sequence" identifies the cluster(s) with the best HMM match to the user-specified sequence; "GND Lookup" provides the GND(s) for the user-specified UniProt ID(s); "Find by Taxonomy" provides a list of the UniProt accession IDs and their clusters for the user-specified genus/species.

RADICALSAM.ORG EXPLORE SEARCH **SUBMIT** TUTORIALS ABOUT CONTACT

Submit

We encourage user-submitted annotations for cluster or individual sequence. Upon review and approval, these annotation will be included on results pages for individual clusters.

Your name
Enter name

Your email
Enter email
Your email address will never be shared.

Cluster ID
Enter cluster ID
If you don't know this, then please provide details below.

Function/Annotation

Provide the protein function that is associated with the described function.

Accession ID
Enter UniProt/NCBI accession ID
Enter the UniProt (preferred) or NCBI accession ID that is associated with your submission. If this is unknown, please provide details below.

Sequence

Provide the protein sequence that is associated with the described function.

Publication DOI
Enter DOI/publication identifier
If the publication DOI is not available, provide a link to the publication, or provide details below.

Details

Provide additional details regarding the sequence, cluster, publication, annotation, or other information.

[Read Terms of Service and Disclaimer](#)
 I have read and agree to the terms of service.

[Submit Annotation](#)

Figure 15.
The Submit page for community submission of enzymatic activities and metabolic functions.

Send Feedback or Questions about RadicalSAM.org

Feedback, questions, or requests can be submitted to the team below.

Your name

Your email

Your email address will never be shared.

Your institution

Comments:

Fill in comments.

Figure 16.

The Contact page form for submitting feedback or questions.

Table 1.SFLD Subgroups, Names, and (Mega)Clusters in [RadicalSAM.org](https://www.RadicalSAM.org)

Subgroup	Subgroup Name	RS.org (Mega)cluster
1	7-carboxy-7-deazaguanine synthase-like	Megacluster-3-1
2	Coproporphyrinogen III oxidase-like	Megacluster-2-2
3	Antiviral proteins (viperin)	Megacluster-1-5
4	Avilamycin synthase	Megacluster-1
5	B12-binding domain containing	Megacluster-2-1
6	Biotin and thiazole synthase domain containing	Megacluster-4
7	DesII-like	Megacluster-1-8
8	ELP3/YhcC	Megacluster-2-4, -2-5
9	F420, menaquinone cofactor biosynthesis	Megacluster-4-2
10	FeMo-cofactor biosynthesis protein	Megacluster-1-4
11	Lipoyl synthase like	Cluster-8
12	Methylthiotransferase	Megacluster-2-3
13	Methyltransferase Class A	Cluster-6
14	Methyltransferase Class D	Megacluster-1-3
15	Organic radical activating enzymes	Megacluster-3
16	PLP-dependent	Megacluster-7
17	SPASM/twitch domain containing	Megacluster-1-1
18	Spectinomycin biosynthesis	Megacluster-1-1
19	Spore photoproduct lyase	Megacluster-5-3
20	tRNA wybutosine-synthesizing	Cluster-10
	Protein MJ0683-like	Megacluster-5-1
	Uncharacterized protein family UPF0313	Cluster-10
	DUF5131	Megacluster-5-2
	3',8-Cyclase/Mo cofactor synthesis	Megacluster-1-2