Genomic Enzymology: Discovery of Novel Enzymes in Novel Pathways

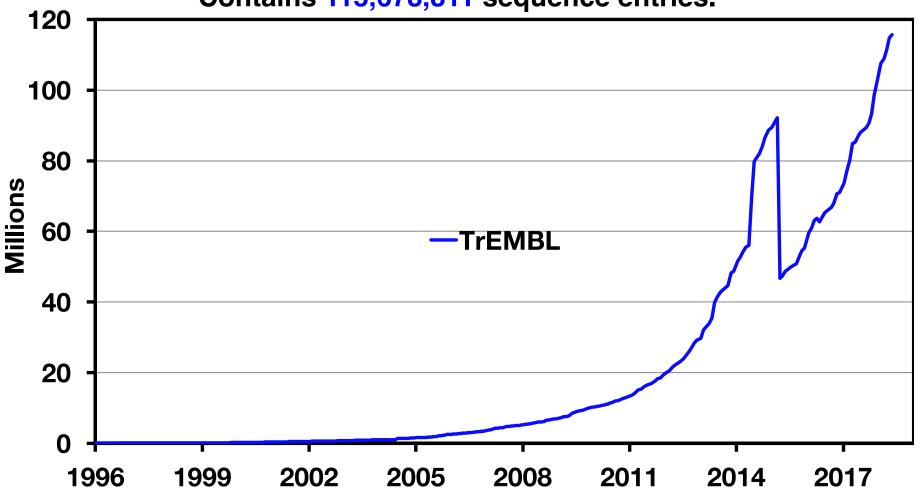
John A. Gerlt
University of Illinois, Urbana-Champaign

Penn State Bioinorganic Workshop June 2, 2018

The number of protein sequences is "exploding"!



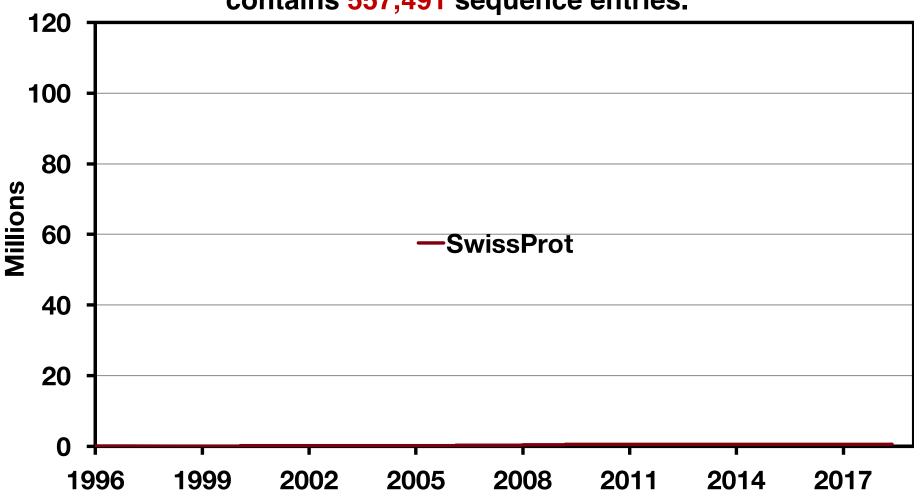
Release 2018_05 of UniProtKB/TrEMBL Contains 115,678,811 sequence entries.



But the number of curated annotations is lagging!



Release 2018_05 of UniProtKB/SwissProt contains 557,491 sequence entries.

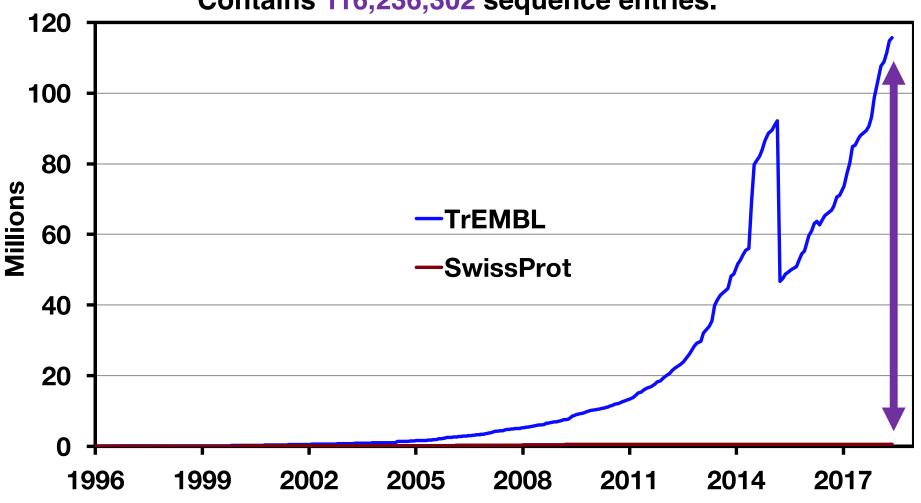


Perhaps 50% have unknown or uncertain functions



Release 2018_05 of UniProtKB

Contains 116,236,302 sequence entries.

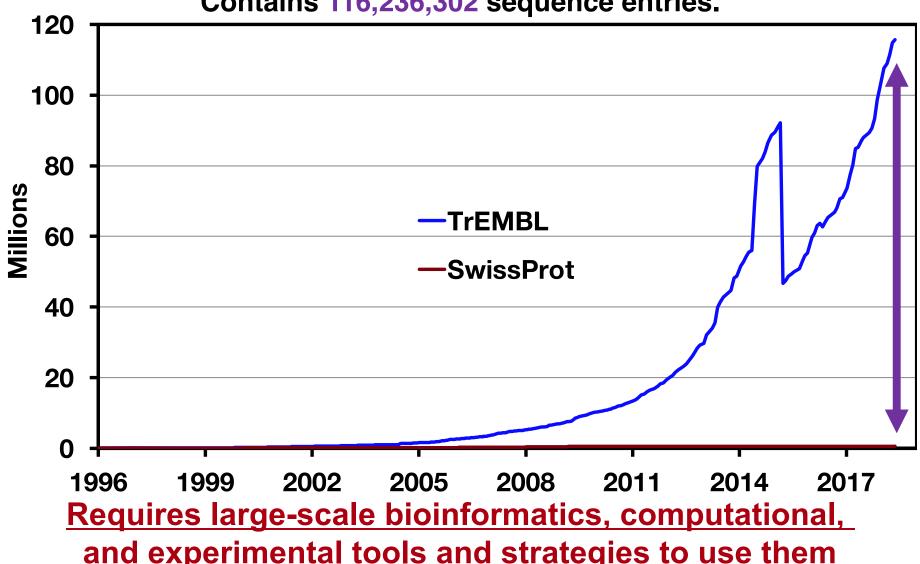






Release 2018_05 of UniProtKB

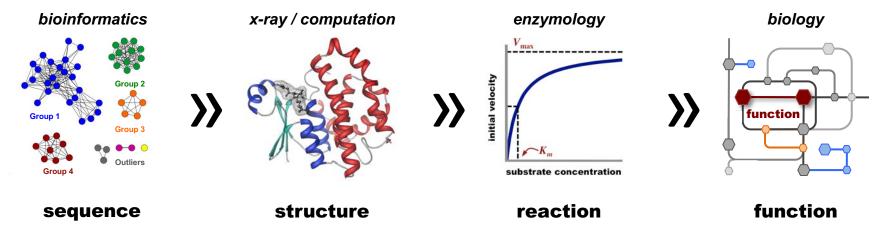
Contains 116,236,302 sequence entries.



and experimental tools and strategies to use them

U54GM093342: Enzyme Function Initiative





Albert Einstein
Steve Almo

Boston University
Karen Allen

EMBL-EBI
Alex Bateman

University of Illinois

John Gerlt

John Cronan

Pennsylvania State
Squire Booker

Sanford-Burnham MRI
Andrei Osterman
Dmitry Rodionov

UCSF
Matthew Jacobson
Andrej Sali
Brian Shoichet

University of UtahC. Dale Poulter

University of VirginiaWladek Minor

"Genomic Enzymology" Web Tools



Biochemistry

Perspective

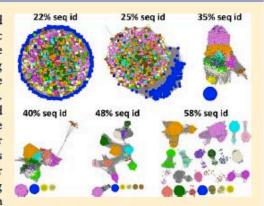
pubs.acs.org/biochemistry

Genomic Enzymology: Web Tools for Leveraging Protein Family Sequence—Function Space and Genome Context to Discover Novel Functions

John A. Gerlt*0

Departments of Biochemistry and Chemistry, Institute for Genomic Biology, University of Illinois, Urbana-Champaign Urbana, Illinois 61801, United States

ABSTRACT: The exponentially increasing number of protein and nucleic acid sequences provides opportunities to discover novel enzymes, metabolic pathways, and metabolites/natural products, thereby adding to our knowledge of biochemistry and biology. The challenge has evolved from generating sequence information to mining the databases to integrating and leveraging the available information, i.e., the availability of "genomic enzymology" web tools. Web tools that allow identification of biosynthetic gene clusters are widely used by the natural products/synthetic biology community, thereby facilitating the discovery of novel natural products and the enzymes responsible for their biosynthesis. However, many novel enzymes with interesting mechanisms participate in uncharacterized small-molecule metabolic pathways; their discovery and functional characterization also can be accomplished by leveraging information in protein and nucleic acid databases. This Perspective focuses on



two genomic enzymology web tools that assist the discovery novel metabolic pathways: (1) Enzyme Function Initiative-Enzyme Similarity Tool (EFI-EST) for generating sequence similarity networks to visualize and analyze sequence—function space in protein families and (2) Enzyme Function Initiative-Genome Neighborhood Tool (EFI-GNT) for generating genome neighborhood networks to visualize and analyze the genome context in microbial and fungal genomes. Both tools have been adapted to other applications to facilitate target selection for enzyme discovery and functional characterization. As the natural products community has demonstrated, the enzymology community needs to embrace the essential role of web tools that allow the protein and genome sequence databases to be leveraged for novel insights into enzymological problems.

Genomic enzymology strategy for functional assignment



1. Explore sequence function-space in a protein family, segregating family into isofunctional clusters using known functions (SwissProt, literature)

Sequence similarity networks (SSNs) using EFI-EST

2. Predict functions of isofunctional clusters in novel metabolic pathways using genome context (types of reactions enabled by proximal enzymes)

Genome neighborhood networks (GNNs) using EFI-GNT

3. Prioritize clusters for functional characterization by (metagenome) abundance

Chemically guided functional profiling (CGFP)

Today



- 1. Target selection using sequence similarity networks: EFI-EST web tool
- 2. Pathway discovery using genome neighborhood networks: EFI-GNT web tool
- 3. Example: using SSNs and GNNs to discover catabolic pathways for D-apiose
- 4. Target prioritization using chemically guided functional profiling: CGFP-ShortBRED beta tool

Today



- 1. Target selection using sequence similarity networks: EFI-EST web tool
- 2. Pathway discovery using genome neighborhood networks: EFI-GNT web tool
- 3. Example: using SSNs and GNNs to discover catabolic pathways for D-apiose
- 4. Target prioritization using chemically guided functional profiling: CGFP-ShortBRED beta tool

Sequence Similarity Networks (SSNs): Atkinson et al.



PLoS one

OPEN & ACCESS Freely available online



Relationships Across Diverse Protein Superfamilies

Holly J. Atkinson^{1,2}, John H. Morris³, Thomas E. Ferrin^{2,3,4}, Patricia C. Babbitt^{2,3,4}*

1 Graduate Program in Biological and Medical Informatics, University of California San Francisco, San Francisco, Califomia, United States of America, 2 Institute for Quantitative Biosciences, University of Califomia San Francisco, San Francisco, Califomia, United States of America, 3 Department of Pharmaceutical Chemistry, University of California San Francisco, San Francisco, California, United States of America, 4 Department of Biopharmaceutical Sciences, University of California San Francisco, San Francisco, California, United States of America

Abstract

The dramatic increase in heterogeneous types of biological data—in particular, the abundance of new protein sequences requires fast and user-friendly methods for organizing this information in a way that enables functional inference. The most widely used strategy to link sequence or structure to function, homology-based function prediction, relies on the fundamental assumption that sequence or structural similarity implies functional similarity. New tools that extend this approach are still urgently needed to associate sequence data with biological information in ways that accommodate the real complexity of the problem, while being accessible to experimental as well as computational biologists. To address this, we have examined the application of sequence similarity networks for visualizing functional trends across protein superfamilies from the context of sequence similarity. Using three large groups of homologous proteins of varying types of structural and functional diversity—GPCRs and kinases from humans, and the crotonase superfamily of enzymes—we show that overlaying networks with orthogonal information is a powerful approach for observing functional themes and revealing outliers. In comparison to other primary methods, networks provide both a good representation of group-wise sequence similarity relationships and a strong visual and quantitative correlation with phylogenetic trees, while enabling analysis and visualization of much larger sets of sequences than trees or multiple sequence alignments can easily accommodate. We also define important limitations and caveats in the application of these networks. As a broadly accessible and effective tool for the exploration of protein superfamilies, sequence similarity networks show great potential for generating testable hypotheses about protein structure-function relationships.

Citation: Atkinson HJ, Morris JH, Ferrin TE, Babbitt PC (2009) Using Sequence Similarity Networks for Visualization of Relationships Across Diverse Protein Superfamilies. PLoS ONE 4(2): e4345. doi:10.1371/journal.pone.0004345

Editor: I. King Jordan, Georgia Institute of Technology, United States of America

Received September 10, 2008; Accepted December 10, 2008; Published February 3, 2009

This is an open-access article distributed under the terms of the Creative Commons Public Domain declaration which stipulates that, once placed in the public domain, this work may be freely reproduced, distributed, transmitted, modified, built upon, or otherwise used by anyone for any lawful purpose.

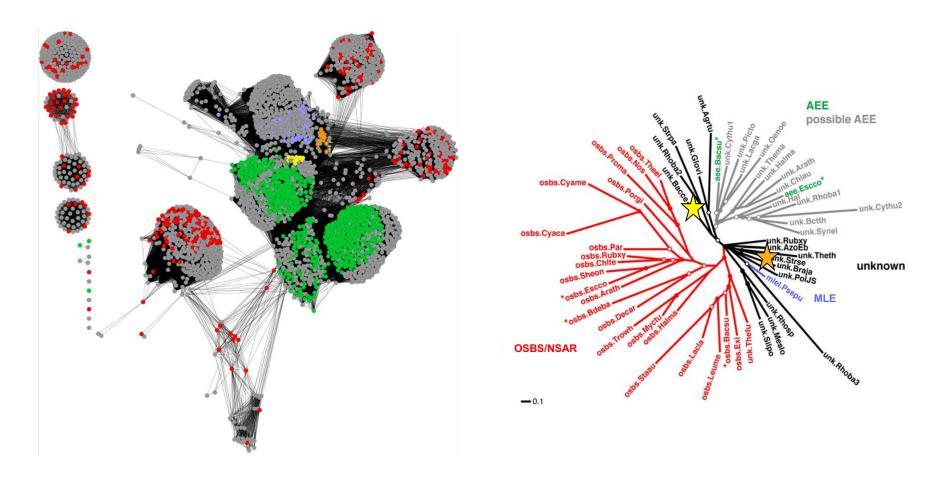
Funding: This work was supported by NIH grant GM60595 and NSF grant DBI 0640476 to P.C.B. and P41 RR01081 to T.E.F.. H.J.A. received support from NIH grant T32 GM067547. Initial exploration of sequence similarity networks used the enclases and amidohydrolases superfamilies as example data sets, and was supported by P01 GM071790 to P.C.B.. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* E-mail: babbitt@cgl.ucsf.edu

Sequence similarity networks (SSNs) vs. dendrograms: large-scale visualization tool for target selection

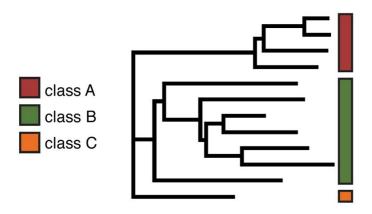




Families are easier to visualize and connect with SSNs, so function predictions are easier to formulate and explore

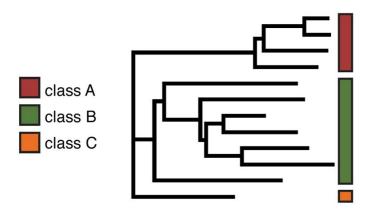
Dendrograms/trees for sequence relationships





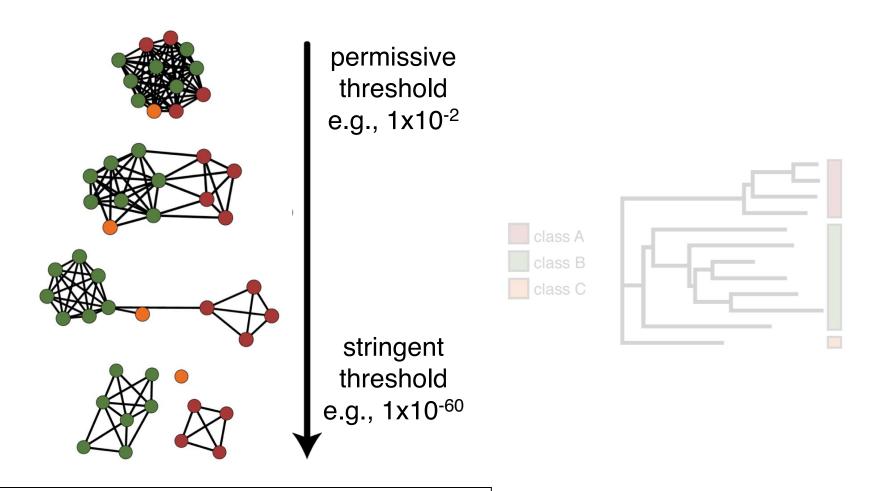


<u>Connectivity: structure-based</u> <u>multiple sequence alignments (MSAs)</u>



Sequence similarity networks



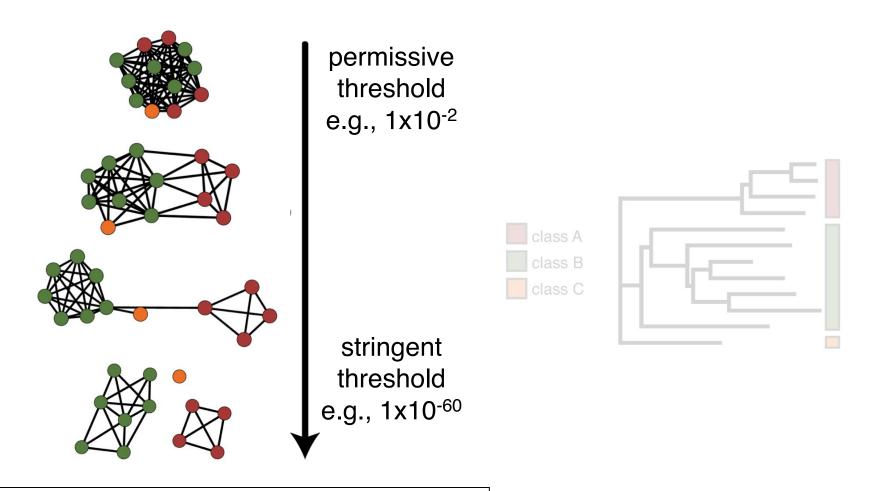


node (circle) = sequence

edge (line) = connection less than a
user-defined similarity score (e-value)

Connectivity: all-by-all BLASTP e-values



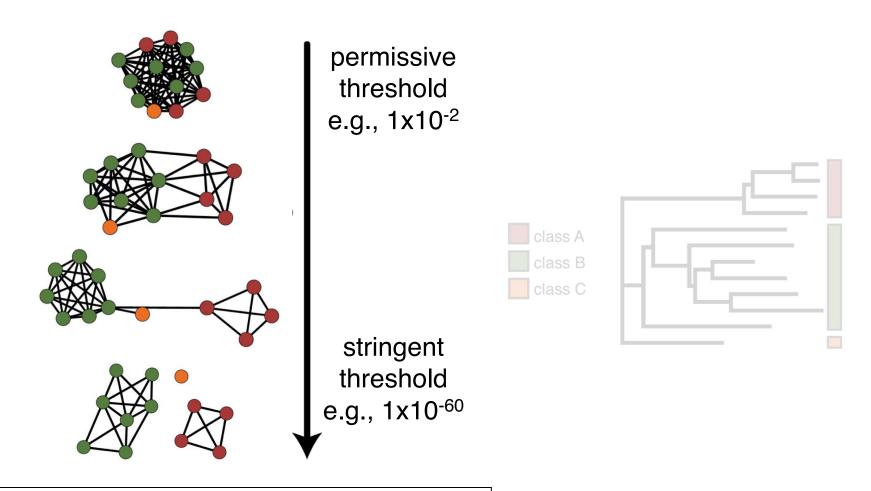


node (circle) = sequence

edge (line) = connection less than a
user-defined similarity score (e-value)

Faster to calculate than MSAs



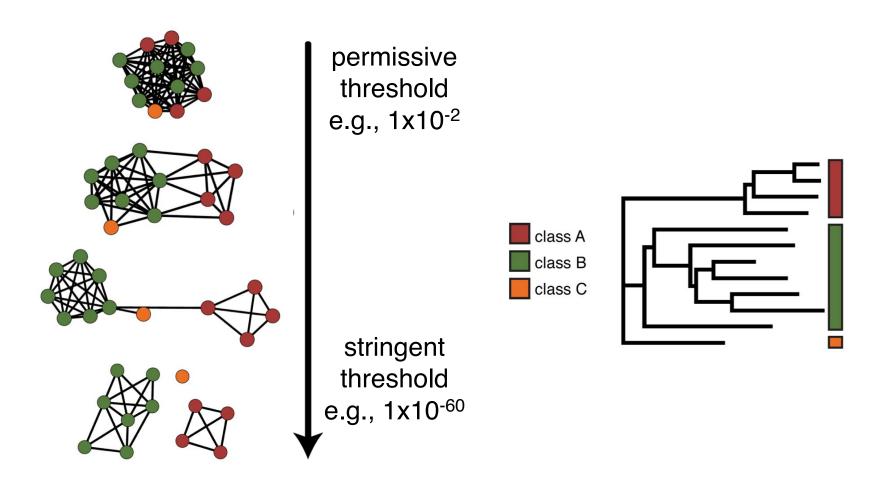


node (circle) = sequence

edge (line) = connection less than a
user-defined similarity score (e-value)

Qualitatively similar results



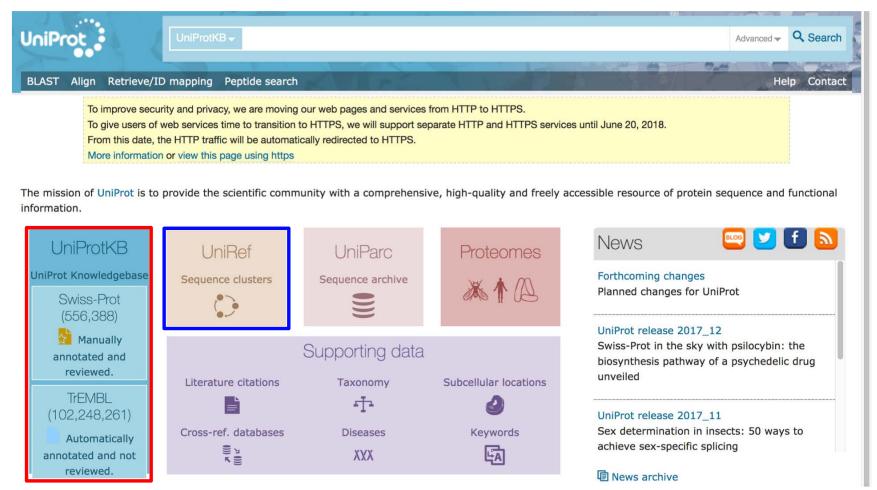


node (circle) = sequence

edge (line) = connection less than a
user-defined similarity score (e-value)

EFI-EST uses sequences from UniProt (EMBL/EBI)





http://www.uniprot.org/

UniProt databases used by EFI-EST



UniProt Knowledgebase (UniProtKB): curated protein information, including function, classification, and cross-reference.

UniProtKB/Swiss-Prot: manually annotated

UniProtKB/TrEMBL: automatically annotated

UniProt Reference Clusters (UniRef): clustered sets of sequences from UniProtKB and UniProt Archive records to obtain complete coverage of sequence space at several resolutions while hiding redundant sequence.

UniRef90: 90% sequence identity over 80% of the length of the longest sequence (used by Option B of EFI-EST). On average, a two-fold reduction in number of sequences, providing a four-fold reduction in running time.

Why UniProt, not NCBI? Annotations can be changed



Contact

Submit updates or corrections to UniProt

- > We are actively seeking any type of updates and/or corrections whether they have been published or not.
- > Please see the guidelines for submitting updates.

Other ways to contact us...

Send us general questions and suggestions Submit new protein sequence data Email us on help@uniprot.org

Name :									
E-mail:									
Subject:	UniProtKB/TrEMBL D2JA13 entry update r								
This is currently an unreviewed entry. We are going to update it and integrate it as a reviewed entry in UniProtKB/Swiss-Prot. Please post your update in the text field below.									
Message:									

https://www.uniprot.org/update

EFI-EST uses families from Pfam and InterPro





Pfam: conserved protein families based a seed alignment of representative sequences that is used to generate a profile hidden Markov model (HMM). Includes families and clans (superfamilies of multiple families).

16,712 families in Pfam 31.0 (March 2017)



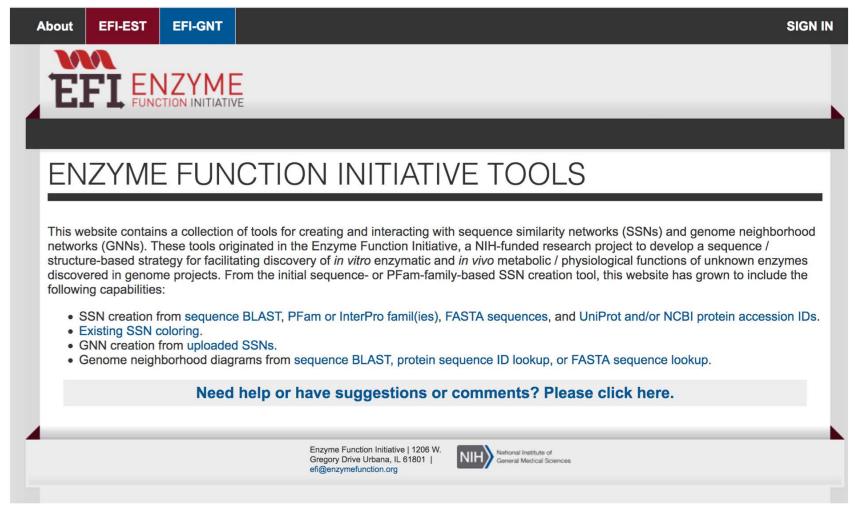
InterPro: aggregates of families, domains, and sites curated by 14 databases, e.g., Pfam, CATH-Gene3D, ProSite, SFLD, TIGRFAMs, etc.

33,947 domains, families, and sites in InterPro 68 (April 28, 2018)

http://pfam.xfam.org/; www.ebi.ac.uk/interpro/

EFI-EST and EFI-GNT Web Tools

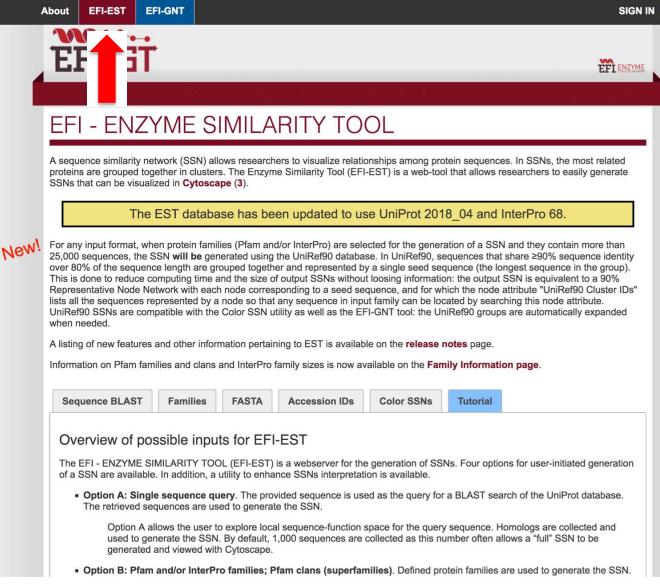




https://efi.igb.illinois.edu/

EFI-EST

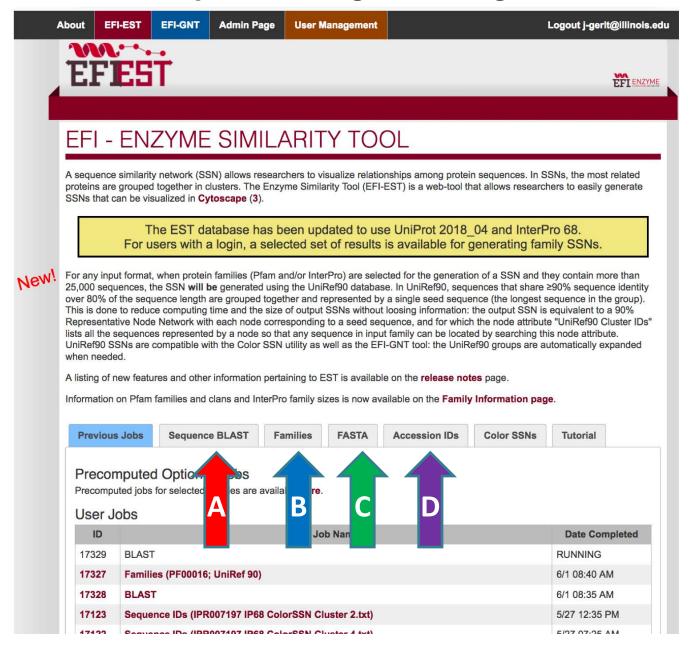




https://efi.igb.illinois.edu/efi-est/

Four options for generating SSNs





Four options for generating SSNs



- Option A: Homologues in UniProt to user-supplied sequences (collected by BLAST)
- Option B: User-supplied Pfam and/or InterPro families from UniProt, using full length sequences or domains
- Option C: User-supplied FASTA file; FASTA header can be "read" for UniProt and/or NCBI IDs
- Option D: User-supplied file of UniProt and/or NCBI IDs

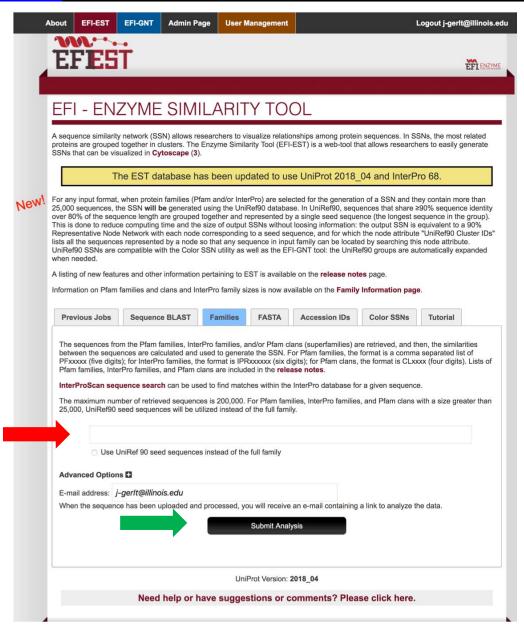
Option A: Homologues collected by BLAST



Ab	out E	FI-EST	EFI-GNT	Admin Pag	ge User N	lanagement		Î	Logout j-gerl	t@illinoi		
1	EF	ES.	ř _							EFI ENZ		
	EFI - ENZYME SIMILARITY TOOL A sequence similarity network (SSN) allows researchers to visualize relationships among protein sequences. In SSNs, the most relat proteins are grouped together in clusters. The Enzyme Similarity Tool (EFI-EST) is a web-tool that allows researchers to easily gener SSNs that can be visualized in Cytoscape (3). The EST database has been updated to use UniProt 2018_04 and InterPro 68.											
p												
F li U V	25,000 sequences, the SSN will be generated using the UniRef90 database. In UniRef90, sequences that share ≥90% sequence identity over 80% of the sequence length are grouped together and represented by a single seed sequence (the longest sequence in the group). This is done to reduce computing time and the size of output SSNs without loosing information: the output SSN is equivalent to a 90% Representative Node Network with each node corresponding to a seed sequence, and for which the node attribute "UniRef90 Cluster IDs lists all the sequences represented by a node so that any sequence in input family can be located by searching this node attribute. UniRef90 SSNs are compatible with the Color SSN utility as well as the EFI-GNT tool: the UniRef90 groups are automatically expanded when needed. A listing of new features and other information pertaining to EST is available on the release notes page. Information on Pfam families and clans and InterPro family sizes is now available on the Family Information page.											
	Previous Jobs Sequ		Sequenc	rence BLAST Families FAST			Accession IDs	Color SSNs	Tutorial			
•	sequence	es are calc	ulated and of retrieved	used to general sequences is 1	ate the SSN. S	Submit only or	JniProt database an	without FASTA hea	der. The defa	ult		
	default:		ery E-value:	5			egative log of e-value	e for retrieving simil	ar sequences	(≥ 1;		
	Maximum Blast Sequences: 1000 If desired, include Pfam and/or InterPro families, in the analysis of your sequence. For Pfam families, the form list of PFxxxxx (five digits); for InterPro families, the format is IPRxxxxxx (five digits); for Pfam clans, the format For Pfam families, InterPro families, and Pfam clans with a size greater than 25,000, UniRef90 seed sequence of the full family.									eparated		
	☐ Use UniRef 90 seed sequences instead of the full family Advanced Family Options □											
			gerlt@illin	ois.edu								
		Section of the sectio			processed, yo	ou will receive	an e-mail containing	a link to analyze th	ne data.			

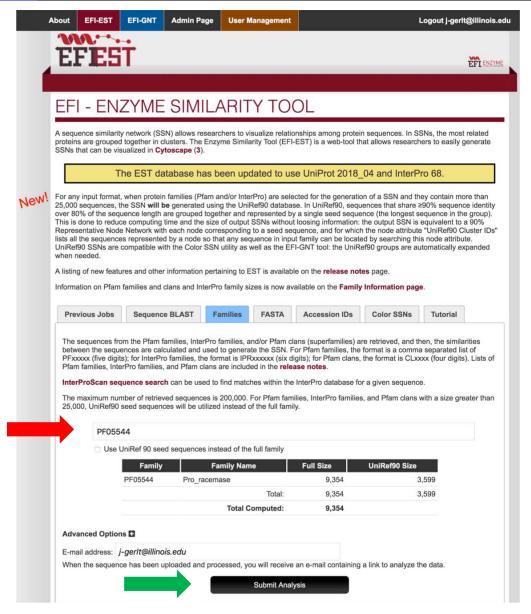
Option B: Pfam and InterPro families, Pfam clans





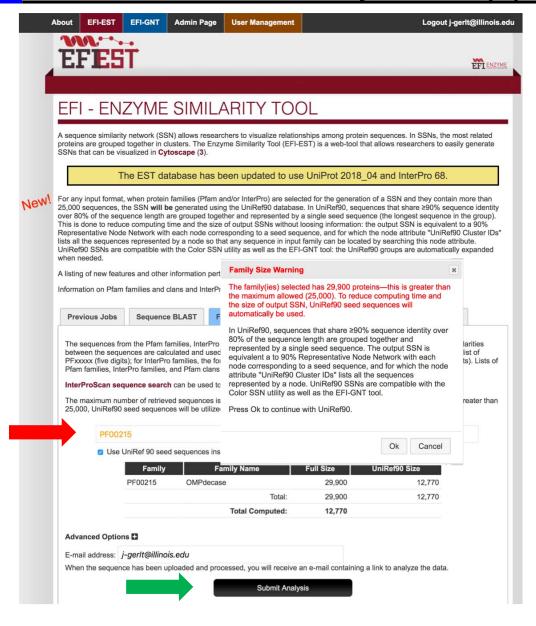
Option B: UniRef for <25K sequences (15,547 Pfams)





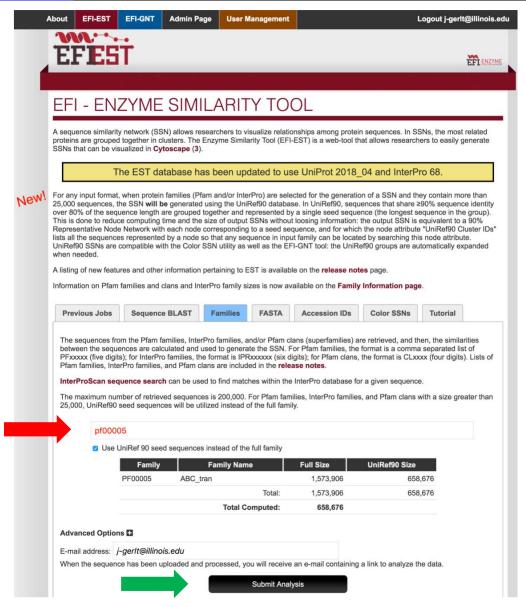
Option B: UniRef90 for ≥ 25K sequences (1,157 Pfams)





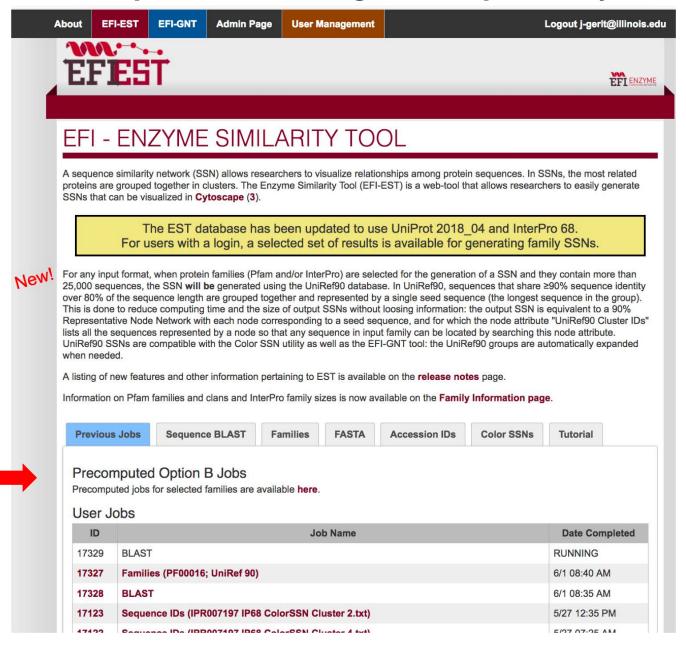
Option B: >200K UniRef90 sequences (8 Pfams)





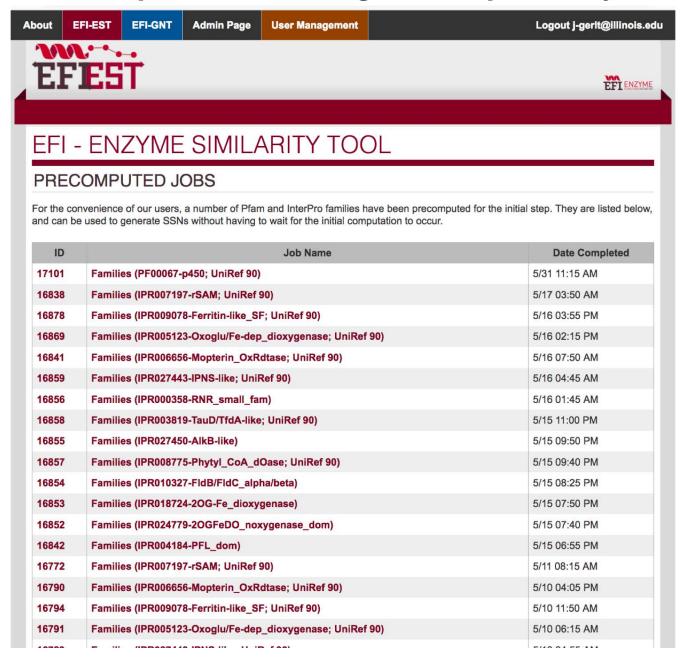
Precomputed "bioinorganic" Option B jobs





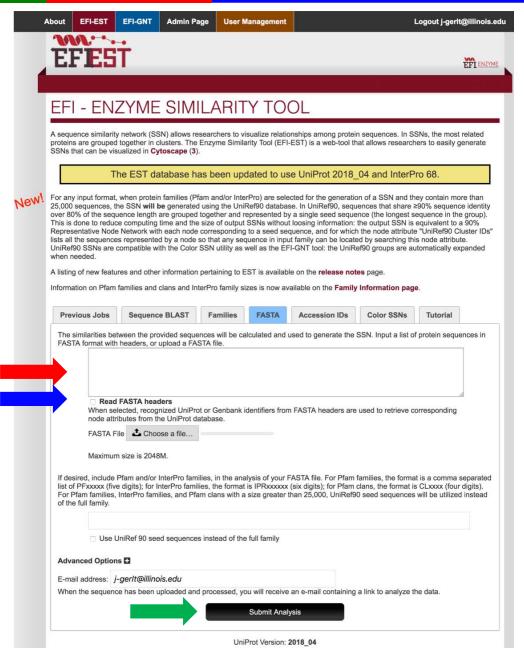
Precomputed "bioinorganic" Option B jobs





Option C: FASTA file; header can be "read" for IDs





Option C: FASTA header with UniProt IDs



>sp|A6WXX7|4HPE2_OCHA4 4-hydroxyproline 2-epimerase 2 OS=Ochrobactrum anthropi (strain ATCC 49188 / DSM 6882 / JCM 21032 / NBRC 15819 / NCTC 12168) GN=Oant_1111 PE=1 SV=1

MARHSFFCVDGHTCGNPVRLVAGGGPNLEGSTMMEKRAHFLREYDWIRTGLMFEPRGHDM MSGSILYPPTRPDCDVAVLFIETSGCLPMCGHGTIGTVTMAIEQGLVTPKTPGKLNLDTP AGLVAIEYEQNGQYVERVRLTNVPAFLYAEGLEVECPDLGNLKVDVAYGGNFYAIVEPQE NYTDMEDYSALQLIAWSPILRERLNEKYKFQHPLLPDINRLSHILWTGKPKHPEAHARNA VFYGDKAIDRSPCGTGTSARMAQLAAKGKLKPGDEFVHESIIGSLFHGRVERATEVVGQD RTLPAIIPSIAGWARMTGYNTIFIDDRDPFAHGFTVA

Option C: FASTA header with NCBI IDs



>WP_012091266.1 hydroxyproline-2-epimerase [Ochrobactrum anthropi]

>A6WXX7.1 RecName: Full=4-hydroxyproline 2-epimerase 2; Short=4Hyp 2-epimerase 2; Short=4HypE 2

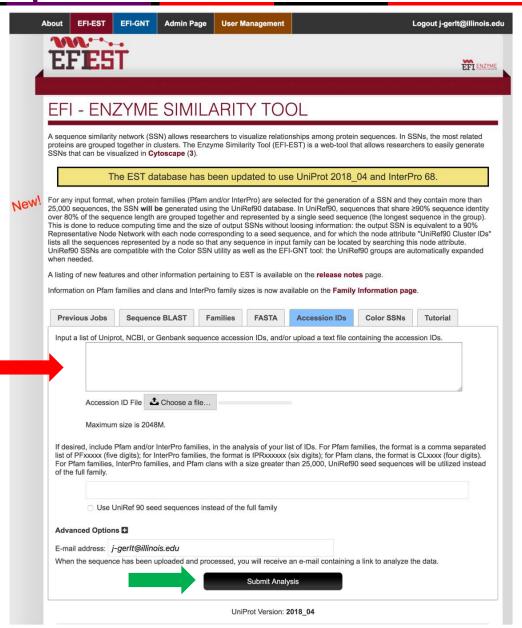
>ABS13831.1 proline racemase [Ochrobactrum anthropi ATCC 49188]

>AIK43759.1 4-hydroxyproline epimerase [Ochrobactrum anthropi]

MARHSFFCVDGHTCGNPVRLVAGGGPNLEGSTMMEKRAHFLREYDWIRTGLMFEPRGHDMMSG SILYPPTRPDCDVAVLFIETSGCLPMCGHGTIGTVTMAIEQGLVTPKTPGKLNLDTPAGLVAIEYEQNG QYVERVRLTNVPAFLYAEGLEVECPDLGNLKVDVAYGGNFYAIVEPQENYTDMEDYSALQLIAWSPILR ERLNEKYKFQHPLLPDINRLSHILWTGKPKHPEAHARNAVFYGDKAIDRSPCGTGTSARMAQLAAKG KLKPGDEFVHESIIGSLFHGRVERATEVVGQDRTLPAIIPSIAGWARMTGYNTIFIDDRDPFAHGFTVA

Option D: UniProt and/or NCBI IDs









EFEST						E	FI ENZYME
EFI - ENZY	ME SI	MIL	_ARI1	Y TC	OL	٠,	
TL - FOT 1-	4-b b b		1-1-11	HaiDart 20	47 44 1	-D 00 0	-
The EST da	labase has b	een upc	ialed to use	UniProt 20	17_11 and Inte	PIO 66.0	
INPUT S	B ENERATE ATA SET	>>	C ANALYSIS	>>>	D GENERATE NETWORKS	>>	E DOWNLOAD FILES
DATA SET COMPL	ETED						
Network Information							
Generation Summary Table Down	nload						
Date Completed Database Version					7-11 / Interpro: 66.0	n	
Input Option				Option B	7-117 Interpro: 00.0		-
Job Number				14151			
PFam/Interpro Families				PF00215			
E-Value				5			
Fraction				1			
Domain				off			
Sequence Identity				1.0			
Sequence Overlap				1			
UniRef Version				90			
Number of IDs in PFAM/InterPro	o Family			10,413			
Total Number of Nodes				10,413			
Convergence Ratio ⁺				0.323			
The convergence ratio is a mea- tetained from the BLAST (e-value rom 1.0 for sequences that are vi Parameters for SSN I To finalize the generation of an SS tetwork is needed. This will deter	s less than the spery similar (identification) Finalization SN, a similarity th	pecified the cal) to 0.0	reshold; defaul for sequences at defines which	t 5) to the total r that are very di	number of sequenc fferent (unrelated).	e pairs. The va	lue decreases
Analyze your data set /iew plots and histogram to deter	?				continuing.		
Number of Edges Histogram	Download 🖪	Preview					
ength Histogram	Download 🗈	Preview					
Alignment Length Quartile Plot	Download 🖪	Preview					
Percent Identity Quartile Plot	Download 🗈	Preview					

Data Set Completed page



Parameters for SSN Finalization

To finalize the generation of an SSN, a similarity threshold that defines which protein sequences should or should not be connected in a network is needed. This will determine the segregation of proteins into clusters.

Analyze your data set ?

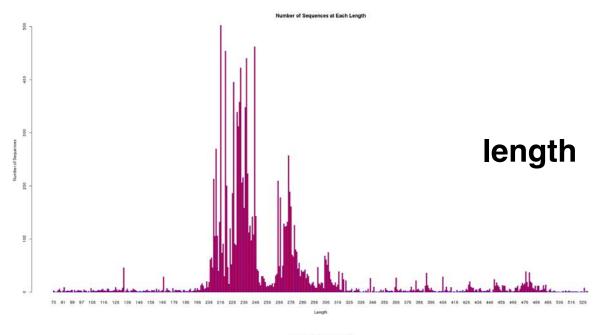
View plots and histogram to determine the appropriate lengths and alignment score before continuing.

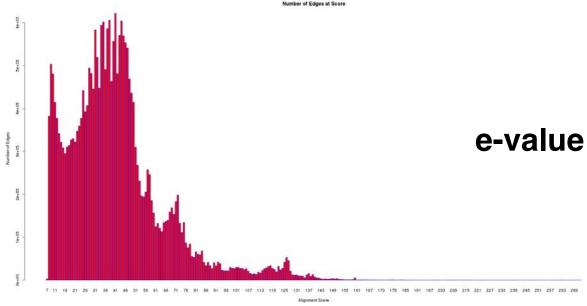


Number of Edges Histogram	Download 🛂	Preview
Length Histogram	Download D	Preview
Alignment Length Quartile Plot	Download D	Preview
Percent Identity Quartile Plot	Download 🖪	Preview

Two histograms: length and e-value

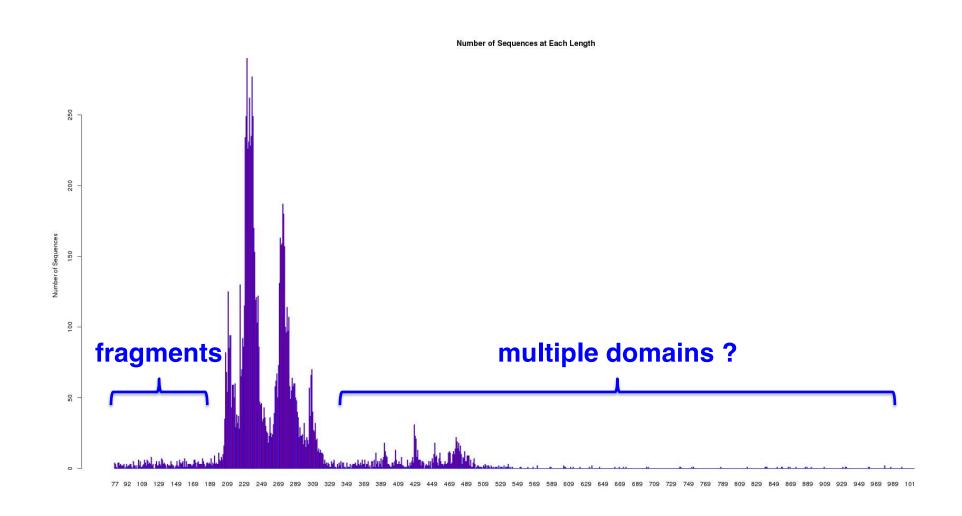






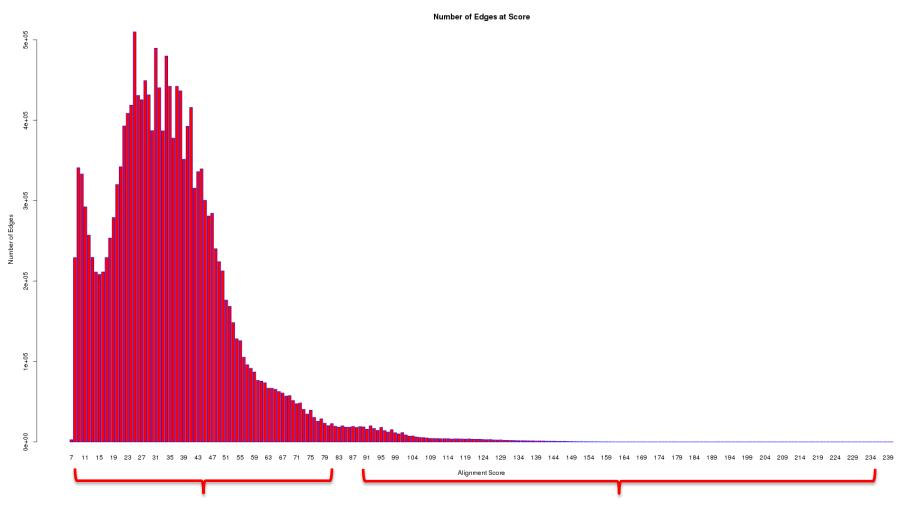
Length histogram: fragments and multiple domains





Edge histogram: sequence divergence



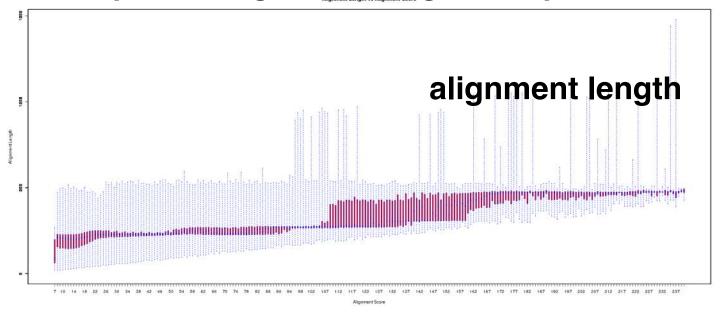


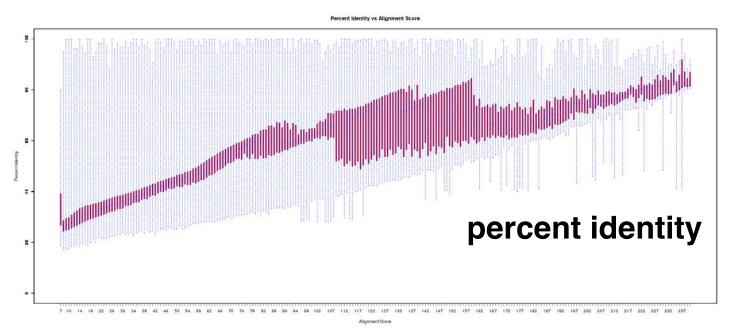
divergent connections between clusters

convergent, connections within clusters

Quartile plots: alignment length and percent identity

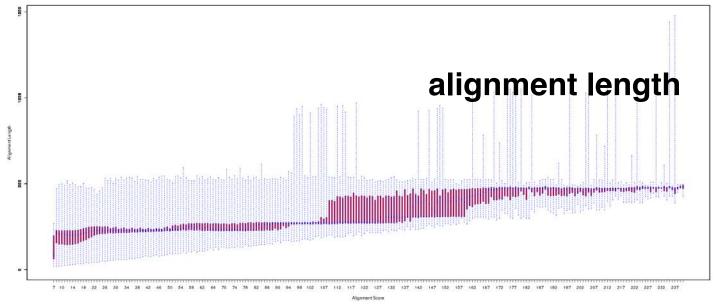


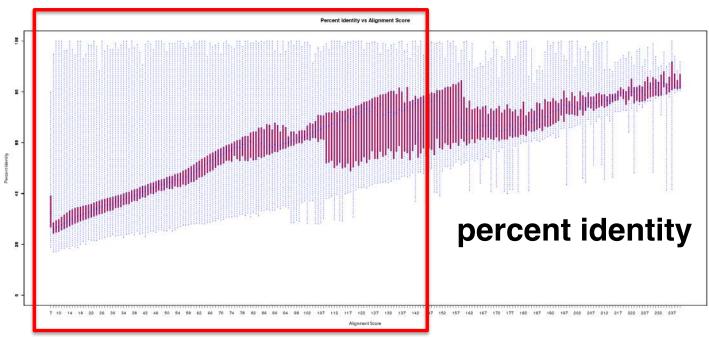








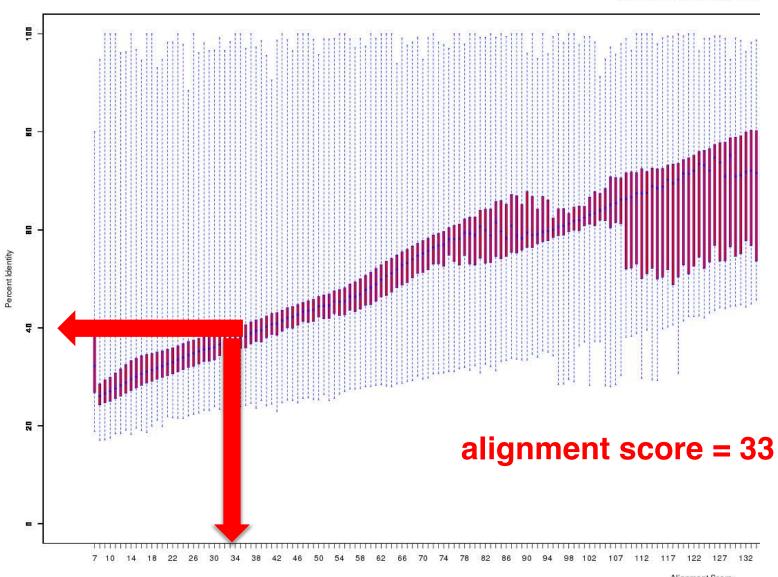




Alignment score for initial SSN: ~40% identity

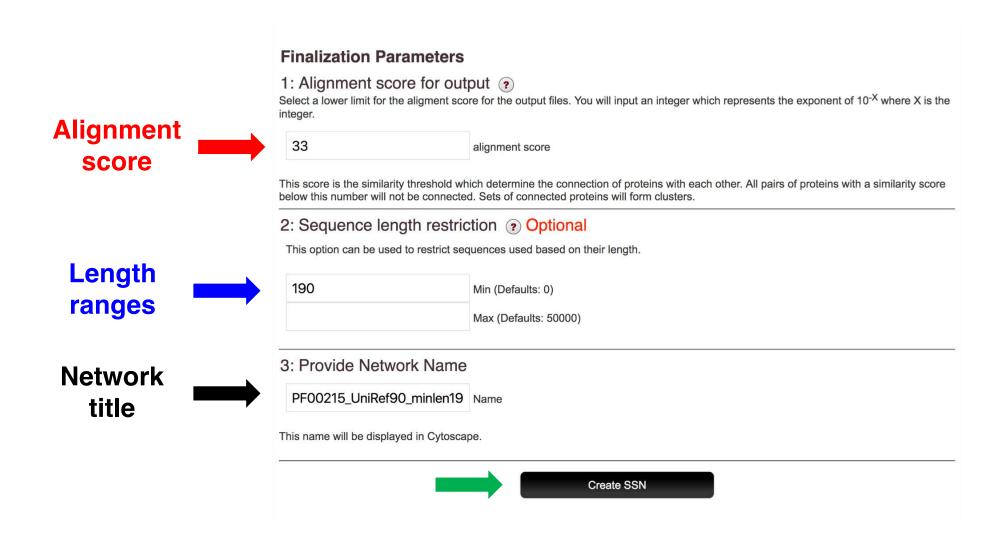


Percent Identity vs Alignment



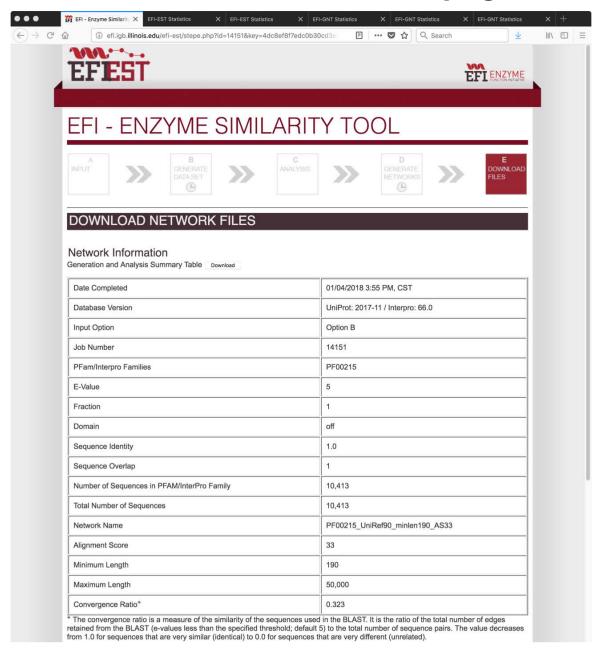
THE Analysis step: input alignment score and length ranges





Download Network File page





Network statistics and downloads



Full network

	# Nodes	# Edges	File Size (MB)
Download ZIP	10,045	8,264,486	2,201 MB

Representative Node Networks ?

In representative node (RepNode) networks, each node in the network represents a collection of proteins grouped according to percent identity. For example, for a 75% identity RepNode network, all connected sequences that share 75% or more identity are grouped into a single node (meta node).

The cluster organization is not changed, and the clustering of sequences remains identical to the full network.

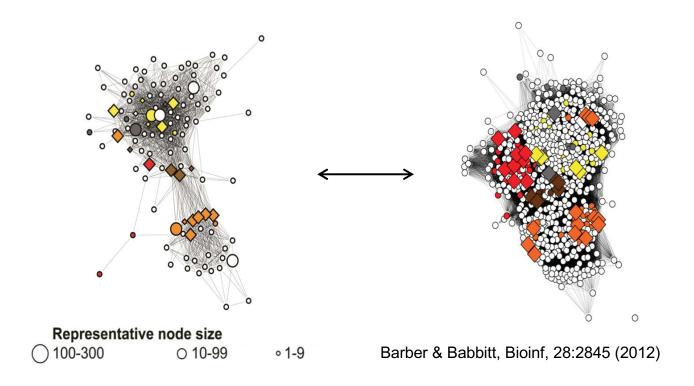
Sequences are collapsed together to reduce the overall number of nodes making for less complicated networks easier to load in Cytoscape.

	% ID	# Nodes	# Edges	File Size (MB)
Download Download ZIP	100	10,045	8,264,486	2,214 MB
Download Download ZIP	95	10,004	8,187,043	2,194 MB
Download Download ZIP	90	9,888	7,993,153	2,143 MB
Download Download ZIP	85	8,987	6,479,231	1,743 MB
Download Download ZIP	80	8,073	5,120,215	1,383 MB
Download Download ZIP	75	7,265	4,105,975	1,114 MB
Download Download ZIP	70	6,608	3,320,182	905 MB
Download Download ZIP	65	5,943	2,596,128	712 MB
Download Download ZIP	60	5,231	1,877,192	521 MB
Download Download ZIP	55	4,459	1,234,862	349 MB
Download Download ZIP	50	3,694	731,963	214 MB
Download Download ZIP	45	2,951	384,690	120 MB
Download Download ZIP	40	2,269	185,689	66 MB
<u></u>				

Rep node networks

Representative node (rep node) networks





Visualization algorithms (Cytoscape) limit the number of edges (and, therefore, indirectly the number of nodes) that can be displayed: <~500K with 8GB RAM, <~1M with 16GB RAM; <~5M with 64GB RAM; <~10M with 128GB RAM

Solution: "bin" similar sequences into the same "metanode", e.g., 80%, 90%, or 100% sequence identity rep node networks

I don't get a commission! But

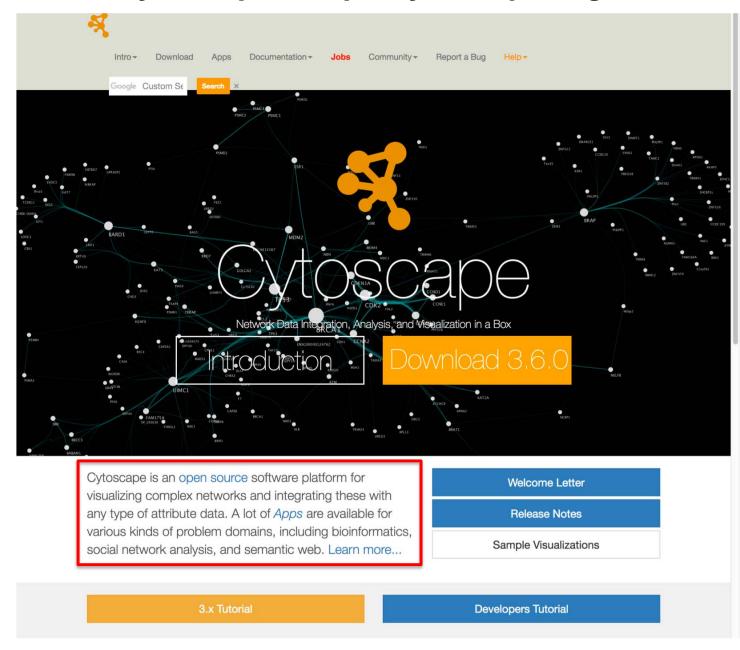
1

	Memory	Processor	RAM	Price
MacBook 12"	512GB	1.3GHz	8GB	\$1,599
MacBook 12"	512GB	1.3GHz	16GB	\$1,799
MacBook Pro 13"	512GB	3.1GHz	8GB	\$1,999
MacBook Pro 13"	512GB	3.1GHz	16GB	\$2,199
iMac 21"	1TB	3.4GHz	8GB	\$1,499
iMac 21"	1TB	3.4GHz	16GB	\$1,699
iMac 21"	1TB	3.4GHz	32GB	\$2,099
iMac 27"	2TB	3.8 GHz	8GB	\$2,299
iMac 27"	2TB	3.8 GHz	16GB	\$2,499
iMac 27"	2TB	3.8 GHz	32GB	\$2,899
iMac 27"	2TB	3.8 GHz	64GB	\$3,699
iMac Pro	1TB	3.2GHz	32GB	\$4,999
iMac Pro	1TB	3.2GHz	64GB	\$5,799
iMac Pro	1TB	3.2GHz	128GB	\$7,399



Cytoscape: http://cytoscape.org/

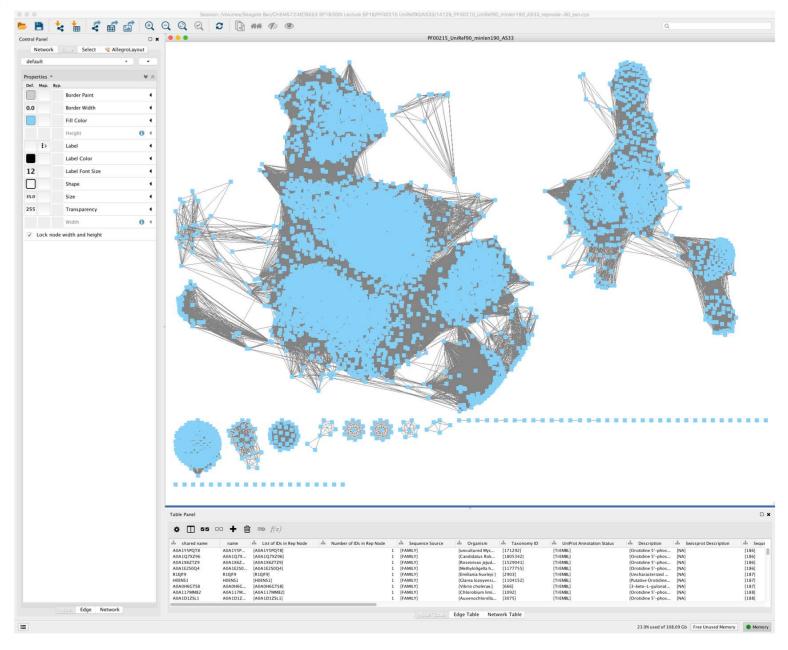






Cytoscape 3.3.0

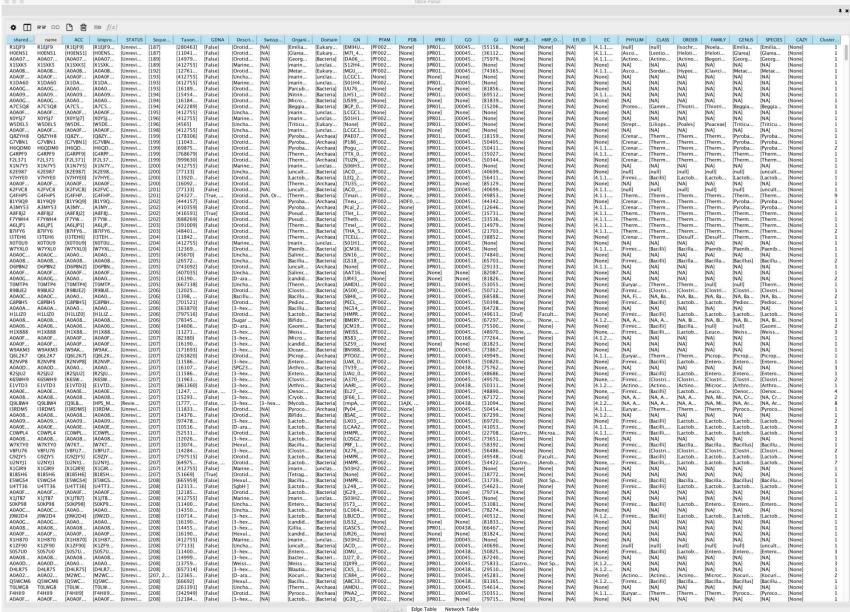






Data panel: node attributes







Node attributes: part 1



shared name – UniProt accession ID

name – UniProt accession ID

Organism – Organism name

Taxonomy ID – NCBI Taxnomy ID

UniProt Accession Status – TrEMBL or SwissProt

Description – UniProt description

Swissprot Description – Swiss Prot description

Sequence Length – number of residues

Gene Name - Gene name

NCBI IDs – RefSeq IDs, UniProt/Swiss-Prot IDs, GenBank IDs, PDB IDs,

"retired" GI numbers

Superkingdom – Phylogenetic superkingdom of the organism (from UniProt)

Kingdom – Phylogenetic kingdom of the organism

Phylum – Phylogenetic phylum of the organism

Class – Phylogenetic class of the organism

Order – Phylogenetic order of the organism

Family – Phylogenetic family of the organism

Genus – Phylogenetic genus of the organism

Species – Phylogenetic species of the organism



Node attributes: part 2



EC – EC (Enzyme Commisssion) number

PFAM – Pfam family number

IPRO – InterPro domain/family number

BRENDA ID — BRENDA ID: https://www.brenda-enzymes.org/

CAZY – Carbohydrate-Active enZYmes: http://www.cazy.org/

GO Term – Gene ontology (from UniProt)

KEGG ID – KEGG ID: http://www.genome.jp/kegg/

PATRIC ID – Genome annotation and assembly: https://www.patricbrc.org/

STRING ID – Protein-protein interactions: https://string-db.org/

HMP Body Site – Human Microbiome Project organisms

HMP Oxygen – Human Microbiome Project organisms

P01 gDNA - availability of gDNA in P01GM118303 Protein Core

UniRef90 Cluster IDs – UniProt IDs for seed sequences in UniRef90 clusters

UniRef90 Cluster Size – number of sequences/IDs in UniRef90 clusters

Sequence Source – Options B/C/D: USER or FAMILY

Query IDs – Option C/D: UniProt IDs in FASTA headers or input list

Other IDs — Option C: FASTA headers with unrecognized accession IDs

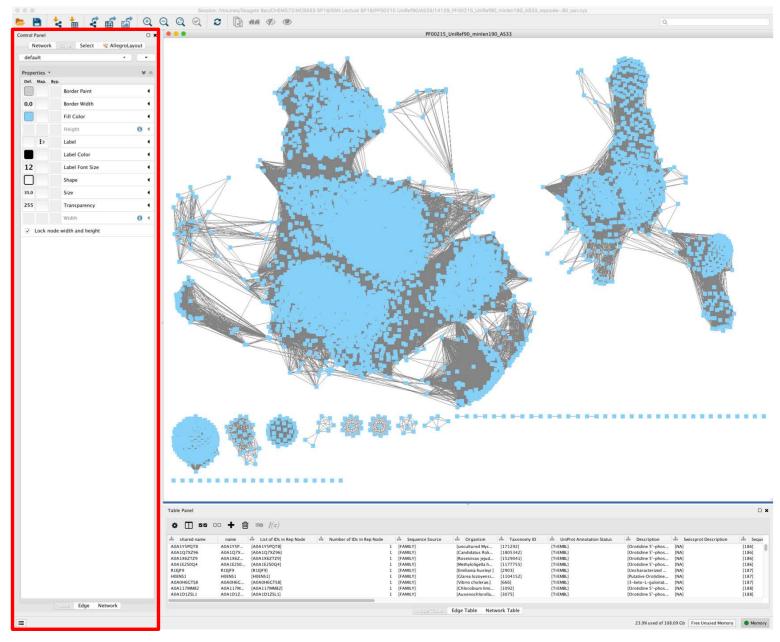
List of IDs in Repnode – UniProt IDs in repnode metanodes

Number of IDs in Repnode – number of sequences in repnode metanodes



Cytoscape 3.3.0: control panels







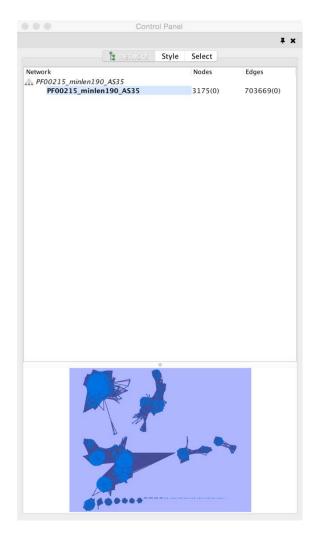
Cytoscape control panels: manipulating networks

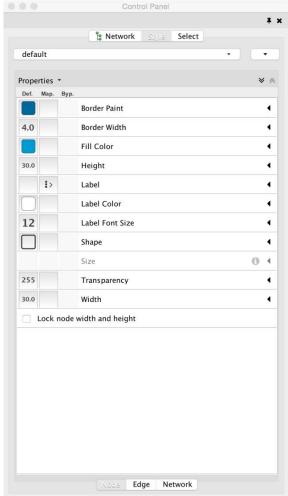


Network

Style

Select







Some "bioinorganic" enzyme superfamilies



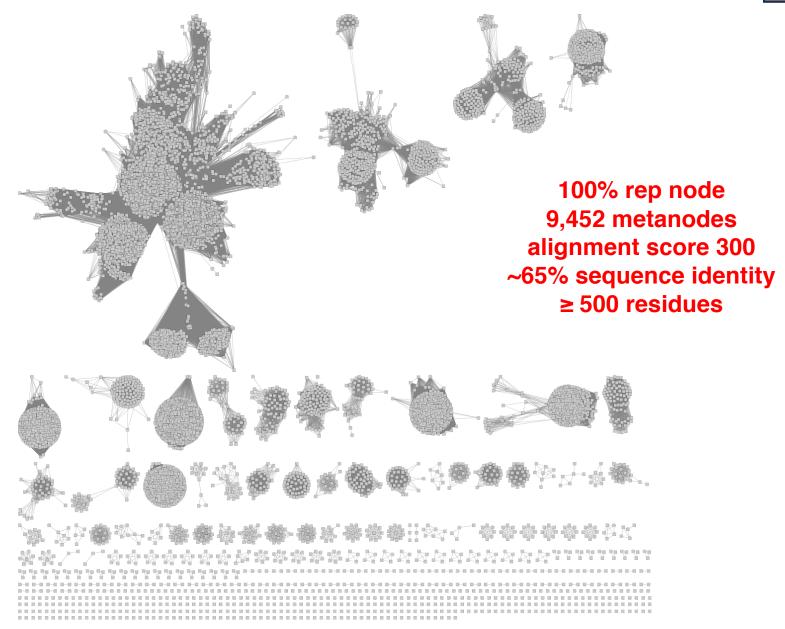
Name	Pfam Family	InterPro Family	Family Size
Radical SAM	PF04055	IPR007197	360,840
P450	PF00067	IPR001128	211,400
Ferritin-like		IPR009078	168,246
Molybdopterin oxidoreductase	PF00384	IPR006656	94,758
Oxoglutarate Fe-independent Dioxygenase	PF13640/PF03171	IPR005123	84,800
Isopenicillin N synthase		IPR027443	50,698
TauD/TfdA-like	PF02668	IPR003819	35,996
Phytanoyl-CoA dioxygenase	PF05721	IPR008775	34,121
Ribonucleotide reductase small subunit family	PF00268	IPR000358	23,764
Oxoglutarate AlkB-like	PF13532	IPR027450	18,414
Glycyl radical enzyme	PF02901	IPR004184	14,229
FldB/FldC dehydratase	PF06050	IPR010327	9,079
20G-Fe dioxygenase	PF10014	IPR018724	2,765
20GFeDO oxygenase domain	PF12851	IPR024779	930

Some "bioinorganic" enzyme superfamilies



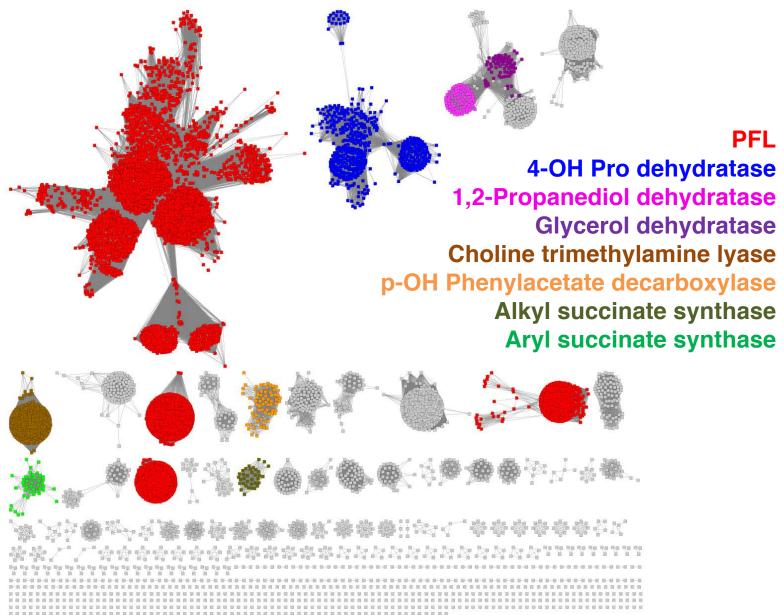
Name	Pfam Family	InterPro Family	Family Size
Radical SAM	PF04055	IPR007197	360,840
P450	PF00067	IPR001128	211,400
Ferritin-like		IPR009078	168,246
Molybdopterin oxidoreductase	PF00384	IPR006656	94,758
Oxoglutarate Fe-independent Dioxygenase	PF13640/PF03171	IPR005123	84,800
Isopenicillin N synthase		IPR027443	50,698
TauD/TfdA-like	PF02668	IPR003819	35,996
Phytanoyl-CoA dioxygenase	PF05721	IPR008775	34,121
Ribonucleotide reductase small subunit family	PF00268	IPR000358	23,764
Oxoglutarate AlkB-like	PF13532	IPR027450	18,414
Glycyl radical enzyme	PF02901	IPR004184	14,229
FIdB/FIdC dehydratase	PF06050	IPR010327	9,079
20G-Fe dioxygenase	PF10014	IPR018724	2,765
20GFeDO oxygenase domain	PF12851	IPR024779	930

Glycyl radical enzyme (GRE) superfamily: 14,229 sequence



Clusters with known functions





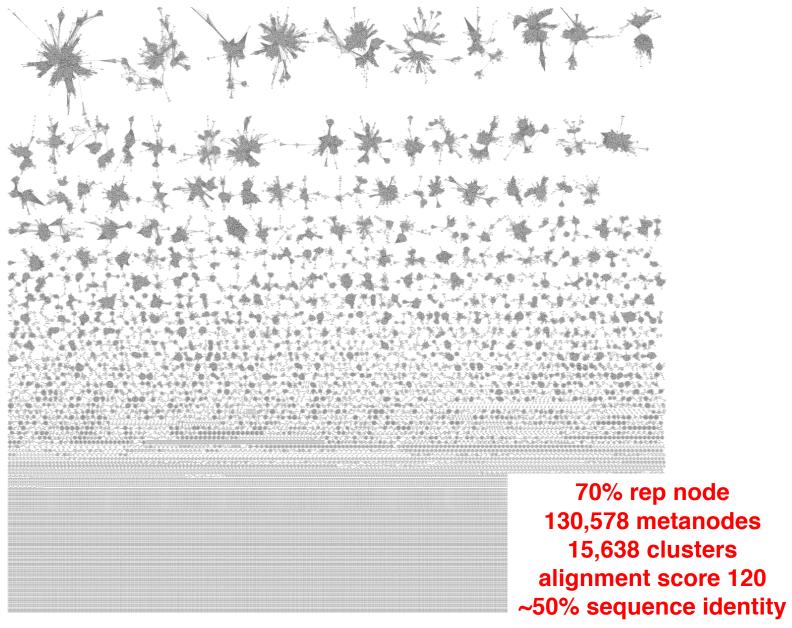
Some "bioinorganic" enzyme superfamilies



Name	Pfam Family	InterPro Family	Family Size
Radical SAM	PF04055	IPR007197	360,840
P450	PF00067	IPR001128	211,400
Ferritin-like		IPR009078	168,246
Molybdopterin oxidoreductase	PF00384	IPR006656	94,758
Oxoglutarate Fe-independent Dioxygenase	PF13640/PF03171	IPR005123	84,800
Isopenicillin N synthase		IPR027443	50,698
TauD/TfdA-like	PF02668	IPR003819	35,996
Phytanoyl-CoA dioxygenase	PF05721	IPR008775	34,121
Ribonucleotide reductase small subunit family	PF00268	IPR000358	23,764
Oxoglutarate AlkB-like	PF13532	IPR027450	18,414
Glycyl radical enzyme	PF02901	IPR004184	14,229
FIdB/FIdC dehydratase	PF06050	IPR010327	9,079
20G-Fe dioxygenase	PF10014	IPR018724	2,765
20GFeDO oxygenase domain	PF12851	IPR024779	930

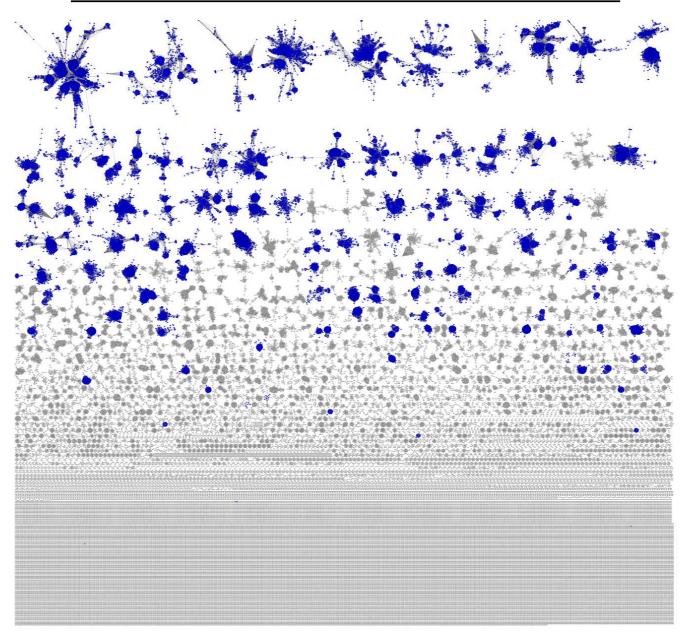
Radical SAM superfamily: 360,891 sequences





Clusters with SwissProt annotations







"These enzymes catalyze a dizzying array of novel transformations, and their discovery has led to a revival of mechanistic enzymology. Indeed, many beginning professors who are interested in chemical biology (i.e. natural product biosynthesis) or mechanistic enzymology find themselves working in the field of radical SAM enzymology. The ability to mine the large space of unknown function within the superfamily would be crucial to uncovering these new reactions and seeding new projects for young professors."

Today



- 1. Target selection using sequence similarity networks: EFI-EST web tool
- 2. Pathway discovery using genome neighborhood networks: EFI-GNT web tool
- 3. Example: using SSNs and GNNs to discover catabolic pathways for D-apiose
- 4. Target prioritization using chemically guided functional profiling: CGFP-ShortBRED beta tool

<u>Use genome context to inform annotation</u> <u>for bacterial and fungal enzymes</u>



High sequence similarity Enzyme function (≥ 70%)

Low sequence similarity ??????? (≤ 70%)

Genome context in Genome Neighborhood Networks (GNNs)

Input: SSN families from EFI-EST



Query families





GNN: bacterial and fungal proteins in gene clusters

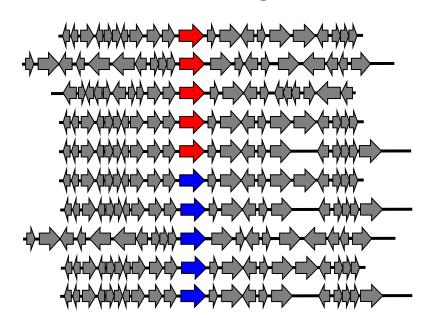


Query families

Genome neighbors







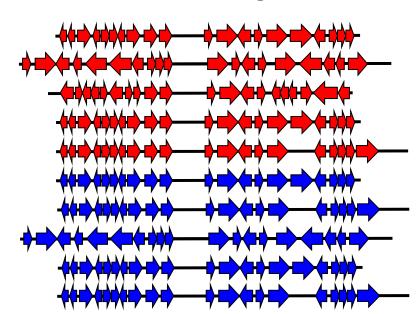
GNN: collect neighbors (±10)

Query families

Genome neighbors







First GNN: Pfam families found by SSN clusters



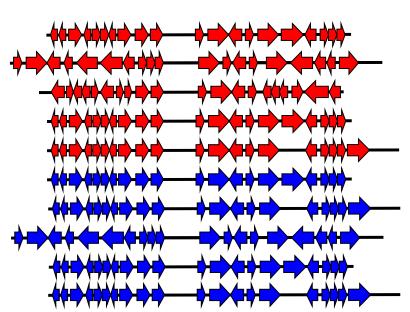
Query families

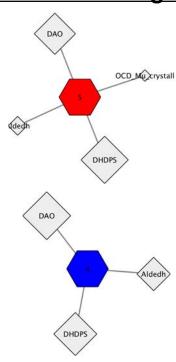
Genome neighbors

SSN cluster neighbors



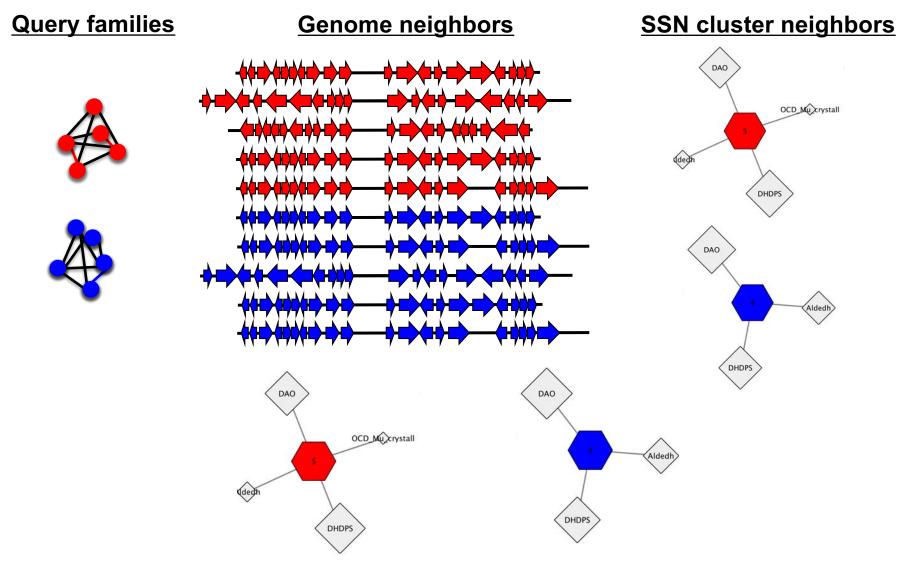






First GNN: "read" pathways from GNN clusters



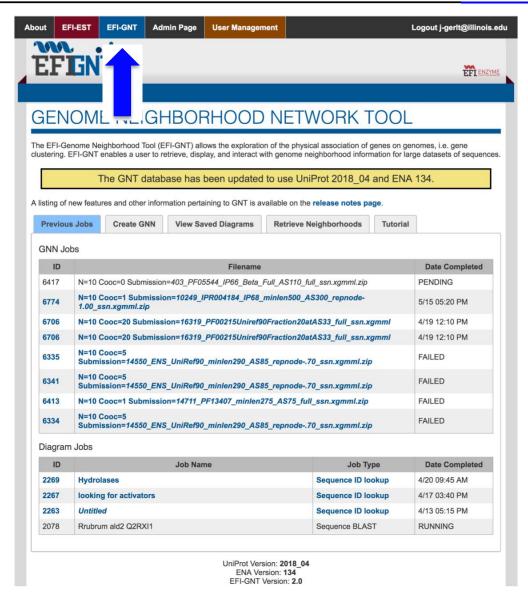


Deduce pathways from identities of neighbor Pfam families



EFI-EST and EFI-GNT Web Tools: EFI-GNT





https://efi.igb.illinois.edu/efi-gnt/



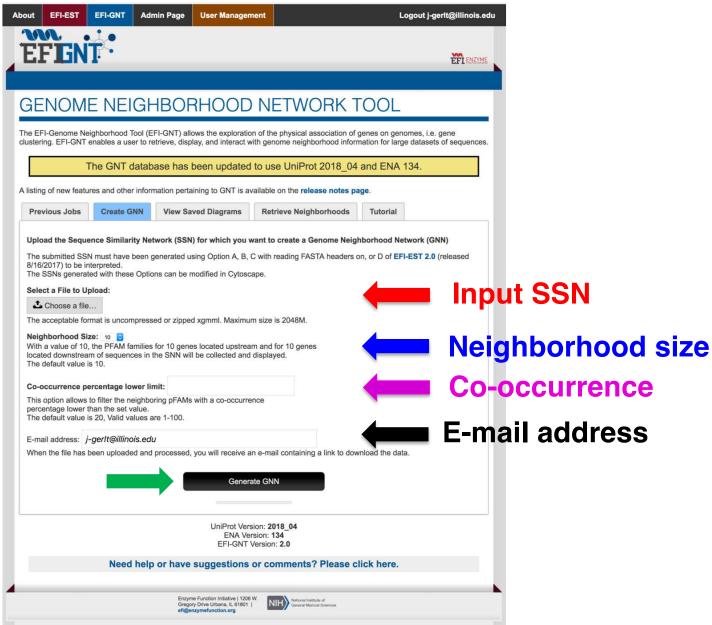
EFI-GNT uses genome sequences from ENA (EMBL/EBI)



EMBL-EBI	Services Research Training About of	
European Nucleotide Archive	Search xamples: BN000065, histone Advanced Sequence	
Search & Browse Submit & Update Software About ENA Support		
European Nucleotide Archive	Popular	
The European Nucleotide Archive (ENA) provides a comprehensive record of the world's nucleotide sequencing information, covering raw sequencing data, sequence assembly information and functional annotation. More about ENA	Submit and update Sequence submissions Genome assembly submissions	
ccess to ENA data is provided though the browser, through search tools, large scale file download and through the API.	Submitting environmental sequences Citing ENA data Rest URLs for data retrieval	
Text Search	Rest URLs to search ENA	
Examples: BN000065, histone	Latest ENA news	
Search Advanced search	19 Mar 2018: ENA Release 135 Release 135 of ENA's assembled/annotated sequences is now available	
Sequence Search	now available	
Enter or paste a nucleotide sequence or accession number		
Search Advanced search		
ENA is part of the ELIXIR infrastructure		

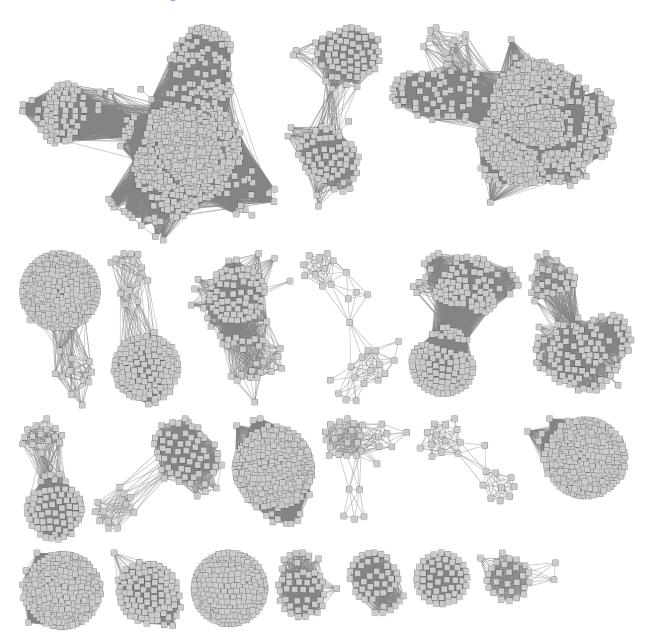
EFI-GNT: create GNN from SSN





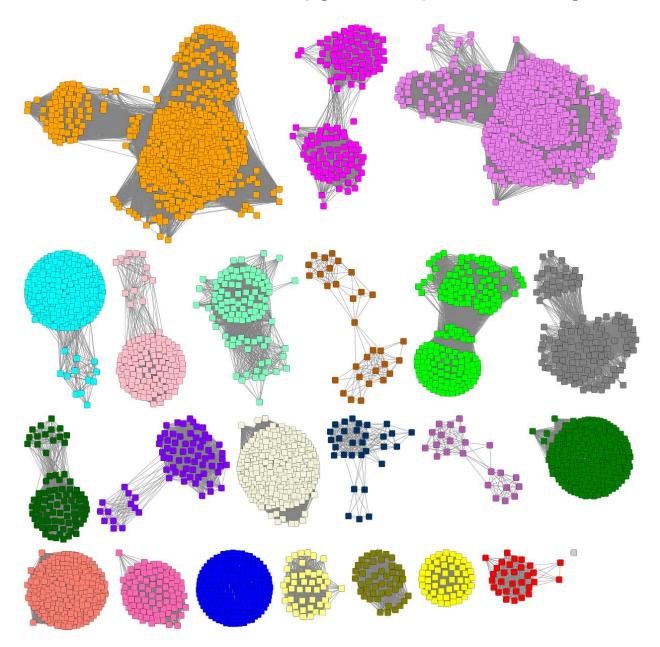
Input SSN: from EFI-EST





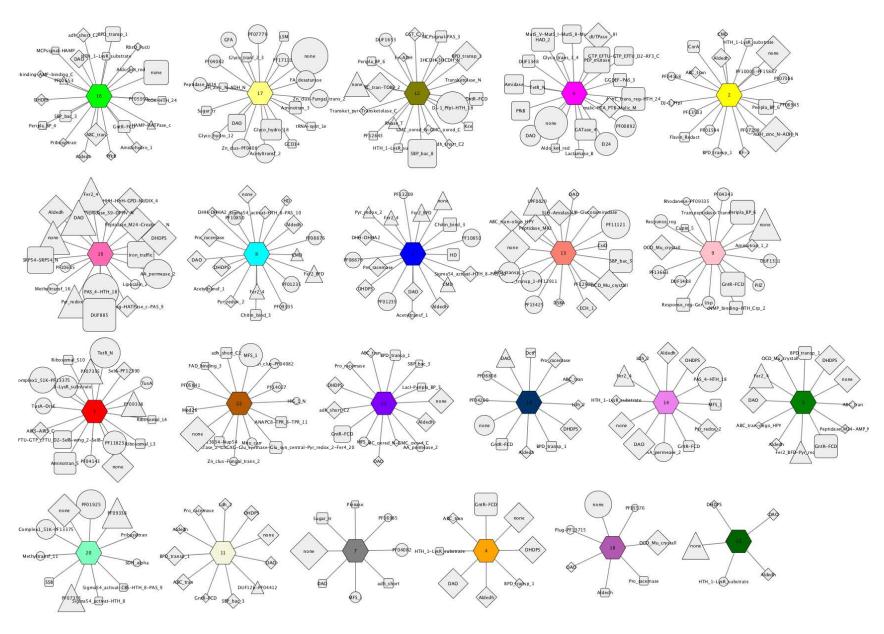
Colored SSN: families (queries) with unique colors





First GNN: Pfam families found by SSN clusters

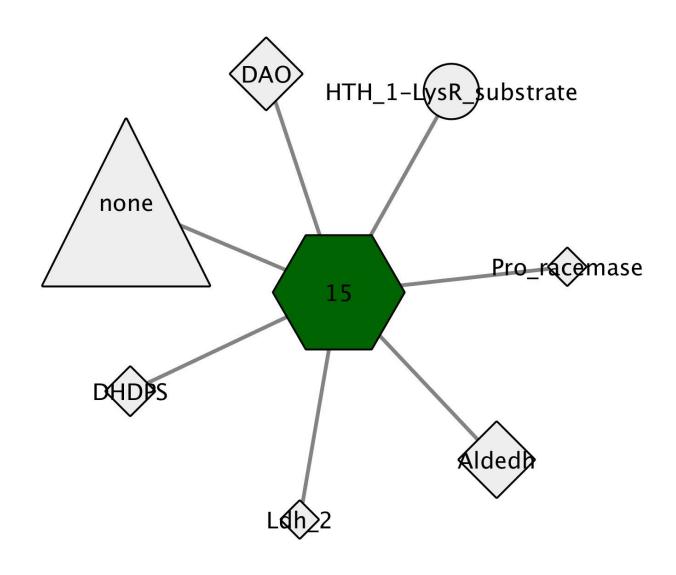




SSN cluster hub-nodes and Pfam family spoke-nodes

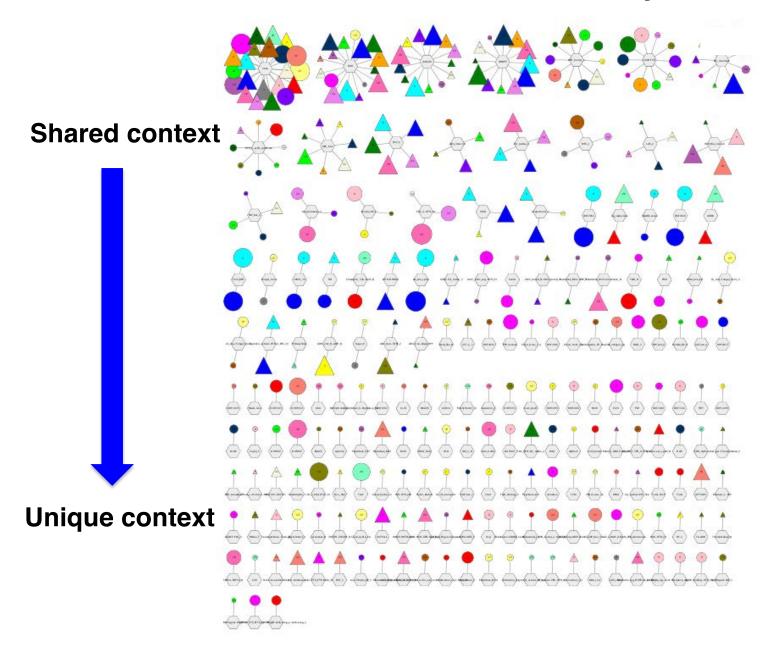


Functionally linked enzymes in pathways



Second GNN: SSN clusters found by Pfam families

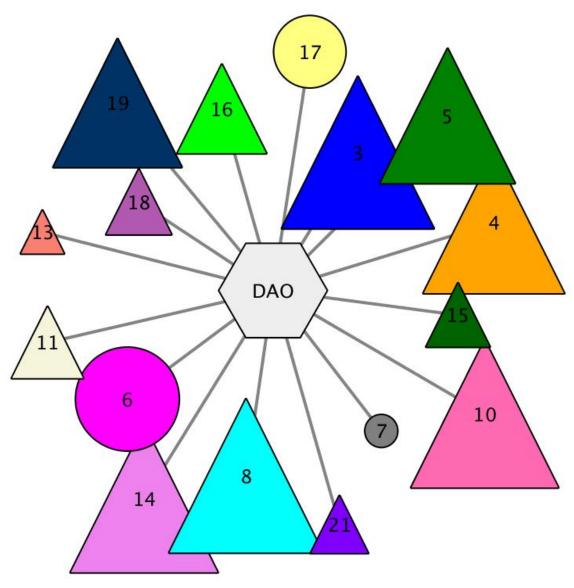




Pfam family hub-nodes and SSN cluster spoke-nodes

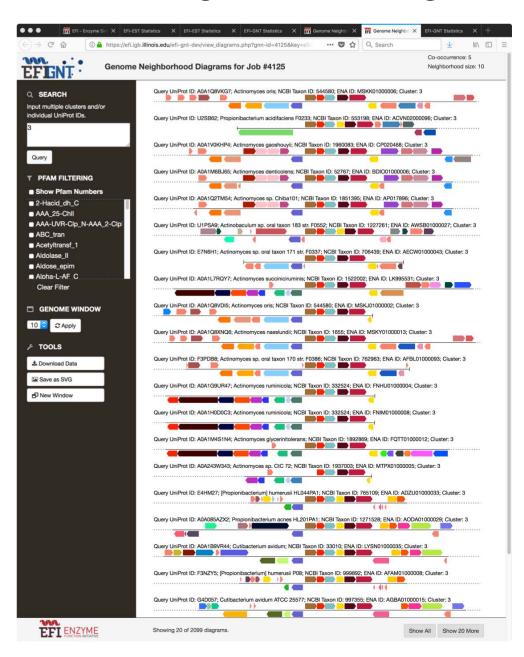


SSN clusters with shared functions



Genome Neighborhood Diagrams





Today



- 1. Target selection using sequence similarity networks: EFI-EST web tool
- 2. Pathway discovery using genome neighborhood networks: EFI-GNT web tool
- 3. Example: using SSNs and GNNs to discover catabolic pathways for D-apiose
- 4. Target prioritization using chemically guided functional profiling: CGFP-ShortBRED beta tool

Functional assignment of multiple catabolic pathways for D-apiose

Nature Chemical Biology, in press DOI: 10.1038/s41589-018-0067-7 Available on line, June 4

A general strategy for discovering pathways





Article

pubs.acs.org/biochemistry

Experimental Strategies for Functional Annotation and Metabolism Discovery: Targeted Screening of Solute Binding Proteins and Unbiased Panning of Metabolomes

Matthew W. Vetting,**,*,¶ Nawar Al-Obaidi,*,¶ Suwen Zhao,*,¶ Brian San Francisco,§,¶ Jungwook Kim,† Daniel J. Wichelecki,§,∥ Jason T. Bouvier,§,∥ Jose O. Solbiati,§ Hoan Vu,# Xinshuai Zhang,§ Dmitry A. Rodionov,O,♥ James D. Love,† Brandan S. Hillerich,† Ronald D. Seidel,† Ronald J. Quinn,*,# Andrei L. Osterman,*,○ John E. Cronan,*,∥,♠ Matthew P. Jacobson,*,‡ John A. Gerlt,*,§,∥,⊥ and Steven C. Almo*,†

[†]Department of Biochemistry, Albert Einstein College of Medicine, Bronx, New York 10461, United States

[‡]Department of Pharmaceutical Chemistry, University of California, San Francisco, California 94158, United States

[§]Institute for Genomic Biology, [∥]Department of Biochemistry, [⊥]Department of Chemistry, [♠]Department of Microbiology, University of Illinois at Urbana—Champaign, Urbana, Illinois 61801, United States

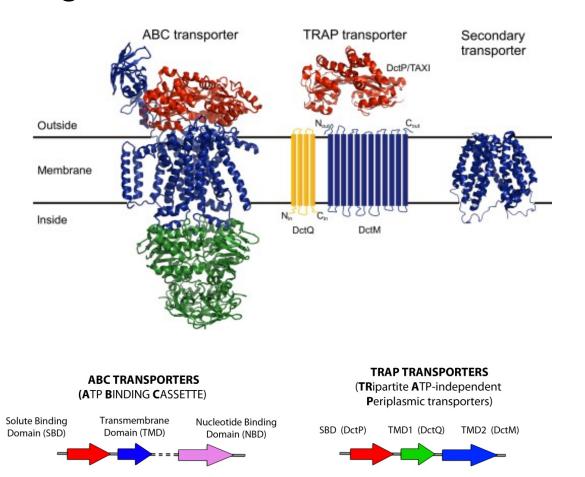
^{*}Eskitis Institute for Drug Discovery, Griffith University, Brisbane, Queensland 4111, Australia

^OSanford-Burnham Medical Research Institute, La Jolla, California 92037, United States

VA.A.Kharkevich Institute for Information Transmission Problems, Russian Academy of Sciences, Moscow 127994, Russia

<u>Transporter Solute Binding Proteins (SBPs):</u> genes colocated with those for catabolic pathways





3B50 TRAP SBP Sialic acid Binding Protein

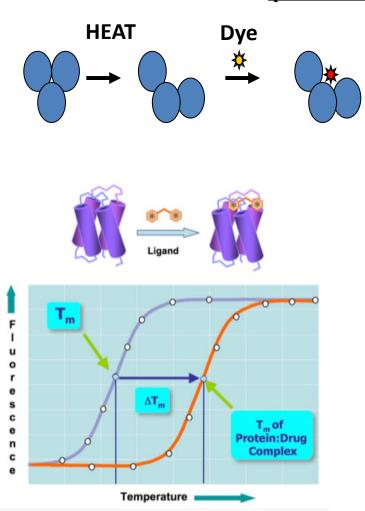


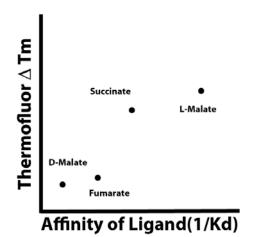
1ANF ABC SBP Maltose Binding Protein

Mulligan, C et al. (2011). "Tripartite ATP-independent periplasmic (TRAP) transporters in bacteria and archaea." <u>FEMS Microbiol Rev</u> 35(1): 68-86.

<u>Ligand screening by differential scanning fluorimetry</u> (DSF or ThermoFluor)







Order of ligand stabilization correlated with affinity of ligand but not absolute

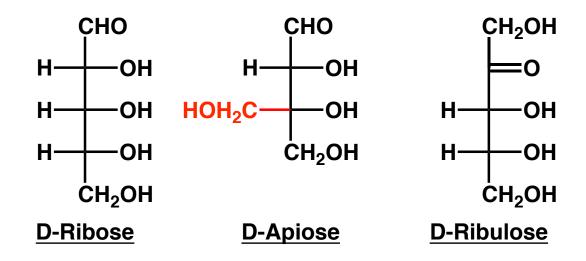
Applied Biosystems 7900HT Fast Real-Time PCR System

1 mM Ligand 10 μm Protein (100X) 100 mM Hepes pH 7.5, 150 mM NaCl Sypro Orange (Indicator Dye)

Library Hits for ABC Transport SBPs (PF13407)

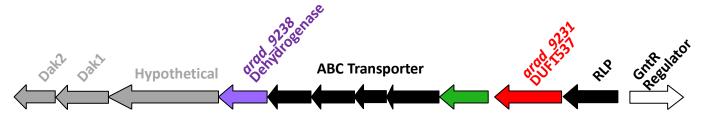


SBP (UniProt)	Organism	D-Ribose (ΔTm °C)	D-Apiose (∆Tm °C)	D-Ribulose (ΔTm °C)
Q2JZQ5	Rhizobium etli CFN 42	12.5	6.8	3.3
B1G898	Burkholderia graminis C4D1M	15.9	8.8	4.7
A6VKQ8	Actinobacillus succinogenes 130Z	16.2	9.7	1.4

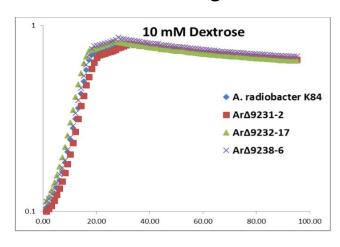


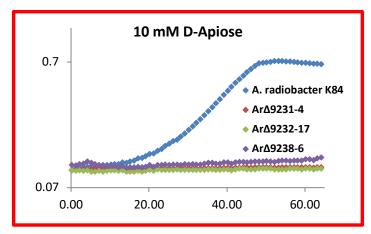
Knockouts in Agrobacterium radiobacter: D-apiose

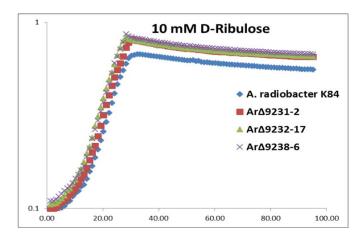


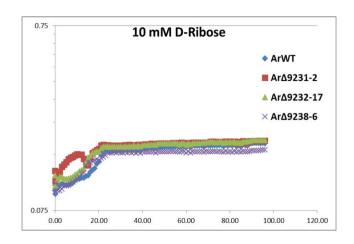


Agrobacterium radiobacter K84







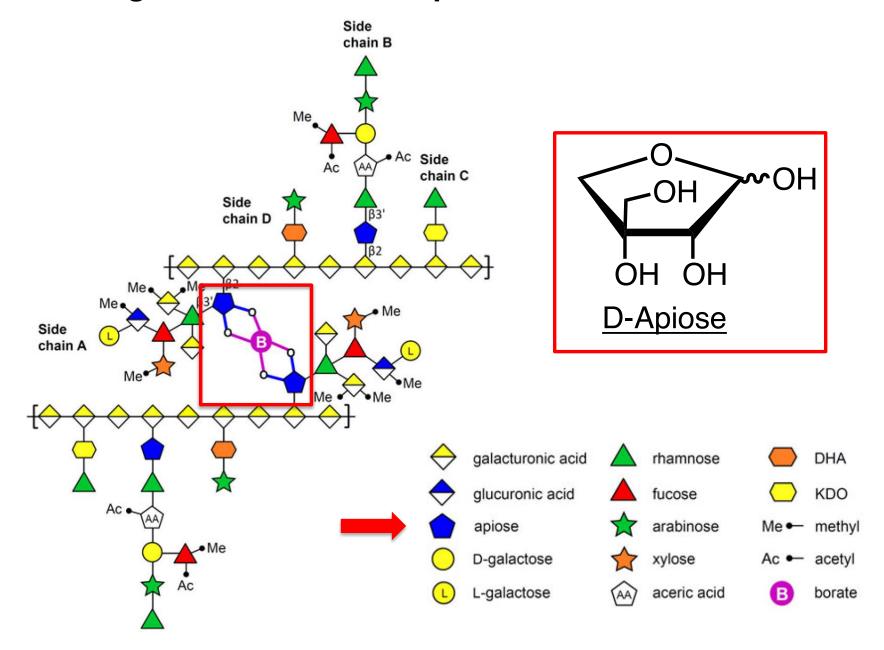


D-Apiose is physiological ligand



SBP (UniProt)	Organism	D-Ribose (ΔTm °C)	D-Apiose (ΔTm °C)	D-Ribulose (ΔTm °C)
Q2JZQ5	Rhizobium etli CFN 42	12.5	6.8	3.3
B1G898	Burkholderia graminis C4D1M	15.9	8.8	4.7
A6VKQ8	Actinobacillus succinogenes 130Z	16.2	9.7	1.4

Rhamnogalacturonan-II: D-apiose cross-links in cell walls



Glycoside hydrolases and lyases: release D-apiose



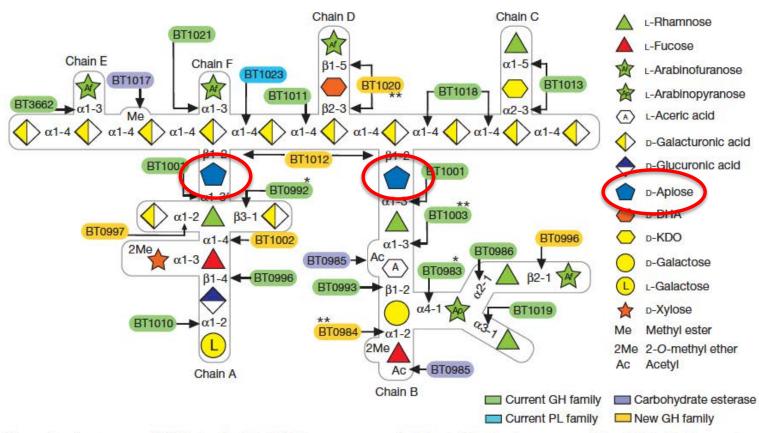


Figure 1 | Schematic of enzymes and PULs involved in RG-II degradation. Sugars shown using the Consortium for Functional Glycomics notation²⁵. Enzymes are appropriately colour-coded. Single asterisks denote a new activity for a glycoside hydrolase (GH) family;

double asterisks denote enzymes with activities that have not previously been observed. D-KDO, 3-deoxy-D-manno-octulosonic acid; D-DHA, 2-keto-3-deoxy-d-*lyxo*-heptulosaric acid; PL, polysaccharide lyase.

But, no catabolic pathway for D-apiose?



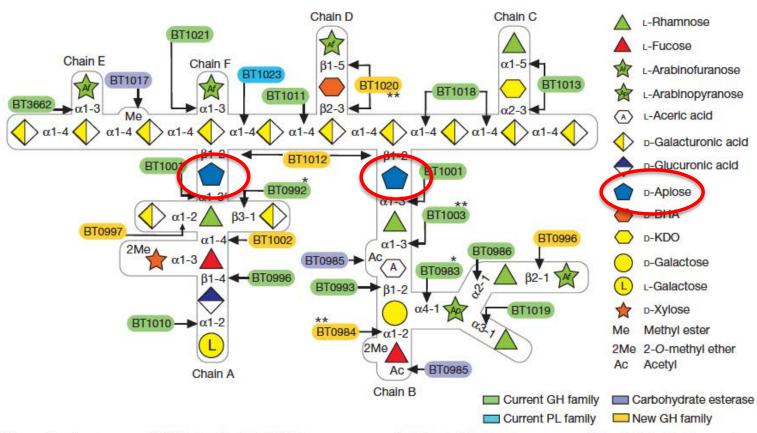
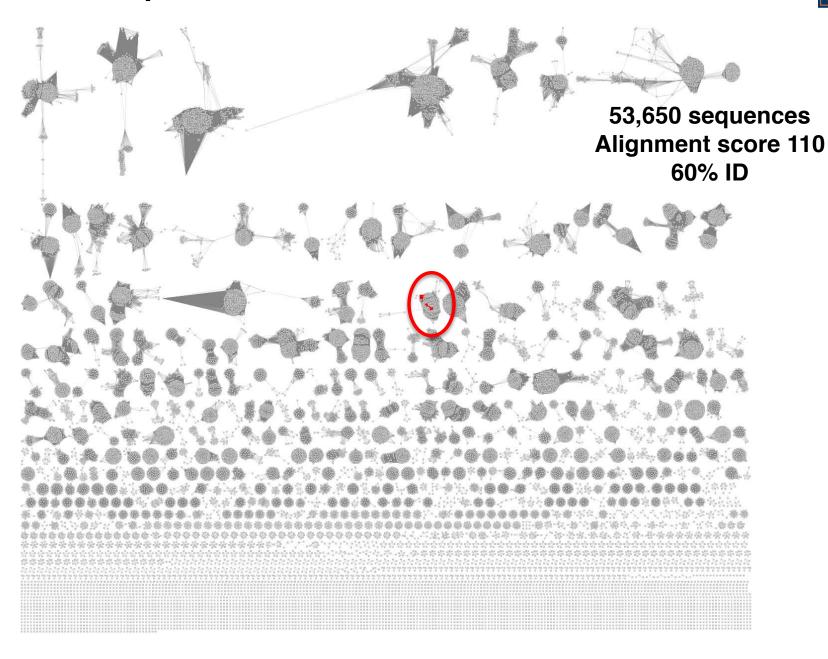


Figure 1 | Schematic of enzymes and PULs involved in RG-II degradation. Sugars shown using the Consortium for Functional Glycomics notation²⁵. Enzymes are appropriately colour-coded. Single asterisks denote a new activity for a glycoside hydrolase (GH) family;

double asterisks denote enzymes with activities that have not previously been observed. D-KDO, 3-deoxy-D-manno-octulosonic acid; D-DHA, 2-keto-3-deoxy-d-*lyxo*-heptulosaric acid; PL, polysaccharide lyase.

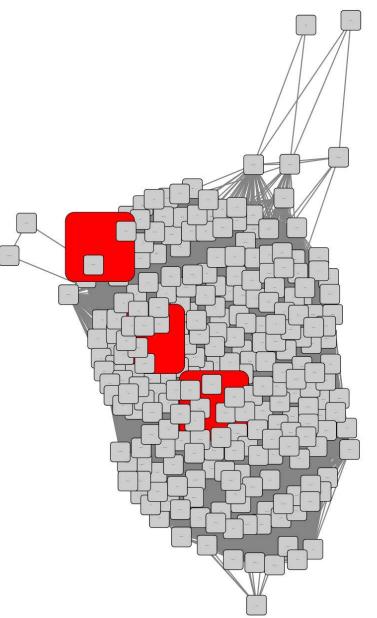
D-Apiose: three hits in one SBP cluster





D-apiose SBP cluster

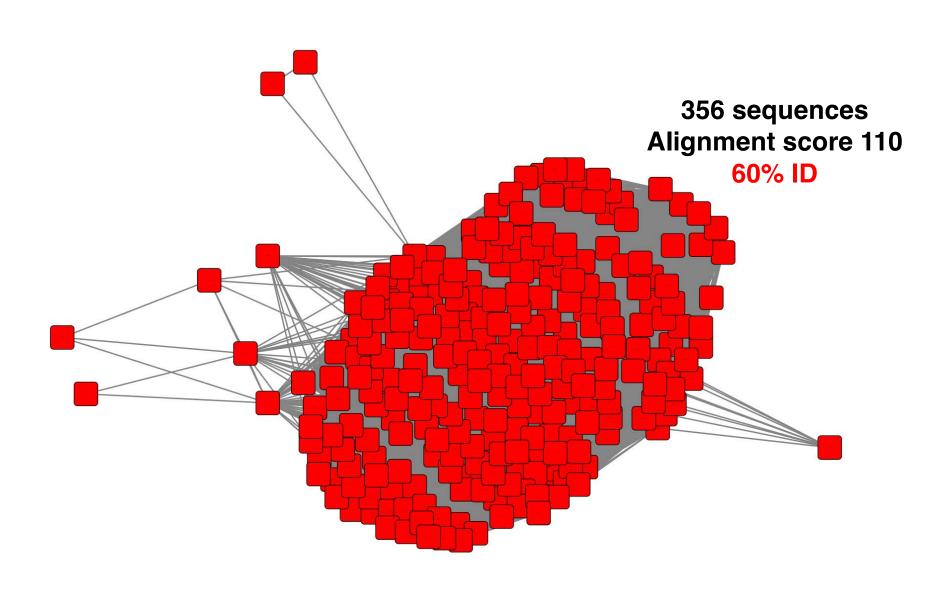




356 sequences
Alignment score 110
60% ID

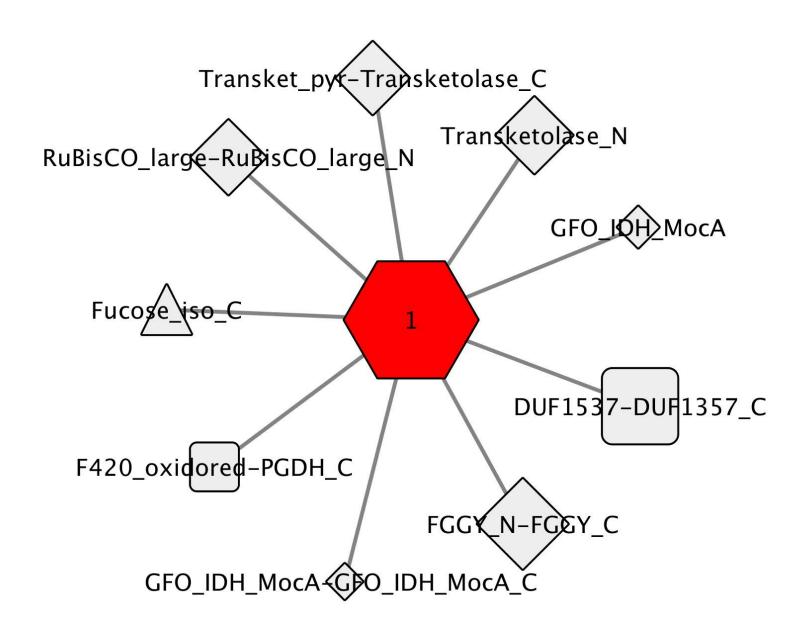
D-Apiose SBP hits: colored SSN





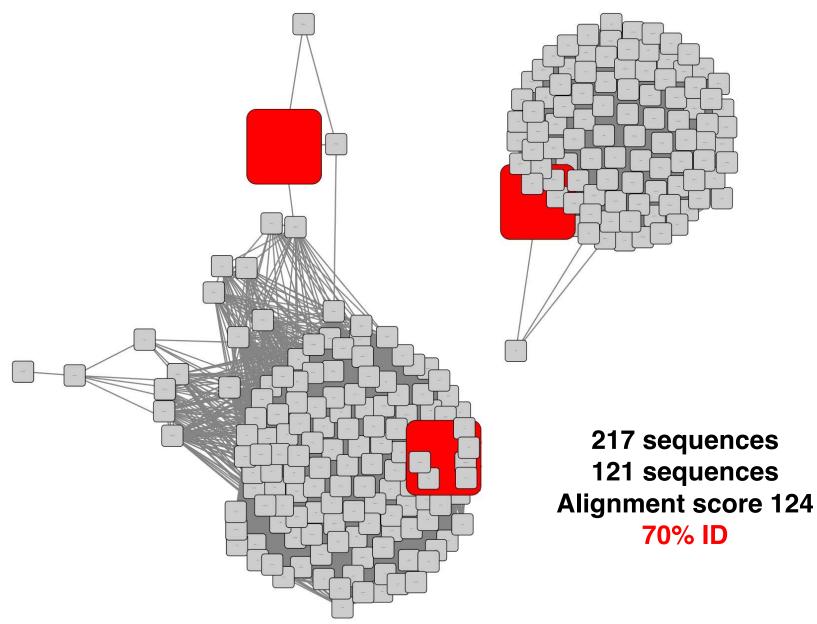
GNN: (too) many enzyme families, multiple pathways?





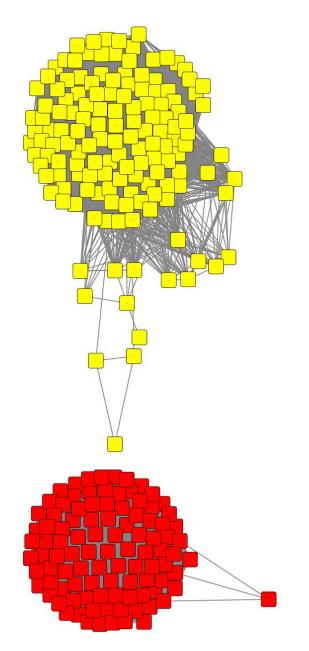
70% ID, alignment score 124: two clusters





70% ID, alignment score 124: colored SSN

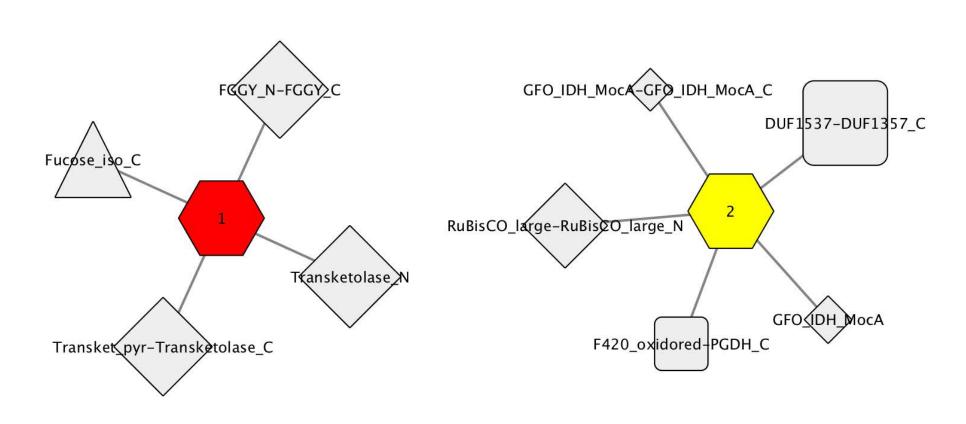




217 sequences 121 sequences Alignment score 124 70% ID

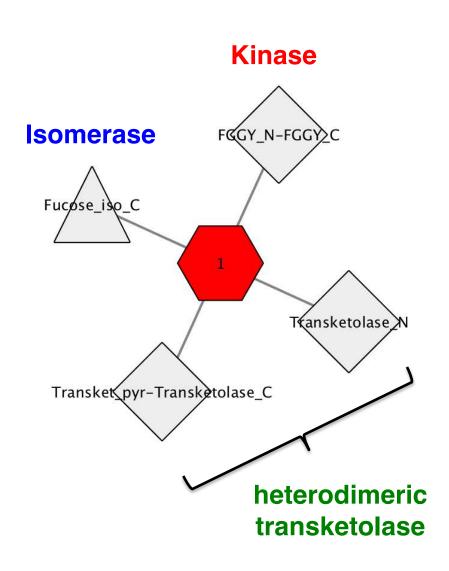
GNN: two pathways?





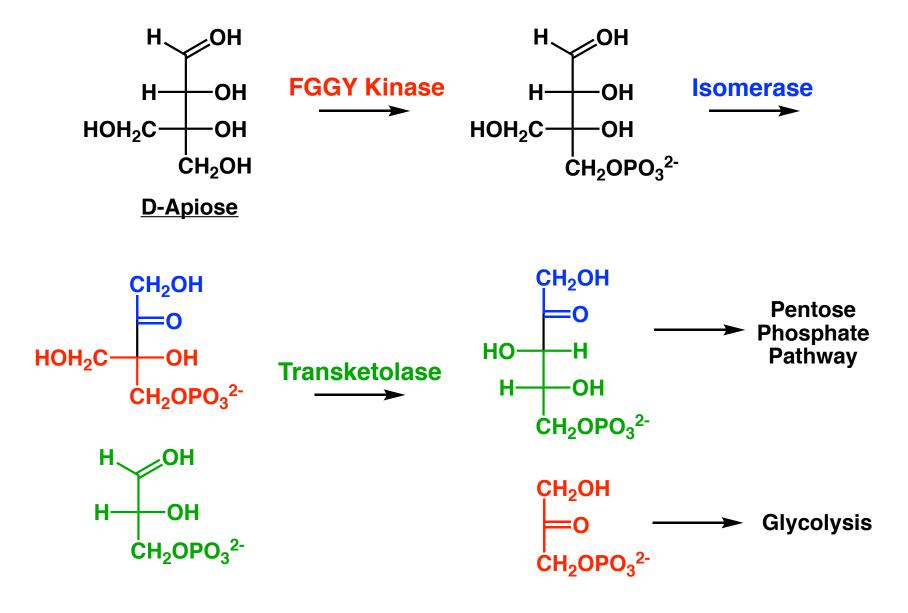
Red cluster: transketolase pathway





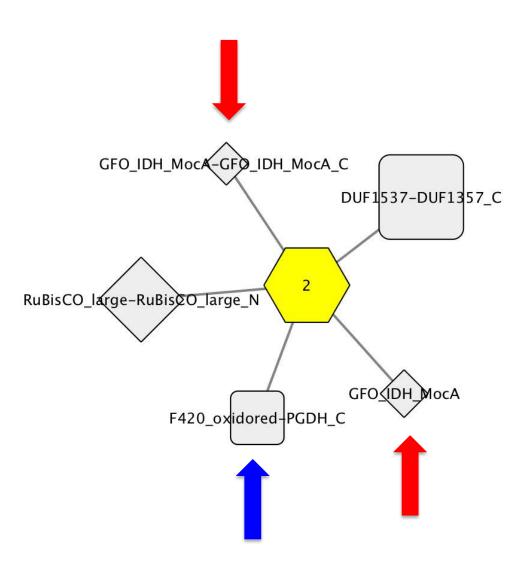
Transketolase pathway





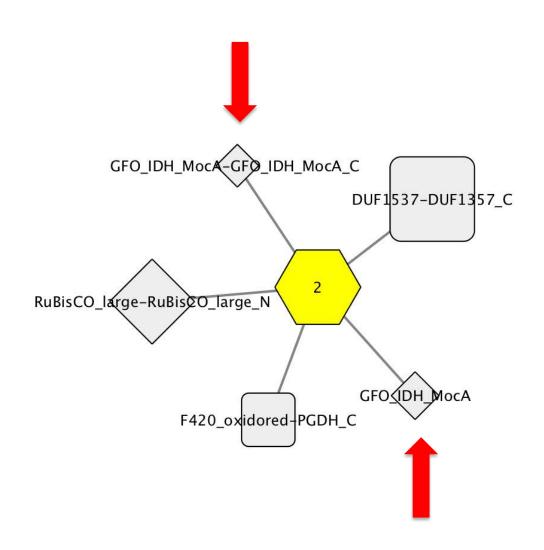
Yellow cluster: oxidative pathway(s) with two dehydrogenases





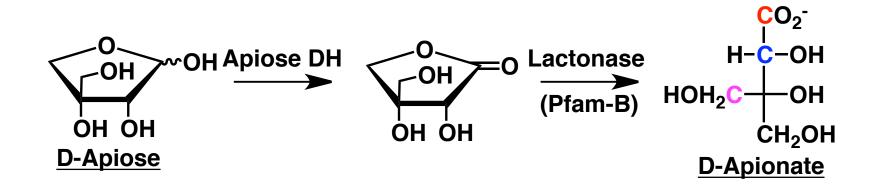
Oxidative pathway: D-apiose dehydrogenase





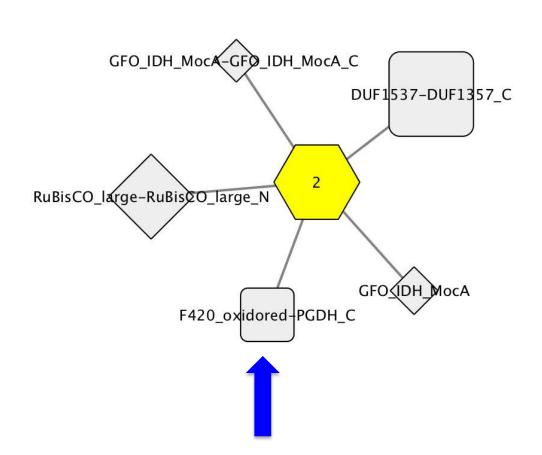
Oxidation of D-apionate





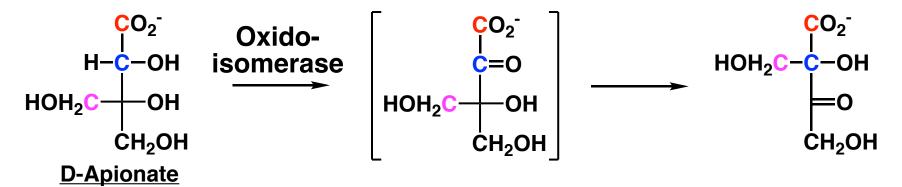
Oxidative pathway: oxidoisomerase





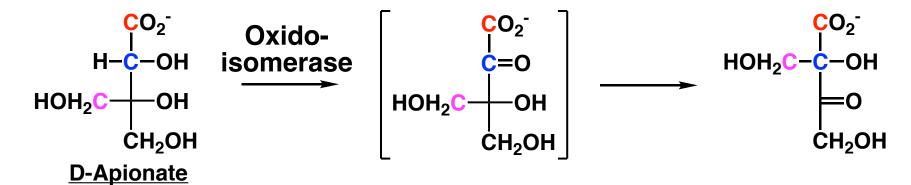
PF16896: novel oxidoisomerase





PF16896: novel oxidoisomerase



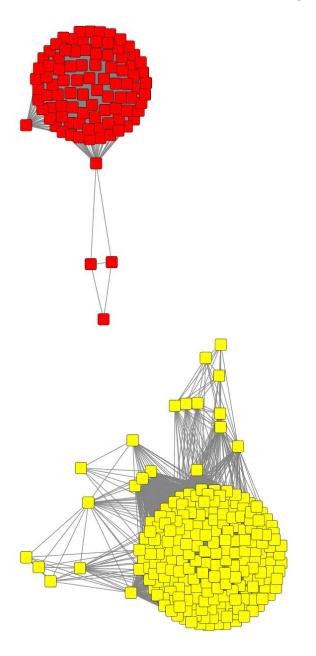


Concerted?

Or, stepwise?

Oxidoisomerase SSN: two major clusters



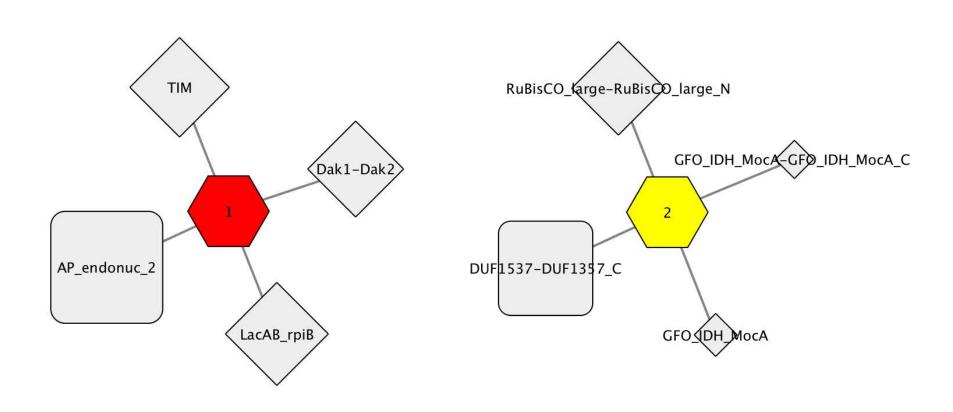


500 total sequences Alignment score 80

70% ID

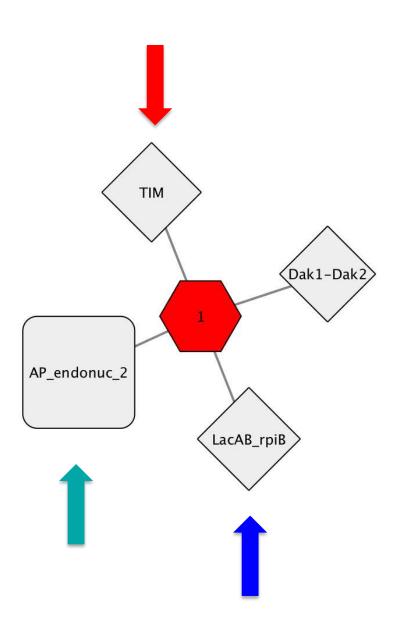
GNN: two (or more) oxidative pathways?





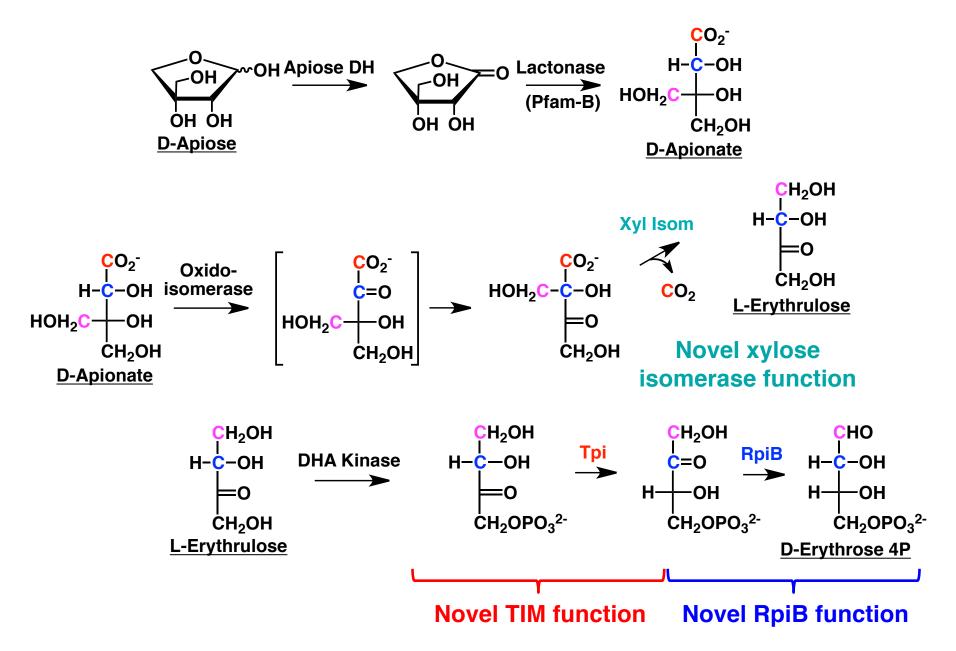
AP endonuclease, TIM, and RpiB families?





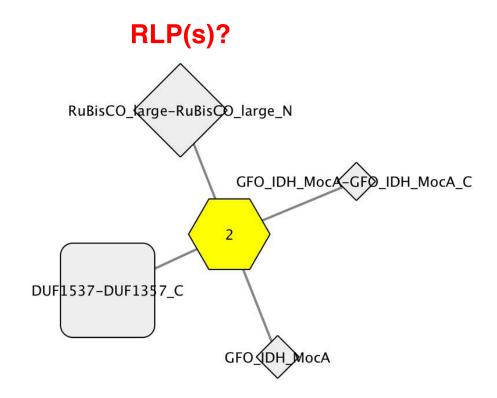
D-Apiose to D-erythrose 4P





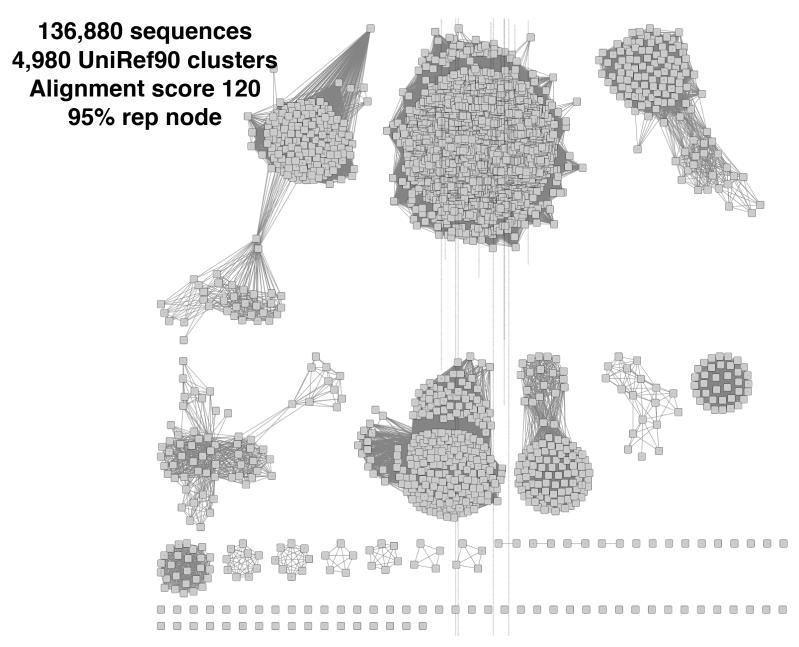
Functions of RuBisCO-Like Proteins (RLPs)?





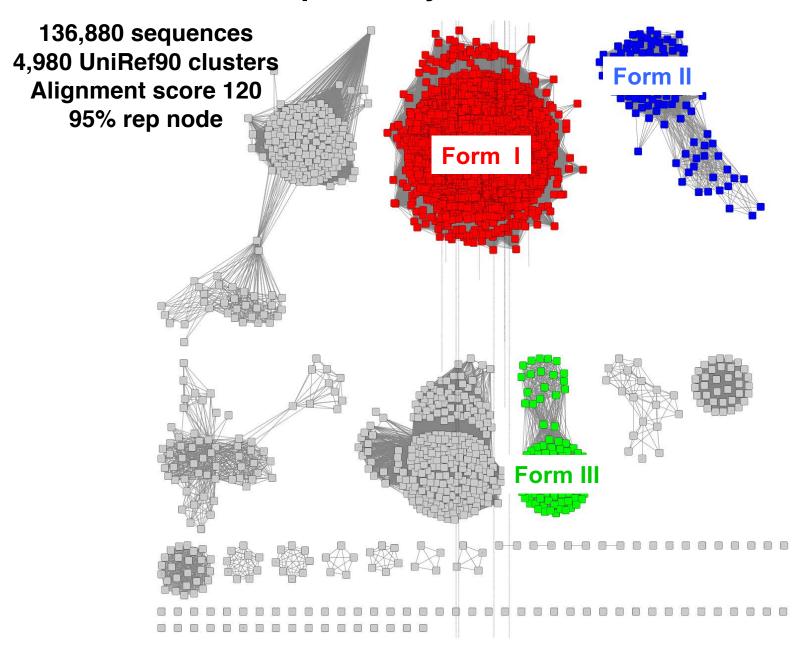
RuBisCO superfamily (PF00016)





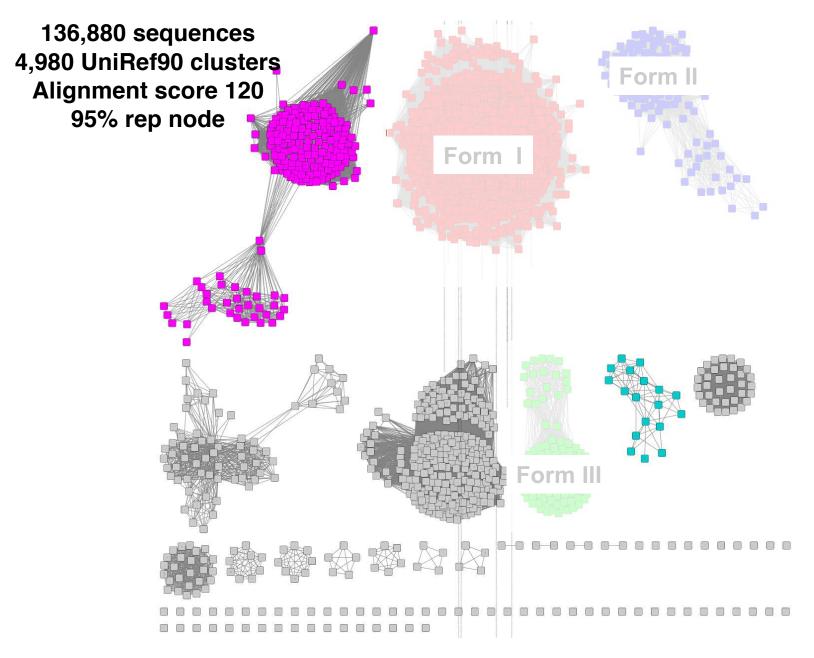
RuBisCO superfamily: three forms of RuBisCO





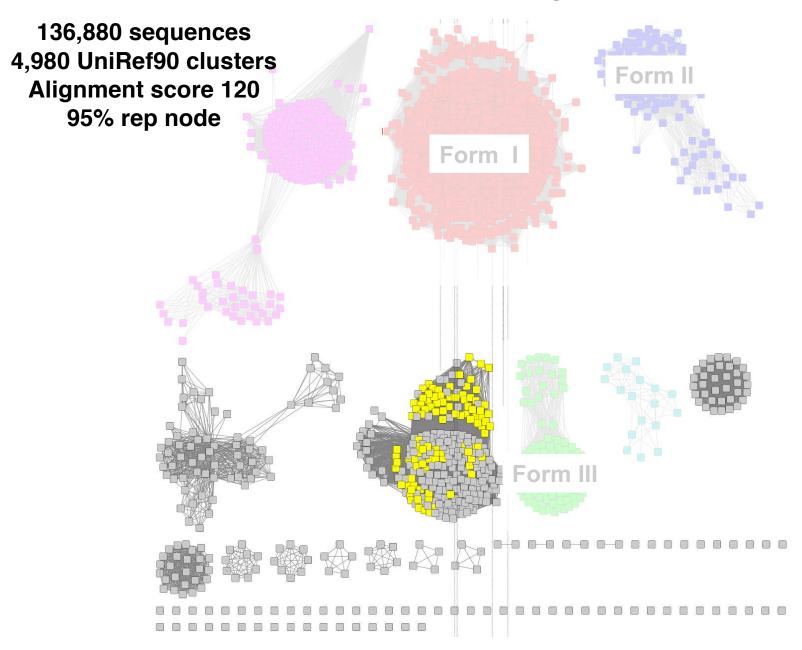
RLPs: Not RuBisCOs, enolate intermediates





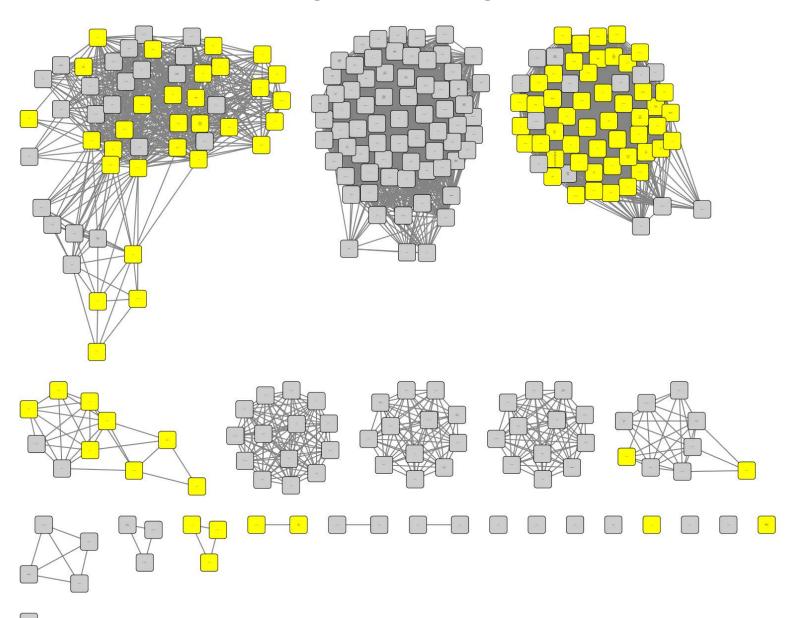
PF16896-RLP Neighbors





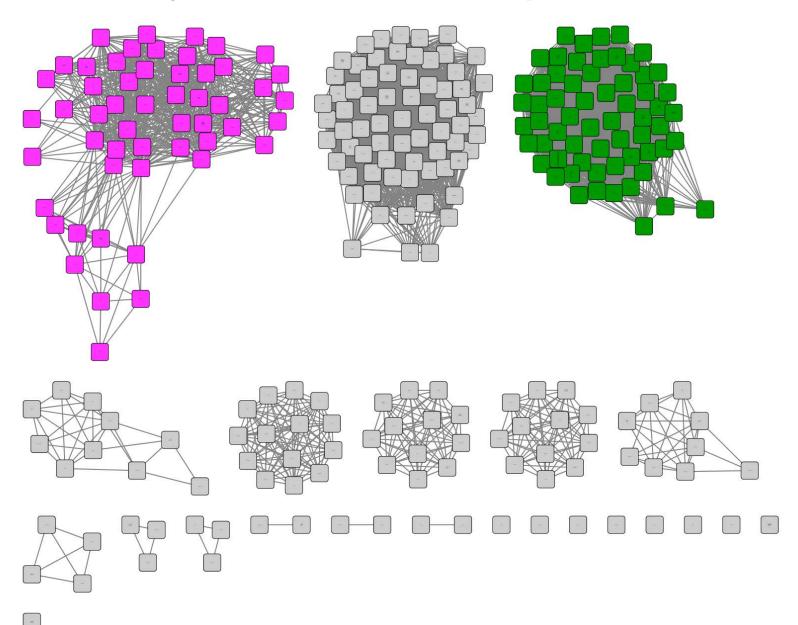
PF16896-RLP Neighbors: alignment score 145





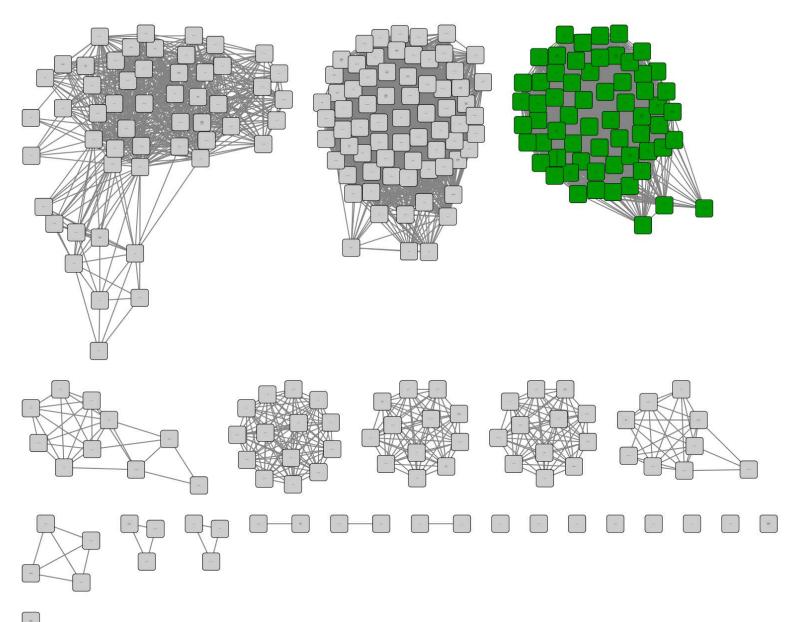
Two major families of PF16896-proximal RLPs





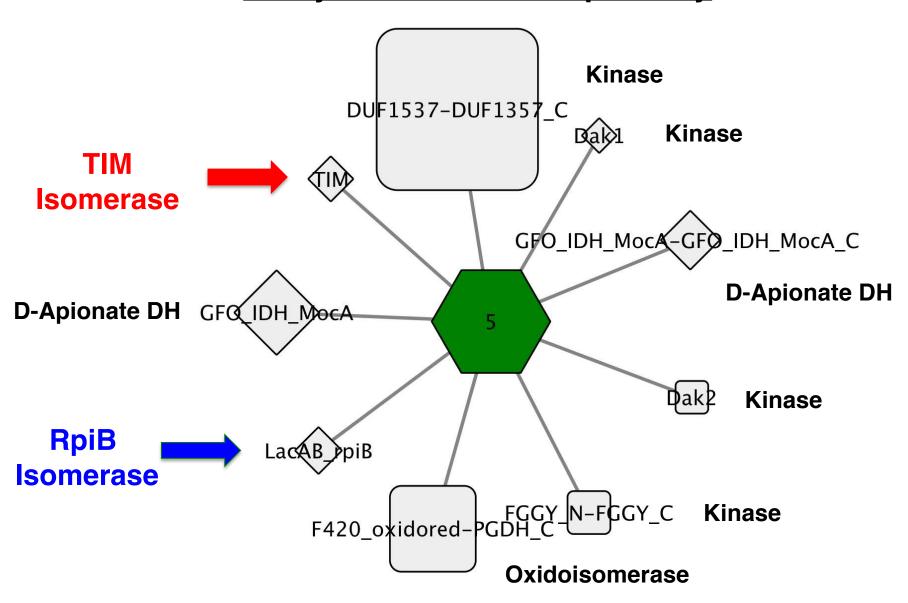
Two major families of PF16896-proximal RLPs





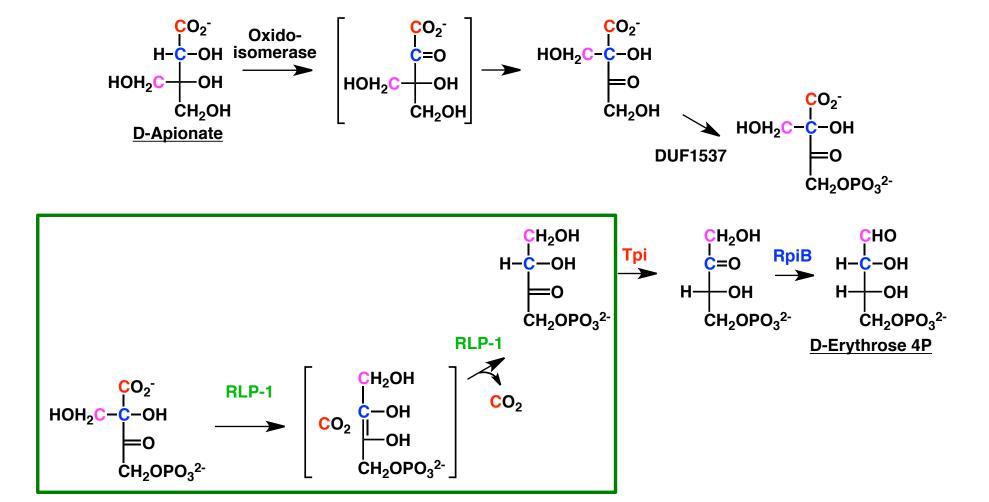


RLP-1: orthologues of TIM and RpiB isomerases in "xylose isomerase" pathway

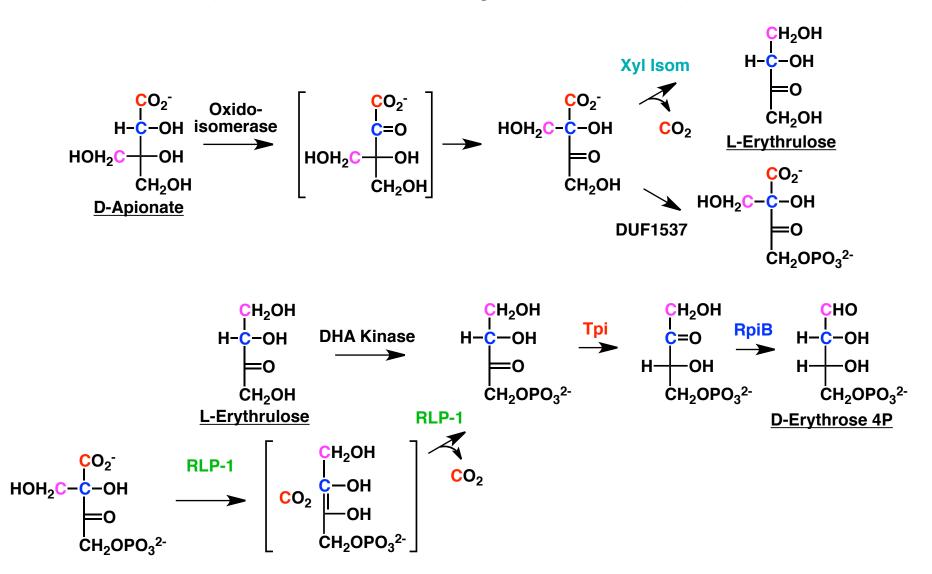


RLP-1: decarboxylase in the RuBisCO superfamily



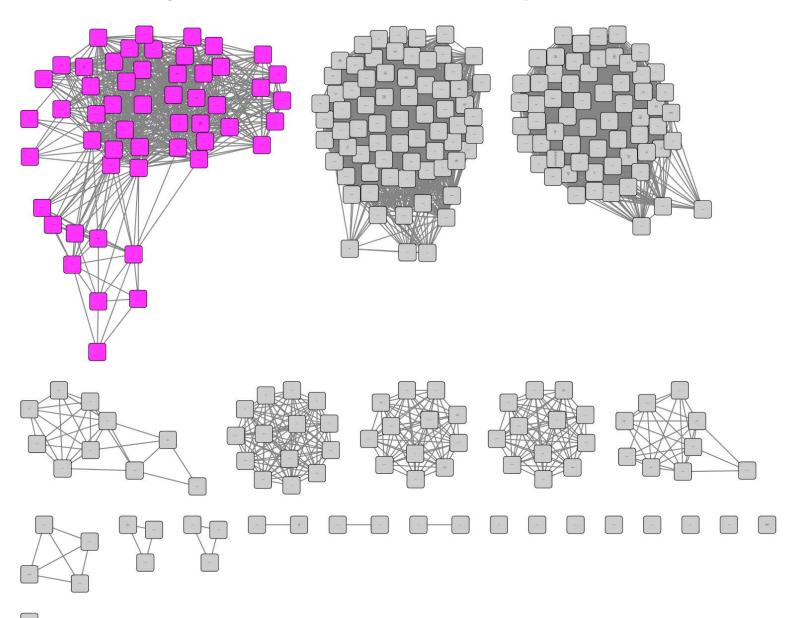


Two pathway variants to D-erythrose 4P (distinct decarboxylase families)



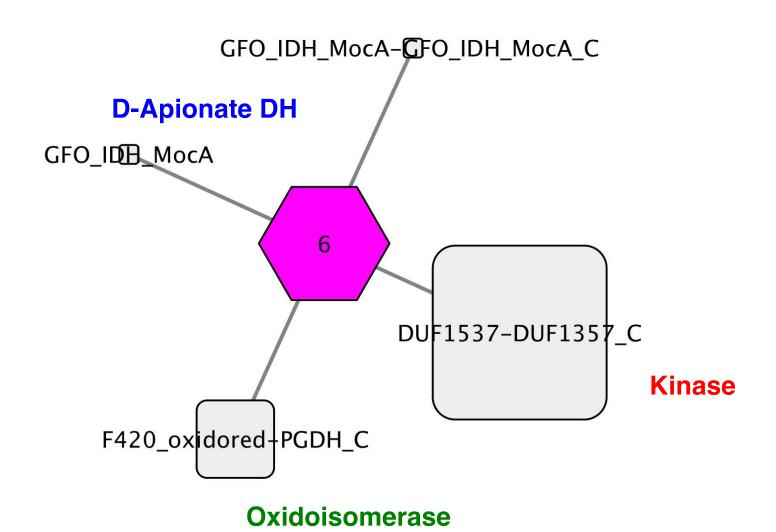
Two major families of PF16896-proximal RLPs





RLP-2: kinase, dehydrogenase, and oxidoisomerase





RLP-2: "transcarboxylase" in the RuBisCO superfamily





RLP-1/RLP-2: conserved decarboxylation, divergent fates of common enediolate intermediate

$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{H} - \text{C} - \text{OH} \\ = \text{O} \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

$$\begin{array}{c} \text{RLP-1} \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{CO}_2 \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

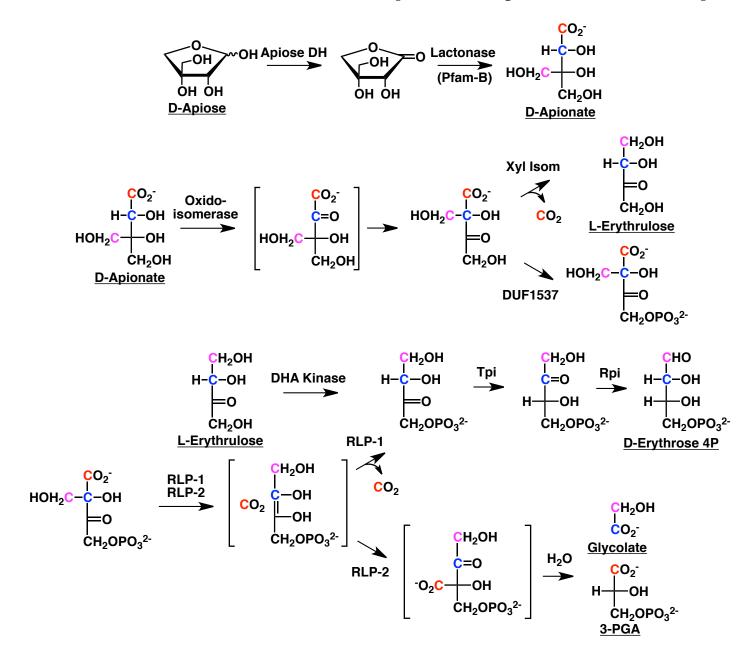
$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{CO}_2 \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{CO}_2 \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

$$\begin{array}{c} \text{CH}_2\text{OH} \\ \text{CO}_2 \\ \text{CO}_2 \\ \text{CH}_2\text{OPO}_3^{2-} \end{array}$$

Three SBPs → four different pathways in >1300 species





Today



- 1. Target selection using sequence similarity networks: EFI-EST web tool
- 2. Pathway discovery using genome neighborhood networks: EFI-GNT web tool
- 3. Example: using SSNs and GNNs to discover catabolic pathways for D-apiose
- 4. Target prioritization using chemically guided functional profiling: CGFP-ShortBRED beta tool

Chemically guided functional profiling



RESEARCH ARTICLE

MICROBIOTA

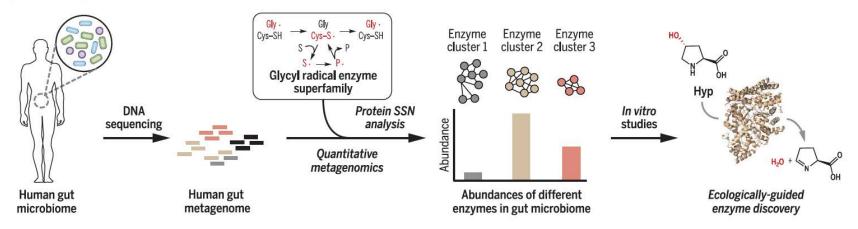
A prominent glycyl radical enzyme in human gut microbiomes metabolizes *trans*-4-hydroxy-L-proline

B. J. Levin, ^{1*} Y. Y. Huang, ^{1*} S. C. Peck, ¹ Y. Wei, ² A. Martínez-del Campo, ¹ J. A. Marks, ¹ E. A. Franzosa, ^{3,4} C. Huttenhower, ^{3,4} E. P. Balskus ^{1,4}†

The human microbiome encodes vast numbers of uncharacterized enzymes, limiting our functional understanding of this community and its effects on host health and disease. By incorporating information about enzymatic chemistry into quantitative metagenomics, we determined the abundance and distribution of individual members of the glycyl radical enzyme superfamily among the microbiomes of healthy humans. We identified many uncharacterized family members, including a universally distributed enzyme that enables commensal gut microbes and human pathogens to dehydrate *trans-4*-hydroxy-L-proline, the product of the most abundant human posttranslational modification. This "chemically guided functional profiling" workflow can therefore use ecological context to facilitate the discovery of enzymes in microbial communities.

Chemically guided functional profiling

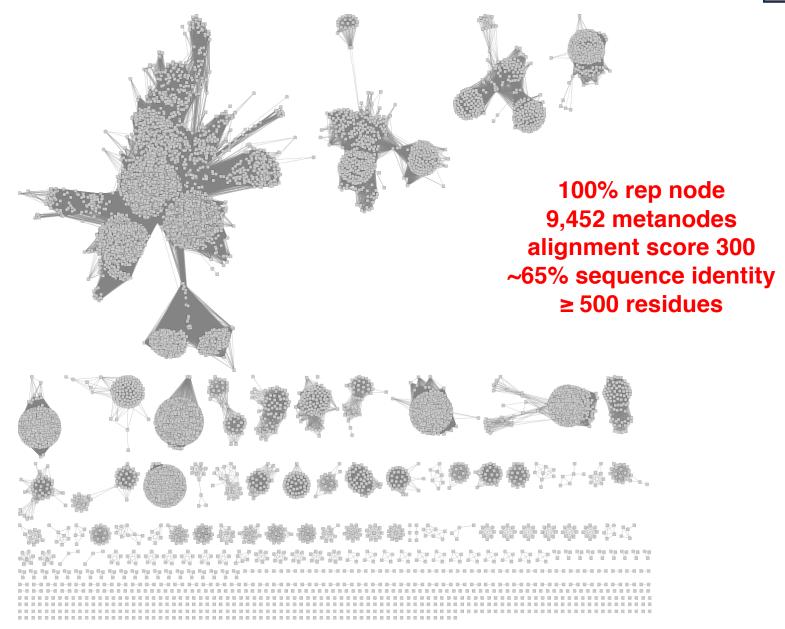
Chemically guided functional profiling enables enzyme discovery in microbiomes. Combining protein sequence similarity network (SSN) analysis with quantitative metagenomics reveals the abundance of both characterized and uncharacterized members of enzyme superfamilies. An analysis of glycyl radical enzymes in healthy human microbiomes facilitated the discovery of *trans*-4-hydroxy-L-proline dehydratase, a universally distributed but previously unknown gut microbial enzyme.



Levin et al., Science 355, 595 (2017) 10 February 2017

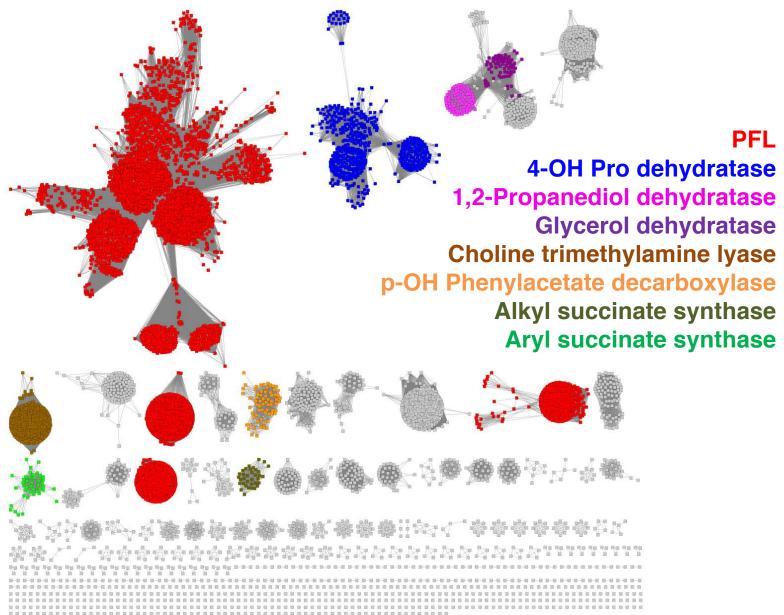
1 of 1

Glycyl radical enzyme (GRE) superfamily: 14,229 sequence



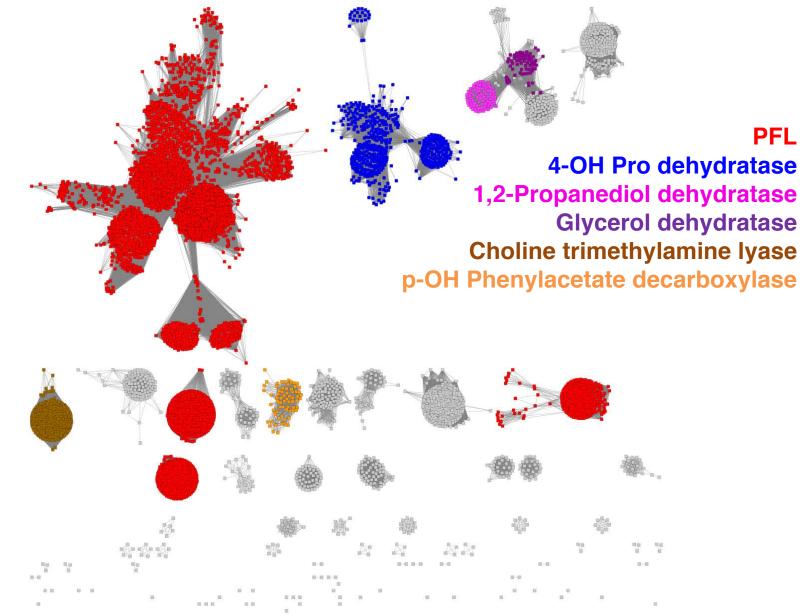
Clusters with known functions





Clusters with orthologues in human gut metagenomes





Beta tool: ShortBRED-CGFP





ENZYME FUNCTION INITIATIVE TOOLS

This website contains a collection of tools for creating and interacting with sequence similarity networks (SSNs) and genome neighborhood networks (GNNs). These tools originated in the Enzyme Function Initiative, a NIH-funded research project to develop a sequence / structure-based strategy for facilitating discovery of *in vitro* enzymatic and *in vivo* metabolic / physiological functions of unknown enzymes discovered in genome projects. From the initial sequence- or PFam-family-based SSN creation tool, this website has grown to include the following capabilities:

- SSN creation from sequence BLAST, PFam or InterPro famil(ies), FASTA sequences, and UniProt and/or NCBI protein accession IDs.
- Existing SSN coloring.
- GNN creation from uploaded SSNs.
- Genome neighborhood diagrams from sequence BLAST, protein sequence ID lookup, or FASTA sequence lookup.

In addition to the EFI utilities, this website also includes high-level access to the ShortBRED program developed by the Huttenhower Lab at the Harvard T.H. Chan School of Public Health. The version of ShortBRED available at this website accepts a SSN as input and develops peptide markers. The markers can be utilized in additional computations (also available on this website) to quantify relative abundances of functional gene families in human metagenomes.

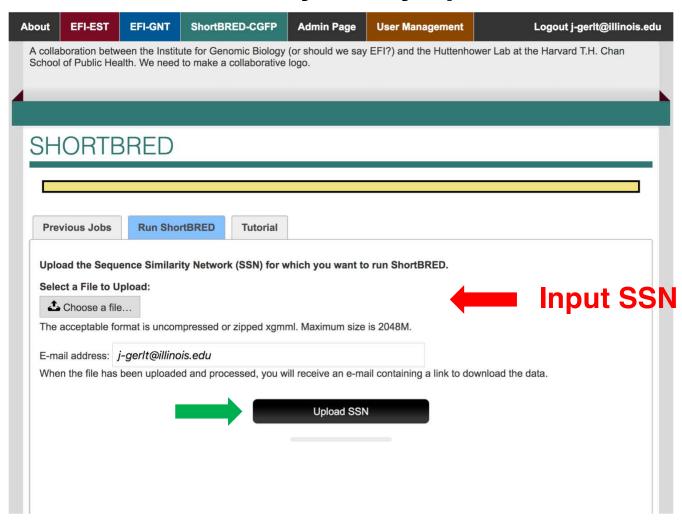
Need help or have suggestions or comments? Please click here.

Enzyme Function Initiative | 1206 W. Gregory Drive Urbana, IL 61801 | efi@enzymefunction.org



ShortBRED: Identify family-specific markers





ShortBRED: Quantitate metagenome abundance

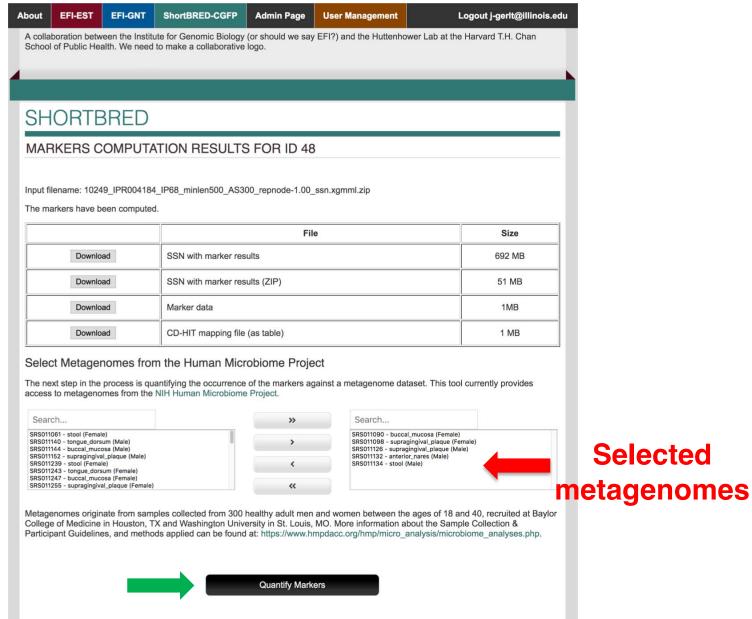


SHORTBRE	D					
	UTATION RESULTS F	OR ID 48				
MAI INEI 10 COM	OTATION NEGOETOT	OITID 40				
put filename: 10249_IPR0	04184_IP68_minlen500_AS300_	_repnode-1.00_ss	n.xgmml.zip			
he markers have been con	nputed.					
		File				
Download	SSN with marker results	SSN with marker results				
Download	SSN with marker results	SSN with marker results (ZIP)				
Download	Marker data	Marker data				
Download	CD-HIT mapping file (as	CD-HIT mapping file (as table)				
Select Metagenomes	s from the Human Microb	oiome Proiect				
he next step in the process	s is quantifying the occurrence of	the markers agair		his tool currently provides		
ccess to metagenomes fro	m the NIH Human Microbiome Pr	oject.				
Search		»	Search			
SRS011061 - stool (Female) SRS011090 - buccal_mucosa (Femal SRS011098 - supragingival_plaque (Female)	>				
SRS011126 - supragingival_plaque (I SRS011132 - anterior_nares (Male) SRS011134 - stool (Male)	Male)	<				
SRS011140 - tongue_dorsum (Male)						

Metagenome library

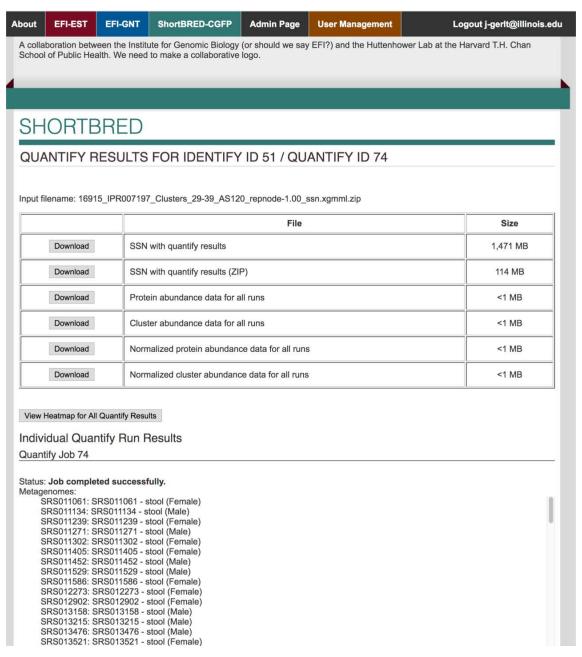
ShortBRED: Quantitate metagenome abundance





ShortBRED: Output files





Genomic enzymology strategy for functional assignment



1. Explore sequence function-space in a protein family, segregating family into isofunctional clusters using known functions (SwissProt, literature)

Sequence similarity networks (SSNs) using EFI-EST

2. Predict novel functions of isofunctional clusters in metabolic pathways using genome context (types of reactions enabled by proximal enzymes)

Genome neighborhood networks (GNNs) using EFI-GNT

3. Prioritize clusters for functional characterization by (metagenome) abundance

Chemically guided functional profiling (CGFP)

Acknowledgements

UCSF

Chakrapani Kalyanaraman Sara Calhoun Magdalena Korczynska Matt Jacobson Andrej Sali Brian Shoichet Patricia Babbitt

Harvard
Emily Balskus
Curtis Huttenhower

NIH U54GM093342 NIH P01GM118303 Albert Einstein
Matt Vetting
Jeff Bonanno
Agni Ghosh
Scott Garforth
Tyler Grove
Art Arcinis
Anthony Gizzi

Steven Almo

Remi Zallot
Dan Davidson
David Slater
Michael Carter
Xinshuai Zhang
Hua Huang
Brian San Francisco
Tyler Stack
John Cronan
John Gerlt

Illinois

Nils Oberg

EFI-EST: efi.igb.illinois.edu/efi-est/ EFI-GNT: efi.igb.illinois.edu/efi-gnt/

Requirements

- Registered account for EFI
- http://efi.igb.illinois.edu
- Java SE Runtime Environment 8
- http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html
- Cytoscape 3.3
- Form experience, this version (not the latest) is the best for our needs
- http://chianti.ucsd.edu/cytoscape-3.3.0/
- http://www.cytoscape.org/download old versions.html
- Within Cytoscape, install:
- BridgeDB from the App manager
- All files generated and used are available in the Workshop Dropbox folder
- Direct link: https://www.dropbox.com/sh/zfy0e6jaa8dttz7/AABkrFkcET5j2WdGtANKid-Pa?dl=0