## Back to the beginning: Which sequences to align?

IPAM MSA Workshop January 2015

Patsy Babbitt <u>babbitt@cgl.ucsf.edu</u> University of California, San Francisco

## **MSAs: Critical for structure-function studies**

|                                           | 10                                                                                           | 151                                                                                                                                            | 176                                                                                 |
|-------------------------------------------|----------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------|
|                                           | *                                                                                            | *                                                                                                                                              | *                                                                                   |
| Cu++ATPase.Ec                             | L DT V V F D KT GT LT E G                                                                    | V I A G VL PD G <mark>K A</mark> E A I KH L                                                                                                    | A M V G D G I N D A P A L                                                           |
| Cu++ATPase.Hs                             | V K V V V F D KT GT IT H G                                                                   | V F A E VL P SH K V A KV K Q L                                                                                                                 | A M V G D G I N D S P A L                                                           |
| Ca++ATPase.At                             | A TT I C S D KT GT LT T N                                                                    | V M A R S S P M D K H T L V R L L                                                                                                              | A V T G D G T N D A P A L                                                           |
| Urf.Mj                                    | K V A I V F D S A GT L V K I                                                                 | E A H Q E L K R D L I R N L                                                                                                                    | I M V G D G A N D V P A M                                                           |
| PhosSerPhos.Hs                            | A D A V C F D V D ST V I R E                                                                 | T A E - S G G K G K VI K L L K E                                                                                                               | I M I G D G A T D M E A C                                                           |
| 2-D0-6-PPhos.Sc                           | V D L C L F D L D GT I V S T                                                                 | I T G F D V K N G K P D P E G Y S                                                                                                              | V V F E D A P V G I K A G                                                           |
| DL-Gly-3-Phos.Sc                          | I N A A L F D V D GT I I I S                                                                 | I T A N D V K Q G K P H P E P Y L                                                                                                              | V V F E D A P A G I A A G                                                           |
| Phosphon.Pa<br>Phosphon.St<br>Phosphon.Bc | L Q A A I L D WA G T V V D F<br>I H A V I L D WA G T T V D F<br>I E A V I F D WA G T T V D Y | A T D E V - P N G R P W P A Q A L<br>A T D D L A A G G R P G P W M A L<br>T P D D V - P A G R P Y P W M S Y<br>T C C E S L P O P K D D A D L A | V K V D D T W P G I L E G<br>V K V D D A A P G I S E G<br>I K V G D T V S D M K E G |
| NtermDom.IGPD.Pp                          | VQALL L D MDGV MAEV                                                                          | LEDCPPKPSPEPIL                                                                                                                                 | AMVG D TVDDIIAG                                                                     |
| B-PhosGlucoMut.Ll                         | FKAVLFDLDGVITDT                                                                              | AEVAASKPAPDIFI                                                                                                                                 | IGLE D SQAGIQAI                                                                     |
| HaloAcidDehal.PspYL                       | IKGIAFDLYGTLFDV                                                                              | LSVDPVQVYKPDRVYE                                                                                                                               | LFVS S NAWDATGA                                                                     |
| NtermDomEpoxHyd.Hs                        | L RAAVI <mark>FULDGV</mark> LALP                                                             | 1 E S C Q V G M V <mark>K P</mark> E P Q I Y K                                                                                                 | L F L S D I H Q E L D A A                                                           |
| EnolasePhos.Ko                            | I RAIV <mark>TD IEGT</mark> TSDI                                                             | F D -  -  T L V G A <mark>K R</mark> E A Q S Y R                                                                                               |                                                                                     |



How have enzymes evolved to catalyze the many different chemical reactions required of living organisms?

# ...by re-using a limited set of "privileged" structural scaffolds



Enolase: 40 rxns Amidohydrolase: ~100 rxns



Haloalkanoic acid dehalogenase: 30-50 rxns



Glutathione transferase:44 major subgroups: ?? rxns



Isoprene synthase I: ?? rxns >50,000 products



Crotonase: 20 rxns

### Each functionally diverse enzyme superfamily (SF) links a conserved active site architecture to a fundamental catalytic capability



Jensen, R.A. Ann. Rev. Microbiol 30:409-425(1976) Petsko, G. A. et al., TIBS 18: 372-376 (1993) Babbitt & Gerlt, JBC 272: 30591 (1997) Gerlt & Babbitt, Ann Rev Biochem 70: 209 (2001)



This "chemistry constrained" model restricts the search space for functional inference for all members of the superfamily



>39,000 sequences, 23+ known reactions, many biological functions

# Predicting reaction & substrate specificity is much harder ...





#### Goals

- > Identify the fundamental chemical capability associated with the conserved features of the structural superfamily and the catalytic machinery shared by all members
- > Determine the patterns by which divergence has altered the structural scaffold and active site architecture to enable many different chemical reactions using the superfamily scaffold

#### **Applications**

- > Inform mechanism in knowns
- > Use a SF's fundamental structure-function mapping to limit the search space for functional inference in unknowns
- > Curation of functionally diverse enzyme superfamilies on a large scale
- > Guide identification of starting scaffolds for enzyme engineering
  - which superfamily scaffolds "know" how to catalyze the thermodynamically difficult step for a reaction of interest?

#### What does this have to do with MSA creation?

- > How to take advantage of the volume of [diverse] sequence and structure data now available
- > How to incorporate more functionally/biologically informed ways to choose which sequences/structures to align

#### This chemistry-constrained model suggests a general classification scheme



Annotation transfer only at the level of granularity at which good supporting evidence exists



#### See description in Nucleic Acids Res. 2014 Jan 1;42:D521-30 2

#### What is the Structure-Function Linkage Database (SFLD)?

- A hierarchical classification of enzymes that relates specific sequence-structure features to specific chemical capabilities
- A collection of tools and data for investigating sequencestructure-function relationships and hypothesizing function
- More...

#### How can I use the SFLD?

(see the tutorials for examples)

- Classify a sequence using Hidden Markov Models or BLAST search
- Browse superfamilies in the SFLD
- · Browse reactions (overall)
- Search for a specific enzyme (by name, sequence database ID, or PDB ID)
- View sequence alignments
- · View structures in Chimera
- Download data: sequence sets, multiple alignments, sequence similarity networks...

#### What makes the SFLD unique?

- Superfamilies are defined by a conserved chemical capability such as a partial reaction, families by a conserved overall reaction (more...)
- Conserved partial reactions are correlated with associated active site similarities
- Large-scale summaries of relationships between and within groups of enzymes are provided as sequence similarity networks

#### **Projects under Development**

- Extended SFLD (XSFLD)
- Identifying Potential Misannotations

|       | Universit     | y of California, San Fran  | icisco   About I | UCSF   UCSF Medical Center | SFLD             | RBV            |
|-------|---------------|----------------------------|------------------|----------------------------|------------------|----------------|
| Home  | About SFLD    | Documentation <del>-</del> | Tutorials        | Contact Us                 |                  | Curator's Entr |
|       |               | Browse by                  | / Superfamil     | ly Browse by Reaction      | Search by Enzyme |                |
| Top L | evel          | Namo                       | )                |                            |                  |                |
| L Su  | perfamily (co | e) Enola                   | ise              |                            |                  |                |
| LS    | ubgroup       | muco                       | nate cycloiso    | omerase                    |                  |                |
| L     | Family        | N-suc                      | cinylamino a     | acid racemase 2            |                  |                |
|       |               |                            |                  |                            |                  |                |



|                                                           | Total           | 100% 🖵                                 | <100% 🖵                   |
|-----------------------------------------------------------|-----------------|----------------------------------------|---------------------------|
| Functional domains                                        | 196             | 0                                      | 196                       |
| UniProtKB                                                 | 255             | 0                                      | 255                       |
| GI                                                        | 694             | 0                                      | 694                       |
| Structures                                                | 3               |                                        |                           |
| Reactions                                                 | 1               |                                        |                           |
| Functional domains of this<br>New functional domains were | family were las | st updated on Jar<br>this family on De | n. 7, 2015<br>c. 31, 2014 |

|                                                                                                                  |               |             |                         | Pairwise % I | dentities  |                          |                  |                |
|------------------------------------------------------------------------------------------------------------------|---------------|-------------|-------------------------|--------------|------------|--------------------------|------------------|----------------|
| 0                                                                                                                | 170           | 180         | 190_                    | 200          | 210        | <u>2</u> 20              | 230              | 240            |
| FKMKVGT                                                                                                          | EVTRDVARI     | KAVRQQVGEDI | AIRV <mark>DVNQG</mark> | WENAATTLÇ    | GLRAMKDLNI | DWL <mark>EQPVDS</mark>  | EDIDGMVEIK       | SKSDVPLN       |
| 'LKI <mark>KVGT</mark>                                                                                           | SIEEDVARI     | KAVRSRVGNDI | AIRV <mark>DVNQG</mark> | WKTSTAALK    | ALKQLEELQI | DWV <mark>E</mark> QPVAA | DDIDGLAEVE       | AKIAIPVN       |
| FKM <mark>KVGTNVKEDVKRIEAVRERVGNDIAIRV</mark> DVNQGWKNSANTLTALRSLGHLNIDWI <mark>E</mark> QPVIADDIDAMAHIRSKTDLPLN |               |             |                         |              |            |                          |                  |                |
| <b>FKM</b> KVGT                                                                                                  |               | GAVINGINUL  | MIN DVNQG               | WINDANTIT    | Анконопин  | Sur Sala                 |                  |                |
| Position                                                                                                         | Amino<br>acid | Function    | AINUDVNQG               |              |            |                          | Curator<br>notes | Eviden<br>Code |



| Position | Amino<br>acid | Function                                                              | Curator<br>notes | Eviden<br>Code                 |
|----------|---------------|-----------------------------------------------------------------------|------------------|--------------------------------|
| 163      | Lys (K)       | base (abstracts alpha proton), acid (donates proton to leaving group) |                  |                                |
| 191      | Asp (D)       | metal binding ligand                                                  |                  | ICS                            |
| 218      | Glu (E)       | metal binding ligand                                                  |                  | ICS                            |
| 243      | Asp (D)       | metal binding ligand                                                  |                  | ICS                            |
| 267      | Lys (K)       | base (abstracts alpha proton), acid (donates proton to leaving group) | Catal            | vzed Reaction(s)               |
| -        |               |                                                                       | race             | emization of n-succinylamino a |

| Protein Name                       | Superfamily | Family                             | Species ↓                           | Databases 🛡                  | UniProtKB  | MicrobesOnline<br>Operon | The SEED              | PDB<br>ID            | Structures | Updated          |
|------------------------------------|-------------|------------------------------------|-------------------------------------|------------------------------|------------|--------------------------|-----------------------|----------------------|------------|------------------|
| N-succinylamino<br>acid racemase 2 | Enolase     | N-succinylamino<br>acid racemase 2 | Bacillus anthracis str.<br>Ames ⊡ ↓ | GI REF UP<br>MO TS MB<br>PDB | A0A084CRN0 | 7611946                  | fig 486623.3.peg.5055 | 2P8B<br>2P8C<br>2P88 | 3          | Jan. 07,<br>2015 |
| N-succinylamino<br>acid racemase 2 | Enolase     | N-succinylamino<br>acid racemase 2 | Bacillus cereus ATCC<br>14579 ⊡ ↓   | GI REF UP<br>MO TS MB<br>PDB | Q811L5     | 357748                   | fig 226900.1.peg.322  | 2P8B<br>2P8C<br>2P88 | 3          | Jan. 07,<br>2015 |
| N-succinylamino<br>acid racemase 2 | Enolase     | N-succinylamino<br>acid racemase 2 | Bacillus cereus ATCC<br>10987 ☐ ↓   | GI REF UP<br>MO TS MB<br>PDB | Q73EC5     | 643701                   | fig 222523.1.peg.433  | 2P8B<br>2P8C<br>2P88 | 3          | Jan. 07,<br>2015 |
| N-succinylamino<br>acid racemase 2 | Enolase     | N-succinylamino<br>acid racemase 2 | Bacillus cereus ☐ ↓                 | GI REF UP<br>TS MB PDB       | Q4MJ00     |                          | fig 269801.1.peg.4979 | 2P8B<br>2P8C<br>2P88 | 3          | Jan. 07,<br>2015 |
| N-succinylamino<br>acid racemase 2 | Enolase     | N-succinylamino<br>acid racemase 2 | Bacillus thuringiensis ⊟<br>↓       | GI REF UP<br>MO MB<br>PDB    | Q6HP62     | 595563                   |                       | 2P8B<br>2P8C<br>2P88 | 3          | Jan. 07,<br>2015 |





#### Core SFLD

| Suprafamily      | Superfamily                                           | Subgroups | Families | Sequences | Structures | Reactions |
|------------------|-------------------------------------------------------|-----------|----------|-----------|------------|-----------|
|                  | Amidohydrolase                                        | 11        | 89       | 79238     | 474        | 41        |
|                  | Aromatic Prenyltransferase                            | 2         | 0        | 339       | 18         | 0         |
|                  | Crotonase                                             | 2         | 27       | 75290     | 172        | 28        |
|                  | Enolase                                               | 7         | 20       | 39661     | 355        | 22        |
|                  | Haloacid Dehalogenase                                 | 25        | 22       | 79778     | 570        | 21        |
|                  | Isoprenoid Synthase Type I                            | 14        | 69       | 16579     | 359        | 65        |
|                  | Isoprenoid Synthase Type II                           | 4         | 8        | 7645      | 202        | 8         |
|                  | Nucleophilic Attack 6-Bladed Beta-<br>Propeller (N6P) | 3         | 3        | 31085     | 85         | 2         |
|                  | Radical SAM                                           | 49        | 93       | 113568    | 52         | 66        |
|                  | RuBisCO                                               | 2         | 2        | 41212     | 73         | 2         |
| Thioredoxin Fold | Glutathione Transferase (cytosolic)                   | 42        | 0        | 13097     | 432        | 0         |
| Suprafamily      | Peroxiredoxin                                         | 6         | 0        | 12239     | 179        | 0         |
|                  | <b>TOTAL:</b> 12                                      | 167       | 333      | 515791    | 2971       | 255       |

#### Next steps: Identification & integration of conserved chemical capabilities in superfamilies & patterns by which their chemistry varies

| ~ II                                              | - н                  | п -п        |                    | ~ <b>1</b> | 1 - Ц            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
|---------------------------------------------------|----------------------|-------------|--------------------|------------|------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 11                                                | t n                  | п, п        |                    | ÷I         | ц., Ц            | <b>A</b>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
|                                                   |                      | n Num       | ber of a           | 71         |                  | Conserved #                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| Superfamily ¤                                     | ID =                 | Reactions   | numbers            | ~ 11       | EC<br>positions  | Substructure =                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
| т II                                              | 1 н                  | 7.1 ≤ Π     | ~ I ~ H            |            | с и              | ц.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| <u>Alkaline</u><br>phosphatase-like ¤             | :<br>c.76.1 ≍        | - 1<br>67 ¤ | ≂ı<br>5 ¤          |            | <u>3.1.x.x</u> ¤ |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| <u>SGNH hydrolase</u> ×                           | c.23.10              | -1<br>19 ¤  | - 1<br><b>2</b> 11 |            | <u>3.1.1.x</u> ¤ | царана и обращание и обращи и обращи и обращание и обращани и обращани и обращани и обращ |
| <u>Metallo-dependent</u><br>phosphatases ≖        | d.159.1              | ~ı<br>30 ¤  | ~ 1<br><b>2</b> ¤  |            | <u>3.1.3.x</u> ¤ | COP(0)(=0)0                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| <u>Carbohydrate</u><br>phosphatase<br>≍           | <sup>⊭</sup> e.7.1 ¤ | -ı<br>17¤   | ""<br>4 ¤          |            | <u>3.1.3.x</u> ¤ |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Cobalamin (vitamin<br>B12)-dependent<br>enzymes ¤ | c.1.19⊭              | - 1<br>6 ¤  | °⊐<br>2 ¤          |            | <u>4.2.1.x</u> ⊭ | о                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Phosphoglycerate<br>mutase-like ¤                 | י<br>c.60.1 ¤        | - I<br>58 ¤ | - 1<br>5 ¤         |            | x.x.x.x ¤        |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| <u>Six-hairpin</u><br>glycosidases =              | a.102.1              | - 1<br>6 ¤  | - 1<br><b>3</b> ¤  |            | <u>X.X.X.X</u> ¤ |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| <u>alpha/beta-</u><br><u>Hydrolases</u> ⊭         | c.69.1 ⊧             | - I<br>13 ¤ | -1<br>3 ¤          |            | <u>3.1.1.x</u> ¤ | - :<br>CCOC(C)=0 ::                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |

Chiang et al, PLoS CB, 4:e1000142 (2008)



Gemma Holliday PhD

### Tackling sequence data on a large scale

Most protein superfamilies are too large to manage & easily explore using multiple sequence alignments & trees



We estimate that ~1/3 of the universe of enzyme superfamilies are functionally diverse Almonacid & Babbitt, Curr Op Biol Chem 15:435 (2011)

### Accessing the larger context: Protein similarity networks



- > Powerful hypothesis generator for structure-function relationships
- > Interactive
- > Fast
- > Easy visualization (Cytoscape)
- > Handles thousands of sequences, structures, etc.
- > Pairwise comparisons don't require a multiple alignment
- > Track well with known structurefunction relationships
- > Capture sequence, structure, ligand, relationships
- > Many metrics & algorithms for comparison & clustering of similarity data
- > Networks not a substitute for phylogenetic trees

Node = sequence (or structure)

Edge = connections between sequences w scores as good as the E-value cutoff threhold Sequence networks in this presentation: all-byall BLAST (*E*-values as scores); structure: all-

by-all FAST/TM-align scores

## Validation

Comparison of distance metrics for generating networks shows that BLAST correlates well with several other metrics

| Uronatel                                            | BLAST | SW          | MA                   | PT                            |
|-----------------------------------------------------|-------|-------------|----------------------|-------------------------------|
| MLE                                                 |       |             |                      |                               |
| BLAST                                               |       | 0.999       | 0.971                | 0.953                         |
| Smith-Waterman (SW)                                 | 0.998 |             | 0.970                | 0.953                         |
| Multiple Alignment (MA)                             | 0.800 | 0.798       |                      | 0.974                         |
| Phylogenetic Tree (PT)                              | 0.731 | 0.731       | 0.777                |                               |
|                                                     |       |             |                      |                               |
| NagA                                                | BLAST | SW          | MA                   | PT                            |
| NagA                                                | BLAST | SW          | MA                   | PT                            |
| BLAST NagA                                          | BLAST | SW<br>0.997 | MA<br>0.841          | PT<br>0.748                   |
| NagA   BLAST   Smith-Waterman (SW)                  | BLAST | SW<br>0.997 | MA<br>0.841<br>0.846 | PT<br>0.748<br>0.753          |
| NagABLASTSmith-Waterman (SW)Multiple Alignment (MA) | BLAST | SW<br>0.997 | MA<br>0.841<br>0.846 | PT<br>0.748<br>0.753<br>0.719 |

R<sup>2</sup> values for linear regressions of distances generated from various metrics for scoring similarity among sequences

Other validation analyses show

- > Network topologies are generally robust to missing data
- > Two-dimensional distances in visualized networks correlate well with the underlying distances in high-dimensional space
- > See Atkinson et al, PLoS ONE, 4: e4345 (2009) for more statistical validation of PSNs

## Layout used in this talk



Organic layout: Edge lengths represent degree of connectivity, track with dissimilarity Colors depict function or other types of functional information



| Suprafamily      | Superfamily                                           | Subgroups | Families | Sequences | Structures | Reactions |
|------------------|-------------------------------------------------------|-----------|----------|-----------|------------|-----------|
|                  | Amidohydrolase                                        | 11        | 89       | 79238     | 474        | 41        |
|                  | Aromatic Prenyltransferase                            | 2         | 0        | 339       | 18         | 0         |
|                  | Crotonase                                             | 2         | 27       | 75290     | 172        | 28        |
|                  | Enolase                                               | 7         | 20       | 39661     | 355        | 22        |
|                  | Haloacid Dehalogenase                                 | 25        | 22       | 79778     | 570        | 21        |
|                  | Isoprenoid Synthase Type I                            | 14        | 69       | 16579     | 359        | 65        |
|                  | Isoprenoid Synthase Type II                           | 4         | 8        | 7645      | 202        | 8         |
|                  | Nucleophilic Attack 6-Bladed Beta-<br>Propeller (N6P) | 3         | 3        | 31085     | 85         | 2         |
|                  | Radical SAM                                           | 49        | 93       | 113568    | 52         | 66        |
|                  | RuBisCO                                               | 2         | 2        | 41212     | 73         | 2         |
| Thioredoxin Fold | Glutathione Transferase (cytosolic)                   | 42        | 0        | 13097     | 432        | 0         |
| Suprafamily      | Peroxiredoxin                                         | 6         | 0        | 12239     | 179        | 0         |
|                  | <b>TOTAL:</b> 12                                      | 167       | 333      | 515791    | 2971       | 255       |

#### For the Core SFLD, networks can be downloaded at every level of the hierarchy: Superfamily, Subgroup, & Family



#### Extended SFLD

| Suprafamily | Superfamily                                         | Subgroups | Families | Sequences | Structures | Reactions |
|-------------|-----------------------------------------------------|-----------|----------|-----------|------------|-----------|
|             | Arginase/Deacetylase                                | 0         | 0        | 10570     | 177        | 0         |
|             | Carbohydrate Phosphatase                            | 0         | 0        | 12278     | 160        | 0         |
|             | Carbon-Nitrogen Hydrolase                           | 0         | 0        | 14974     | 49         | 0         |
|             | Chelatase                                           | 0         | 0        | 5776      | 49         | 0         |
|             | Cytidine Deaminase-Like                             | 0         | 0        | 18803     | 90         | 0         |
|             | Di-trans-poly-cis-decaprenylcistransferase          | 0         | 0        | 4075      | 36         | 0         |
|             | dUTPase-Like                                        | 0         | 0        | 6343      | 156        | 0         |
|             | Ferric Reductase Domain                             | 2         | 2        | 282       | 1          | 0         |
|             | Fumarylacetoacetase, C-terminal-related             | 0         | 0        | 10104     | 32         | 0         |
|             | Glutaminase/Asparaginase                            | 0         | 0        | 2672      | 37         | 0         |
|             | HD-Domain/PDEase-Like                               | 0         | 0        | 33746     | 252        | 0         |
|             | Histidine Phosphatase                               | 0         | 0        | 21228     | 254        | 0         |
|             | Isochorismatase-Like Hydrolases                     | 0         | 0        | 11412     | 42         | 0         |
|             | Kringle-Like                                        | 0         | 0        | 1693      | 627        | 0         |
|             | L-Aspartase-Like                                    | 0         | 0        | 17659     | 100        | 0         |
|             | Metalloproteases, Zincins                           | 0         | 0        | 21437     | 429        | 0         |
|             | Methyltransferase Domain 18                         | 0         | 0        | 4312      | 42         | 0         |
|             | Methyltransferase Domain 9                          | 0         | 0        | 433       | 1          | 0         |
|             | NUDIX Hydrolase Domain-Like                         | 0         | 0        | 38313     | 230        | 0         |
|             | Peptidase M24                                       | 0         | 0        | 17752     | 177        | 0         |
|             | Phosphatidylinositol Phosphodiesterase              | 3         | 5        | 11014     | 97         | 5         |
|             | Phospholipase C/P1 Nuclease                         | 0         | 0        | 1560      | 21         | 0         |
|             | Phosphonate Radical SAM                             | 1         | 1        | 902       | 0          | 1         |
|             | PLP-Binding Barrel                                  | 0         | 0        | 19071     | 95         | 0         |
|             | Proline Racemase                                    | 0         | 0        | 748       | 14         | 0         |
|             | Pyruvoyl-Dependent Histidine/Arginine Decarboxylase | 2         | 2        | 335       | 11         | 2         |
|             | Radical SAM 3-amino-3-carboxypropyl Radical Forming | 1         | 1        | 1644      | 2          | 1         |
|             | Ribulose-Phosphate Binding Barrel                   | 0         | 0        | 25997     | 249        | 0         |
|             | SGNH Hydrolase                                      | 0         | 0        | 19406     | 63         | 0         |
|             | Six-Hairpin Glycosidases                            | 0         | 0        | 28690     | 227        | 0         |
|             | SPOUT Methyltransferase                             | 0         | 0        | 2593      | 5          | 0         |
|             | Subtilisin-Like                                     | 0         | 0        | 15855     | 296        | 0         |
|             | Thioesterase/Thiol Ester Dehydrase-Isomerase        | 0         | 0        | 38098     | 226        | 0         |
|             | Xylose Isomerase-Like                               | 0         | 0        | 18333     | 196        | 0         |
|             | TOTAL: 34                                           | 9         | 11       | 438108    | 4443       | 9         |

Select Task →

Download Network Download Data Set

#### Sequence Similarity Networks

Download a Sequence Similarity Network of this superfamily (XGMML format  $\Box$ ).

Network downloads are XGMML files that are readable by program such as Cytoscape. In these networks, nodes represent proteins and edges represent pairwise similarities better than a given *E*-value cutoff. Additionally, these networks contain several attributes with data from the SFLD.



#### What sequences to align: The search for strictosidine variants from plants



#### Nucleophilic attack 6-bladed β-propeller (N6P) SF



Figure adapted from Ma X et al, The Plant Cell, 18:907 (2006)

### Characterized proteins are outliers in the N6P SF











Active sites of characterized SS enzymes are outliers as well

### Many SSL proteins are not from plants





|                   | 44                    | 151                      | 210                                | 254                     | 309      |
|-------------------|-----------------------|--------------------------|------------------------------------|-------------------------|----------|
| PON1 (1v04.pdb)   | GS <mark>E</mark> DLE | SV <mark>N</mark> DIVAVG | DVRVVAEGFDFA <mark>N</mark> GINISP | LV <mark>D</mark> NISVD | QGSTVAAV |
| Drp35 (2dg1.pdb)  | QL <mark>E</mark> GLN | CI <mark>D</mark> DMVFDS | TVTPIIQNISVA <mark>N</mark> GIALST | GP <mark>D</mark> SCCID | LRSTHPQF |
| DFPase (1pjx.pdb) | GA <mark>E</mark> GPV | GC <mark>N</mark> DCAFDY | QMIQVDTAFQFP <mark>N</mark> GIAVRH | GA <mark>D</mark> GMDFD | EKPSNLHF |
| SS (2fpb.pdb)     | APNSFT                | WLYAVTVDQ                | ETTLLLKELHVPGGAEVSA                | NPGNIKRN                | EHFDQIQE |
| gi 147772032      | GP <mark>E</mark> AIA | FL <mark>N</mark> AVDVDQ | EVTVLLRGLGGAGGVTISK                | TP <mark>D</mark> NIKRN | KTISEVQE |
| gi 22326950       | GP <mark>E</mark> SVA | FT <mark>N</mark> DLDIAD | KAVVLVSNLQFP <mark>N</mark> GVSISR | HP <mark>D</mark> NVRTN | RSVSEVEE |
| gi 125556119      | GP <mark>E</mark> SVA | FT <mark>N</mark> GVDIDQ | QVTVLQSNITYP <mark>N</mark> GVAISA | YP <mark>D</mark> NVRPD | RP-TEVMD |
| gi 24308201       | GP <mark>E</mark> SIA | FV <mark>N</mark> DLTVTQ | EVKVLLDQLRFP <mark>N</mark> GVQLSP | FP <mark>D</mark> NIRPS | TYISEVHE |
| gi 1280434        | GP <mark>E</mark> CLI | IF <mark>N</mark> GVTVSK | VSEVLLDELAFA <mark>N</mark> GLALSP | LP <mark>D</mark> NLTPD | T-ISHVLE |
| gi 125559158      | AP <mark>E</mark> DVY | FA <mark>D</mark> AAIEAS | EASVVLDGLGFA <mark>N</mark> GVALPP | NP <mark>D</mark> NIRLG | NMVTSVTE |
| gi 111017930      | GP <mark>E</mark> DVA | AC <mark>N</mark> NSAVGR | ETDLLAEGLQFA <mark>N</mark> GVGLAS | IP <mark>D</mark> NMTSQ | P-VTGVRE |
| gi 15596490       | GP <mark>E</mark> DTA | FT <mark>D</mark> DLDIAS | KTEVLLKDLYFA <mark>N</mark> GVALSA | LP <mark>D</mark> NLQGD | RMITSAKP |
| gi 52549517       | GP <mark>E</mark> DVA | LT <mark>D</mark> DVDIAA | TTRLVLNNLYFA <mark>N</mark> GVAVSP | FP <b>D</b> GISSN       | Q-ITSVQE |



E-value cut-off = 10<sup>-50</sup> 516 SSL subgroup sequences Median alignment length = 297 residues; median percent identity = 41% Colored by # of metal coordinating residues Red = 4 Yellow = 3 Green = 2 Cyan = 1

Gray = 0

|   |                   | 44                    | 151                      | 210                                | 254                     | 309      |
|---|-------------------|-----------------------|--------------------------|------------------------------------|-------------------------|----------|
| • | PON1 (1v04.pdb)   | GS <mark>E</mark> DLE | SV <mark>N</mark> DIVAVG | DVRVVAEGFDFA <mark>N</mark> GINISP | LV <mark>D</mark> NISVD | QGSTVAAV |
|   | Drp35 (2dg1.pdb)  | QL <mark>E</mark> GLN | CI <mark>D</mark> DMVFDS | TVTPIIQNISVA <mark>N</mark> GIALST | GP <mark>D</mark> SCCID | LRSTHPQF |
|   | DFPase (1pjx.pdb) | GA <mark>E</mark> GPV | GC <mark>N</mark> DCAFDY | QMIQVDTAFQFP <mark>N</mark> GIAVRH | GA <mark>D</mark> GMDFD | EKPSNLHF |
|   | SS (2fpb.pdb)     | APNSFT                | WLYAVTVDQ                | ETTLLLKELHVPGGAEVSA                | NPGNIKRN                | EHFDQIQE |
|   | gi 147772032      | GP <mark>E</mark> AIA | FL <mark>N</mark> AVDVDQ | EVTVLLRGLGGAGGVTISK                | TP <mark>D</mark> NIKRN | KTISEVQE |
|   | gi 22326950       | GP <mark>E</mark> SVA | FT <mark>N</mark> DLDIAD | KAVVLVSNLQFP <mark>N</mark> GVSISR | HP <mark>D</mark> NVRTN | RSVSEVEE |
|   | gi 125556119      | GP <mark>E</mark> SVA | FT <mark>N</mark> GVDIDQ | QVTVLQSNITYP <mark>N</mark> GVAISA | YP <mark>D</mark> NVRPD | RP-TEVMD |
|   | gi 24308201       | GP <mark>E</mark> SIA | FV <mark>N</mark> DLTVTQ | EVKVLLDQLRFP <mark>N</mark> GVQLSP | FP <mark>D</mark> NIRPS | TYISEVHE |
|   | gi 1280434        | GP <mark>E</mark> CLI | IF <mark>N</mark> GVTVSK | VSEVLLDELAFA <mark>N</mark> GLALSP | LP <mark>D</mark> NLTPD | T-ISHVLE |
|   | gi 125559158      | AP <mark>E</mark> DVY | FA <mark>D</mark> AAIEAS | EASVVLDGLGFA <mark>N</mark> GVALPP | NP <mark>D</mark> NIRLG | NMVTSVTE |
|   | gi 111017930      | GP <mark>E</mark> DVA | AC <mark>N</mark> NSAVGR | ETDLLAEGLQFA <mark>N</mark> GVGLAS | IP <mark>D</mark> NMTSQ | P-VTGVRE |
|   | gi 15596490       | GP <mark>E</mark> DTA | FT <mark>D</mark> DLDIAS | KTEVLLKDLYFA <mark>N</mark> GVALSA | LP <mark>D</mark> NLQGD | RMITSAKP |
|   | gi 52549517       | GP <mark>E</mark> DVA | LT <mark>D</mark> DVDIAA | TTRLVLNNLYFA <mark>N</mark> GVAVSP | FP <mark>D</mark> GISSN | Q-ITSVQE |

# The large-scale context suggests the great majority of SSL proteins are not SSs



Vitus vinifera SSL lacks SS activity but shows low levels of hydrolase activity typical of the SSL and arylesterase subgroups

## Challenges

**Technical issues** 

- > Extreme sequence diversity
- > Complex insert patterns and multiple domain architectures within a SF

Experimentally characterized proteins only poorly sample the available sequence and structure space

Especially for divergent proteins of functionally diverse enzyme SFs, clustering by sequence and structure similarity fails to track well with functional boundaries

Uniqueness of structure-function relationships in individual SFs complicates development of general solutions

# Alkaline phosphatase superfamily: Inserts to the common core distinguish known reaction classes 40,000 sequences



Many of these inserts are unrelated

- > Locations vary
- Multiple insertions in a single subgroup
- Structural insert patterns fail to track with variations in reaction or mechanism

#### What sequences to align? How well do characterized proteins sample SF sequence space? Most alkaline phosphatase SF members have never been experimentally or structurally characterized *E*-value threshold = $1 \times 10^{-13.1}$ 4590 nodes represent 14,403 sequences SwissProt family bacterial phospholipase C (40% ID filtered) CDP-alcohol phosphatidyltrasferase class-I metal one LTA synthase Median alignment length = 422 residues; opgB O phosphoenthanolamine transferase Median pairwise ID = 30%sulfatase alkaline phosphatase Small nodes: experimentally characterized nucleotide pyrophosphatase/phosphodiesterase metal O Phosphoglycerate mutase Large nodes: structurally characterized Phosphopentomutase PIGG/PIGN/PIGO OUnknown

Collaboration with Dan Herschlag

## Even very well-studied SFs are minimally characterized < 2% of GSTs experimentally shown to catalyze GST-like reactions



Mashiyama et al, PLoS Biol, 12: :e1001843 (2014)

#### Node color: Swiss-Prot family annotation



#### cytosolic Glutathione transferase (GST) SF

*E*-value threshold =  $1 \times 10^{-13}$ 1,568 nodes represent 13,000 sequences (50% ID filtered) Each node contains 1-930 sequences Median edge *E*-value =  $4 \times 10^{-25}$ Median alignment length = 210 residues Large nodes colored by SwissProt classification if >50% of sequences in each node belong to that class Gray: Not classified Triangles: PDB structure

#### Many new classes yet to be discovered





#### **Reaction types fail to track with similarity clusters**

#### Promiscuity likely more widespread than we know



### **Reaction families within a SF evolve at different rates**

unk.Agr

\*structurally characterized

AEE

possible AEE

unk.Cythu2

unknown

Enolase SF



#### "Pseudo-convergent" evolution of the same reaction from different intermediate ancestors in the SF tree



Unknown



Song et al, Nat Chem Biol, 3:486 (2007)



Sakai et al, Biochem, 48:2569 (2009)

## What sequences to align

## SFLD SF curation guided by sequence, structure, and functional information



## **Final thoughts**

Alignment methods with more sophisticated tools for choosing "representative" sequences to include in an MSA would contribute significantly to these difficult problems

> Especially important for both manual and automated curation communities

Automated function prediction experiments (CAFA) suggest the value of methods that allow for incorporation of functionally relevant features

- > Already included in some phylogenomic methods
- Similarity network methods could supplement current MSA approaches but need more rigorous development
- We can supply manually curated Gold Standard sets for evaluation of new methods

## Acknowledgments

Eyal Akiva, PhD Holly Atkinson, PhD\* Alan Barber, PhD\* Shoshana Brown, PhD Gemma Holliday, PhD Michael Hicks, PhD\* Florian Lauck\* Susan Mashiyama, PhD\* Elaine Meng, PhD **David Mischel** Rebecca Davidson Alexandra Schnoes, PhD\* Doug Stryke\* Jack Yu Jeff Yunes

#### **Collaborators**

USCF Resource for Biocomputing, Visualization & Informatics Tom Ferrin John Morris

*N6P Superfamily* Sarah O'Connor (Danforth Center)

> *Enolase Superfamily* John Gerlt (U of IL) Matt Jacobson (UCSF)

CytGST Superfamily Enzyme Function Initiative Steven Almo (Einstein) Richard Armstrong (Vanderbilt)

Alkaline phosphatase superfamily Dan Herschlag (Stanford) Jonathan Lassila

\*Former members

\$\$ NIH, NSF \$\$