# Sequence Similarity Networks for the Protein "Universe"

## John A. Gerlt

## University of Illinois, Urbana-Champaign

## Blue Waters Symposium
## May 13, 2014

# Personnel and Funding

## University of Illinois, Urbana-Champaign

**John A. Gerlt, PI**
**Victor Jongeneel, CoPI**
**Daniel Davidson, IGB**
**Boris Sadkhin, IGB**
**David Slater, IGB**
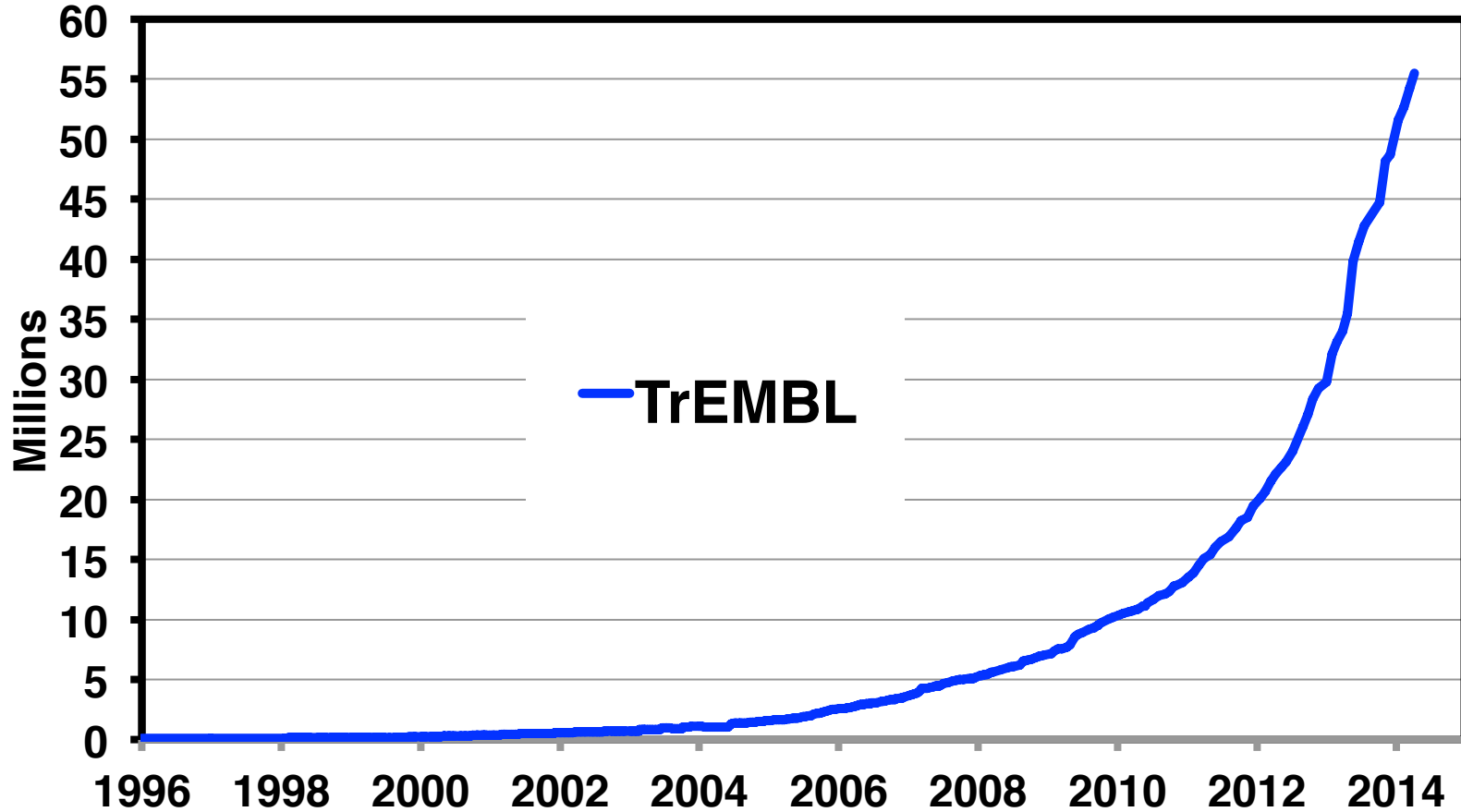
## External Collaborators

**Alex Bateman, EMBL-EBI**
**Matthew Jacobson, UCSF**

# The number of protein sequences is "exploding" !
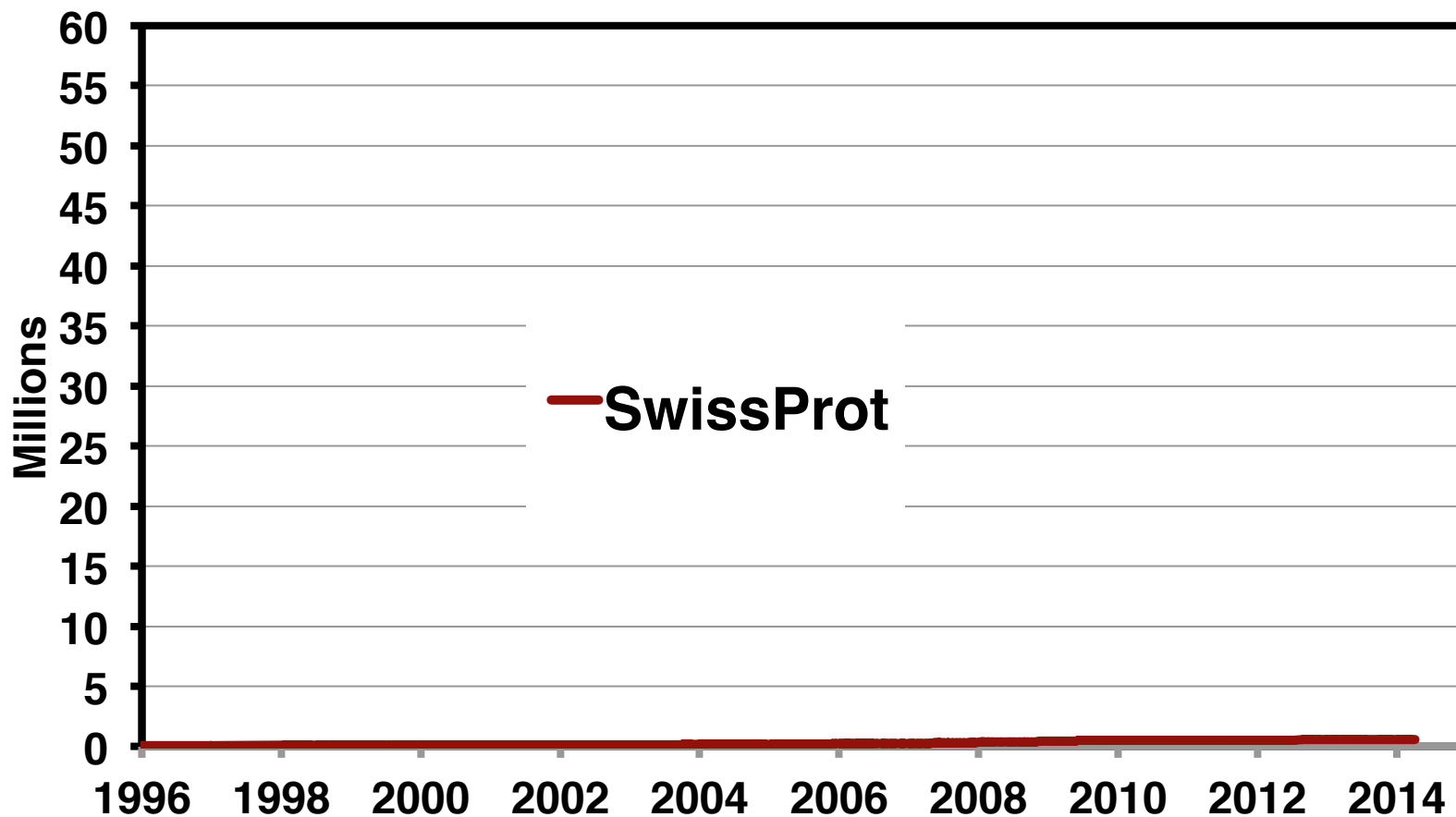
**Release 2014_04 of 16-Apr-2014 of UniProtKB/TrEMBL**

**contains 55,503,547 sequence entries.**

# But the number of curated annotations is lagging !
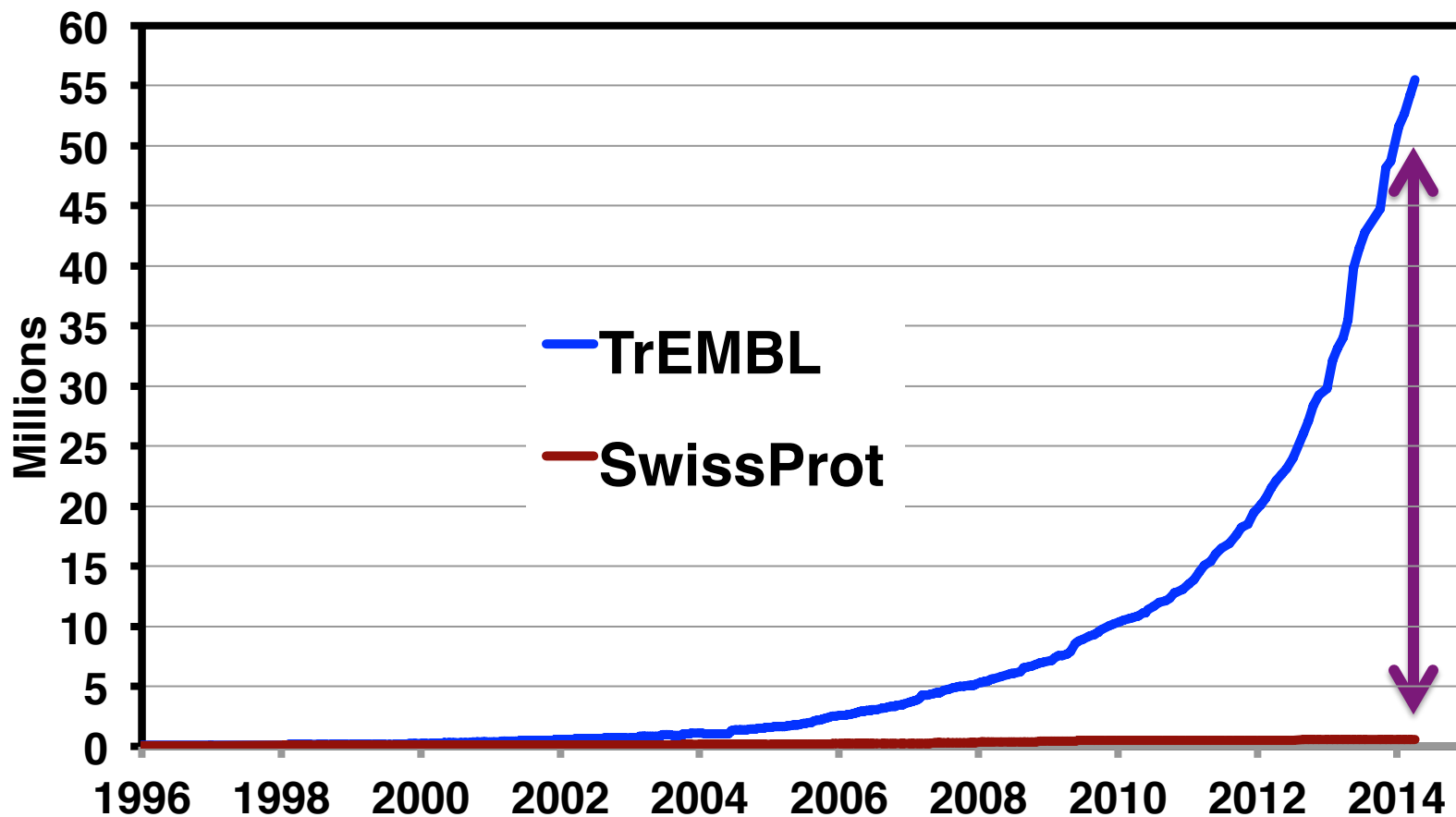
**Release 2014_04 of 16-Apr-2014 of UniProtKB/SwissProt**

**contains 544,996 sequence entries.**

# Perhaps 50% have unknown or uncertain functions
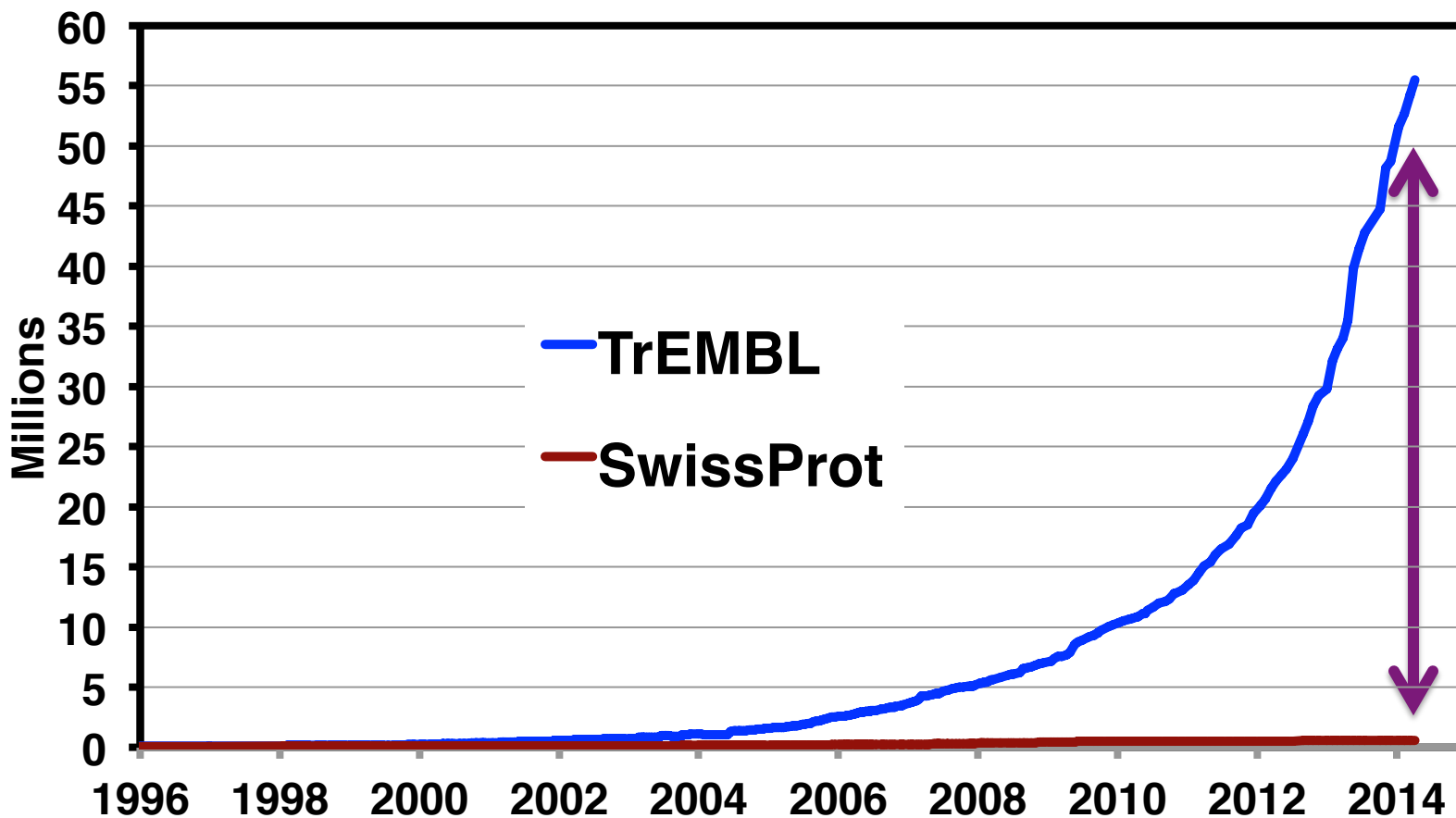
## Release 2014_04 of 16-Apr-2014 of UniProtKB

## contains 55,503,547 sequence entries.

# How do we solve this problem ?
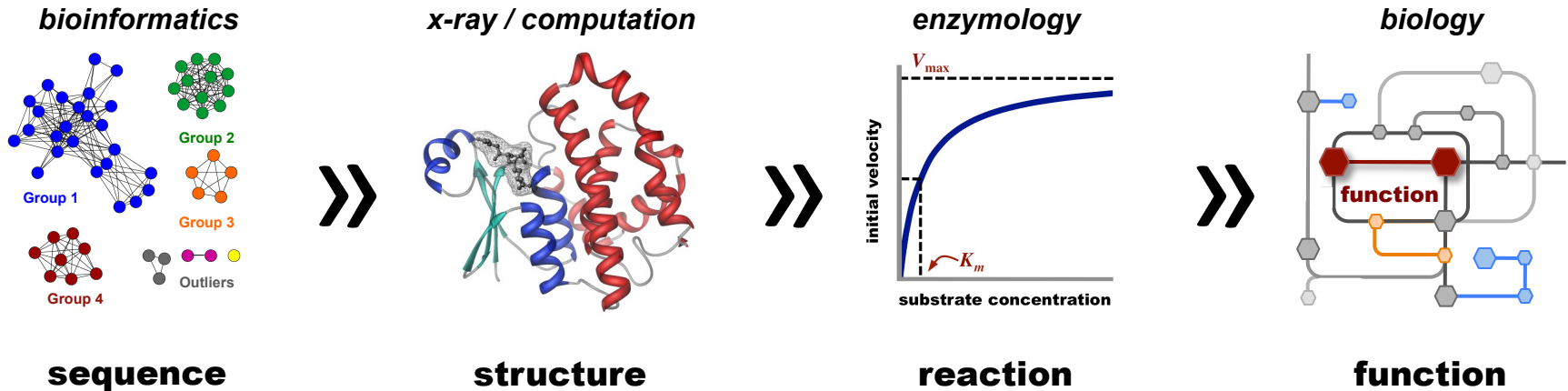
**Release 2014_04 of 16-Apr-2014 of UniProtKB**

**contains 55,503,547 sequence entries.**



**Requires high-throughput experimental
and computational strategies**

# U54 GM093342:  "Enzyme Function Initiative" (EFI)

*bioinformatics*          *x-ray / computation*          *enzymology*          *biology*



sequence                    structure                    reaction                    function

**Albert Einstein**
**Steven Almo**

**Boston University**
**Karen Allen**

**Gladstone Institutes**
**Katherine Pollard**

**University of Illinois**
**John Gerlt**
**John Cronan**
**Jonathan Sweedler**

**UCSF**
**Matthew Jacobson**
**Andrej Sali**
**Brian Shoichet**

**University of New Mexico**
**Debra Dunaway-Mariano**

**Pennsylvania State**
**Squire Booker**

**University of  Virginia**
**Wladek Minor**

**University of Utah**
**C. Dale Poulter**

# Deliverables, not Specific Aims

1.  **Develop robust high-throughput sequence/structure-based tools and strategies** to discover *in vitro* activities and *in vivo* metabolic functions of unknown enzymes

2.  **Disseminate tools to the community** for determining activities and functions of unknown enzymes

3.  **Collaborate with the community** to implement the tools and strategies

4.  **Correct annotations** in the protein databases

# Deliverables, not Specific Aims

1. **Develop robust high-throughput sequence/structure-based tools and strategies** to discover *in vitro* activities and *in vivo* metabolic functions of unknown enzymes

2. **Disseminate tools to the community** for determining activities and functions of unknown enzymes

3. **Collaborate with the community** to implement the tools and strategies

4. **Correct annotations** in the protein databases

# Sequence similarity networks:  Atkinson et al.

# Using Sequence Similarity Networks for Visualization of Relationships Across Diverse Protein Superfamilies

Holly J. Atkinson[1,2], John H. Morris[3], 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, California, United States of America, 2 Institute for Quantitative Biosciences, University of California San Francisco, San Francisco, California, 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.
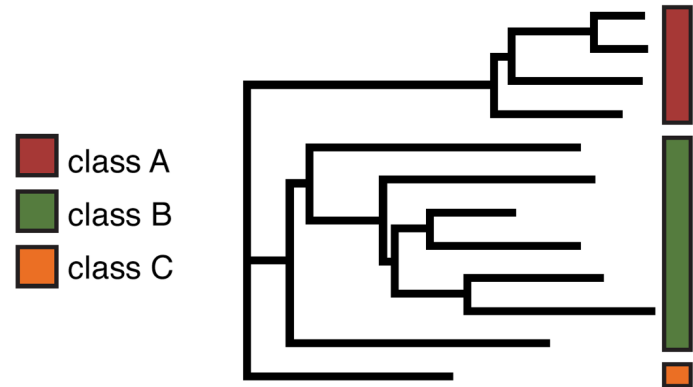
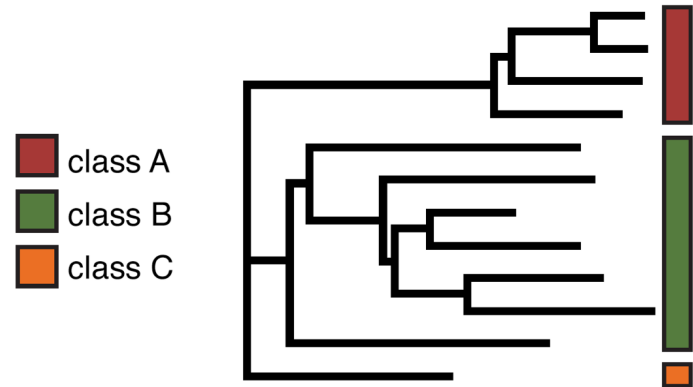# Sequence similarity networks (SSNs) vs dendrograms: enolase superfamily



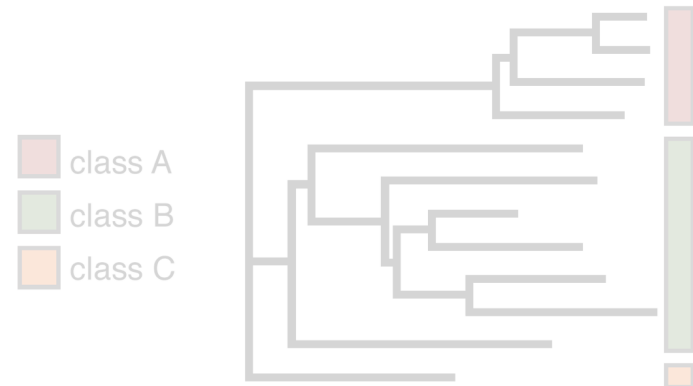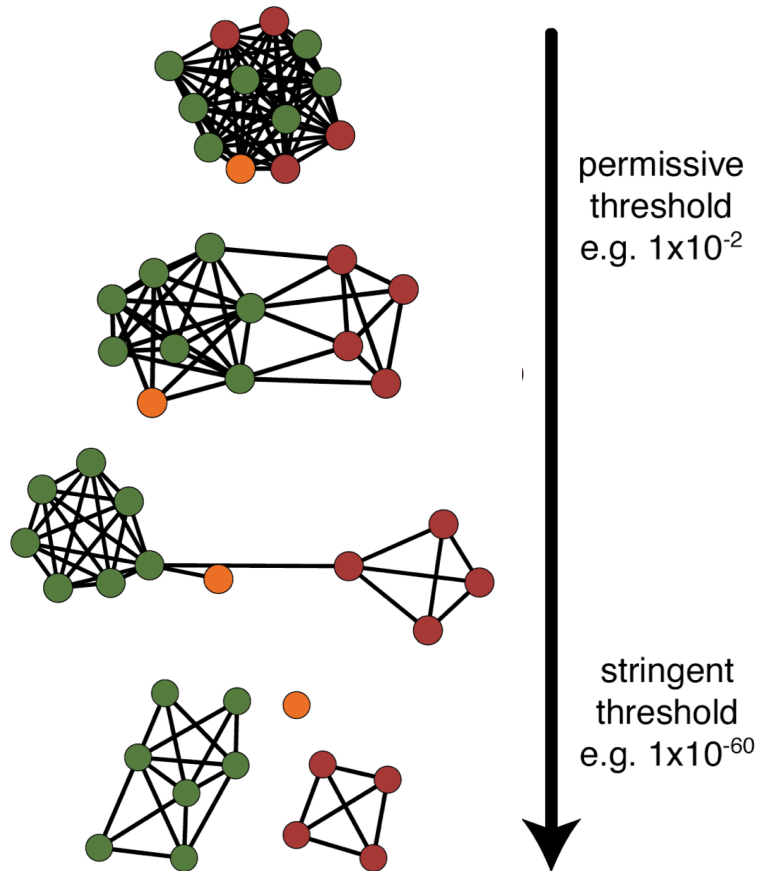**Families are easier to visualize in SSNs, so hypotheses are easier to formulate and explore**

# Dendrograms/trees for sequence relationships



class A
class B
class C

# Connectivity: multiple sequence alignments



class A
class B
class C

# Sequence similarity networks



permissive
threshold
e.g. $1 \times 10^{-2}$

stringent
threshold
e.g. $1 \times 10^{-60}$

class A
class B
class C

node (circle) = sequence

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

# Connectivity:  all-by-all BLASTP e-values



permissive
threshold
e.g. $1 \times 10^{-2}$

stringent
threshold
e.g. $1 \times 10^{-60}$

class A
class B
class C

**node (circle) = sequence**

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

# Faster to calculate than dendrograms



permissive
threshold
e.g. $1 \times 10^{-2}$

stringent
threshold
e.g. $1 \times 10^{-60}$

class A
class B
class C

**node (circle) = sequence**

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

# Qualitatively similar results



permissive
threshold
e.g. $1 \times 10^{-2}$

stringent
threshold
e.g. $1 \times 10^{-60}$

class A
class B
class C

node (circle) = sequence

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

# EFI Home Page



## http://enzymefunction.org/

# EFI Home Page

# EFI-EST: user-initiated sequence similarity networks

# Input sequence for BLAST or Pfam/InterPro family(ies)

# Output **full** and **rep node** networks (.xgmml files)

# Cytoscape 2.8.3: full network, e-110

**User does not have to wait for BLASTs**
**Expedite hypotheses and experiments**

**Pfam**

**Families**:  conserved protein families based a seed alignment of representative sequences that is used to generate a profile hidden Markov model (HMM).

**14,831 families in Pfam 27.0 (March 2013)**

**Clans**:  families (superfamilies) that have a common evolutionary ancestor based on structure and sequence.

**515 clans in Pfam 27.0 containing 4,563 Pfam families**

**http://pfam.sanger.ac.uk/**

# 515 clans (4,563 families): 68,545 sequences/clan

# 10,268 "clanless"-families: 1,909 sequences/family

# BLASTall on Blue Waters:
## 32 processors /node and 64 GB RAM/node

1. **Partition each family into smaller sets of query sequences, e.g., 100 sequences for small families**

2. **Load the entire set of sequences into RAM**

3. **Use BLASTall to calculate e-values for the query sets against all sequences**

4. **Store BLAST results**

5. **Concatentate results**

6. **Filter results to remove redundancy (A-B vs B-A)**

# BLASTall on Blue Waters:
## 32 processors /node and 64 GB RAM/node

**24 hr wall time limits the sizes of the query sets, depending on total number of sequences**

**Solution:  decrease sequences into smaller query sets.  Or, split family into smaller pieces.**

**64GB RAM limits the number of sequences that can be loaded and the number of results that can be stored**

**Solution:  decrease number of processors to make more RAM available per processor but this increases the number of required nodes.  Or, split family into smaller pieces. These are computationally equivalent.**

# 10,268 families:  average 1,909 sequences/family

**To date, 10,264 BLASTs are complete!**

**These are being processed to yield .xgmml files for Cytoscape.  Statistics are being calculated for choice of visualization thresholds.**

**514 BLASTs completed!**

**The successful BLASTs are being processed to yield .xgmml files for Cytoscape as well as statistics to choose visualization thresholds.**

**Families are being extracted from the 514 completed clans (3,049/4,563 to date); these are being processed to yield .xgmml files and statistics.**

**The largest clan (CL0023: 3,066,502 sequences) is too large for Blue Waters, using our current algorithms.**

**The .xgmml files for 514 clans and all 14,831 families will be made available via a Web server.**

**Alternative BLAST approaches will be developed for CL0023.**

**Networks to be calculated on a quarterly refresh cycle to provide current networks to the biological community.**

# Personnel and Funding

## University of Illinois, Urbana-Champaign

John A. Gerlt, PI
Victor Jongeneel, CoPI
Daniel Davidson, IGB
Boris Sadkhin, IGB
David Slater, IGB

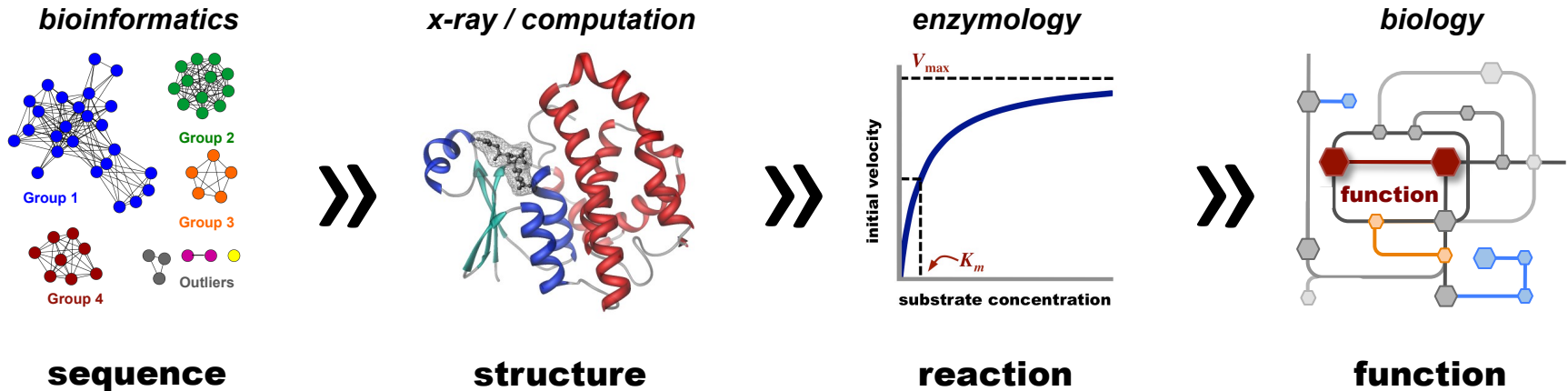## External Collaborators

Alex Bateman, EMBL-EBI
Matthew Jacobson, UCSF

# U54 GM093342: "Enzyme Function Initiative" (EFI)

*bioinformatics*        *x-ray / computation*        *enzymology*        *biology*



sequence                structure                reaction                function

**Albert Einstein**
**Steven Almo**

**University of Illinois**
**John Gerlt**
**John Cronan**
**Jonathan Sweedler**

**University of New Mexico**
**Debra Dunaway-Mariano**

**Boston University**
**Karen Allen**

**Pennsylvania State**
**Squire Booker**

**UCSF**
**Matthew Jacobson**
**Andrej Sali**
**Brian Shoichet**

**Gladstone Institutes**
**Katherine Pollard**

**University of Virginia**
**Wladek Minor**

**University of Utah**
**C. Dale Poulter**