

An Outlook into Ultrascale Visualization of Large-Scale Biological Data

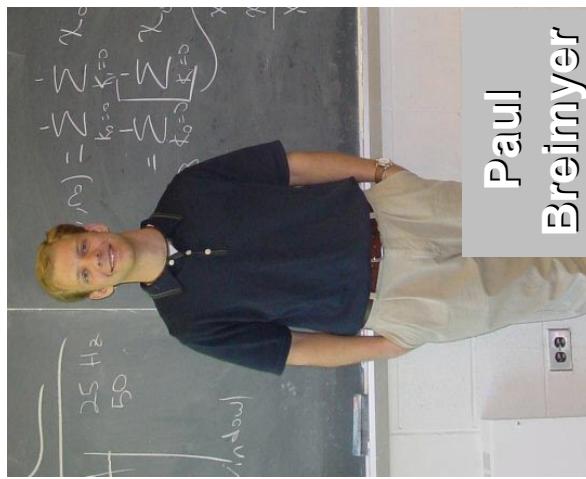
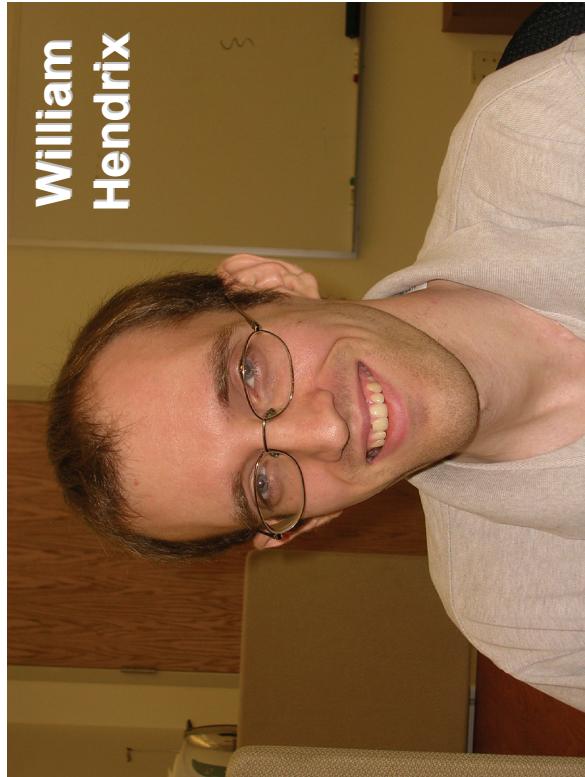
A Customer's Perspective

Dr. Nagiza F. Samatova

Department of Computer Science
North Carolina State University
and

Computer Science and Mathematics Division
Oak Ridge National Laboratory

Contributors



Outline

- Biology 101
- Motivational Problems
- Cross-Cutting Ultrascale Visualization Challenges
 - Top 6 Quests for Next Generation UV of Biological Data

Biology in a Nutshell – CS Perspective

For people with little knowledge but infinite intelligence

Genomes



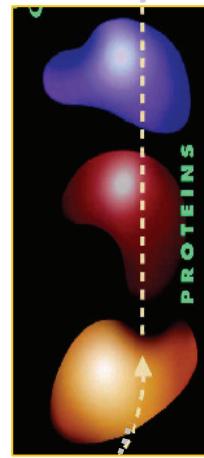
Genome (ROM): assembly code on how to build proteins

Instructions: A, C, T, G

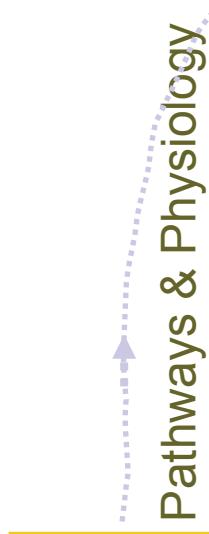
3 variables \Leftrightarrow amino acid

Genome consists of genes

Gene Products



Structure & Function



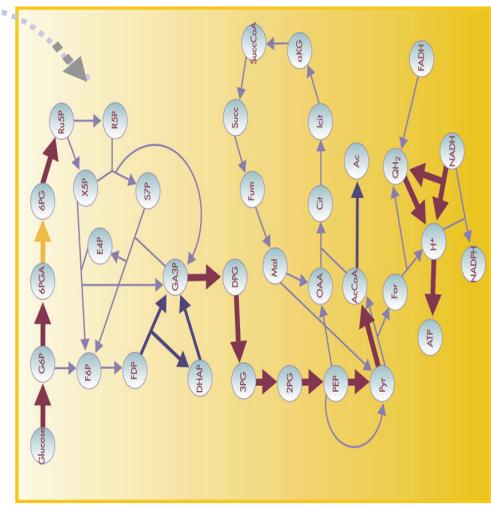
Gene \Leftrightarrow Protein:

Object description \Leftrightarrow Object instantiation

Protein: a sequence of amino acids

Protein \Leftrightarrow Functions

Enzymes: proteins that catalyze biochemical reactions



Pathways & Physiology

- Pathway:** a sequence of reactions
- Network** (directed graph): a set of pathways
(metabolites: nodes and enzymes: edges)

From Genes to Protein Functions

The first and most crucial step in systems biology

Guilt-by-Association goes global

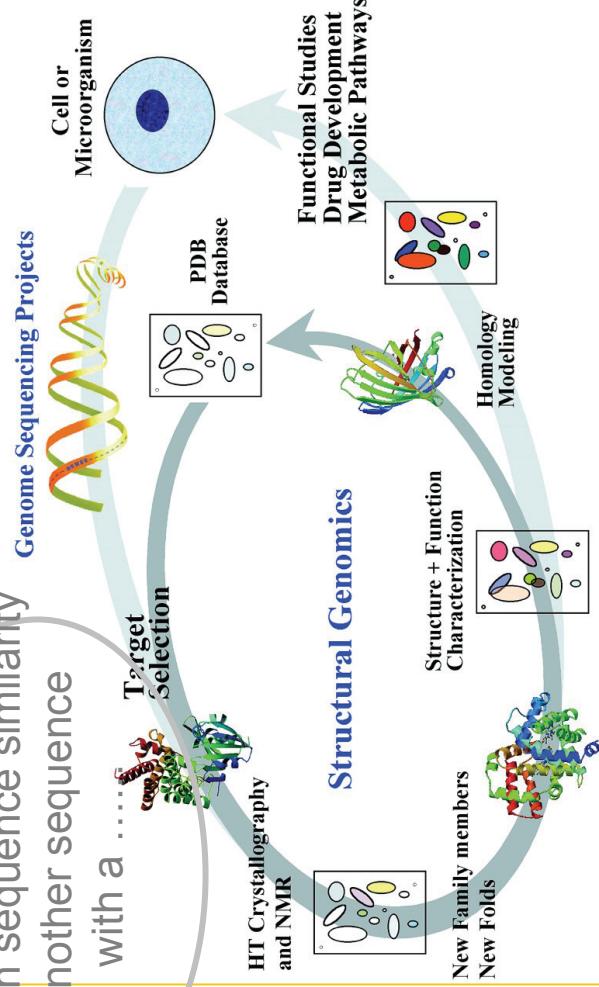
Function assigned
based on sequence similarity
to another sequence
with a

Function assigned
based on sequence similarity
to another sequence
with a

Function assigned
with a

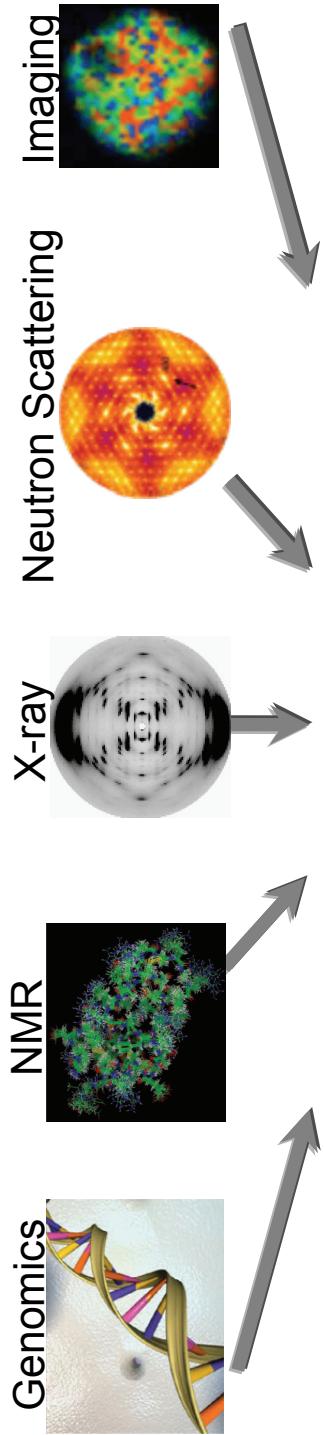
based on sequence similarity
to another sequence
with a

C I I G K G R S Y K R G T V S I T K S G I K C Q P W S S M I P E H S F L P S S Y
C K T G N G K N Y R G T M S K I K N G I T C Q R W S S T S P H R P R F S P A T H
C Y H G D G Q S Y R G I S S T I T T G K K C Q S W S S M T P H R Q K T P E N Y
C A E G V G M N Y R G N V S V I R S G I E C O I W R S R Y P H K P E I N S T T H
C G G Y r G t S T G i C Q W S S P H P

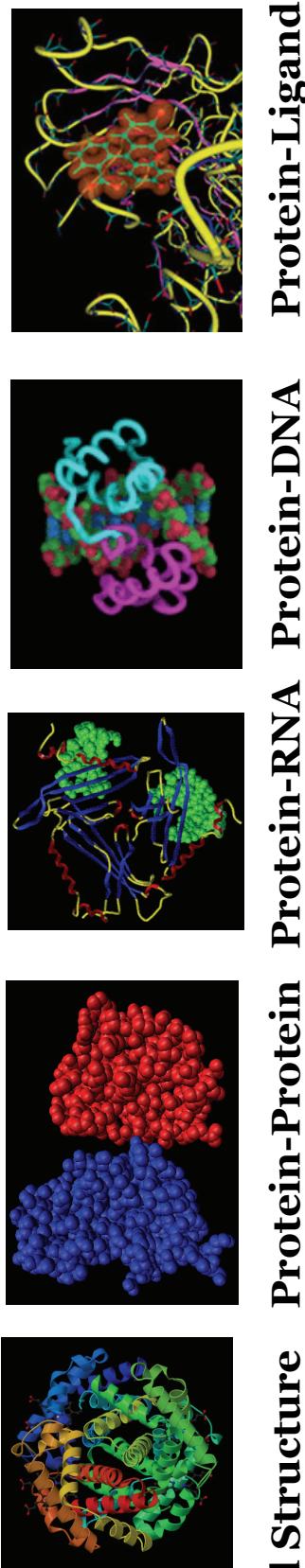


Proteins Function Interactively

What we observe



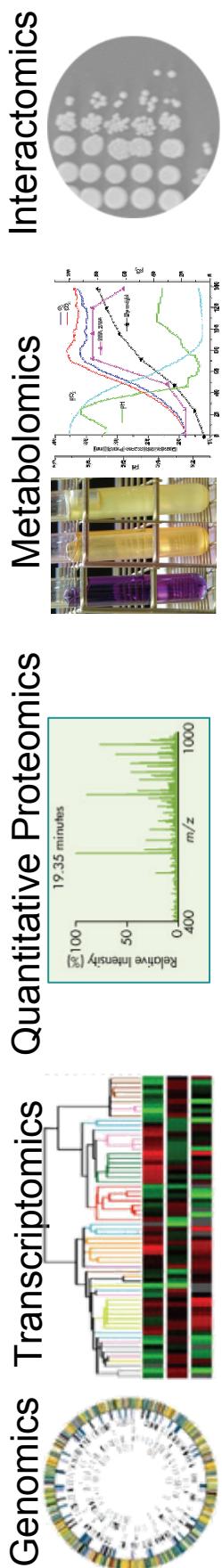
Protein Machines



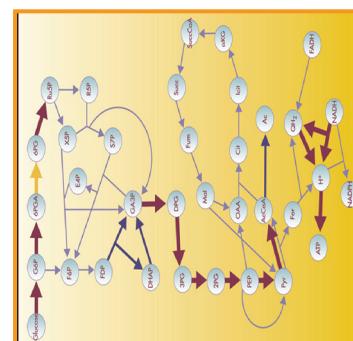
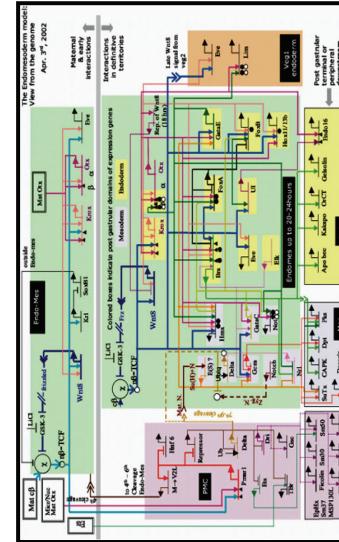
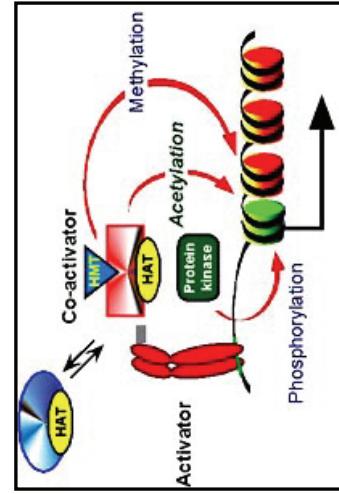
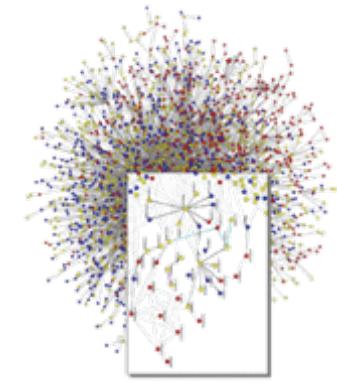
What we want to derive

Proteins Function in Pathways/Networks

What we observe



Networks/Pathways



Metabolic

Regulatory

Signaling

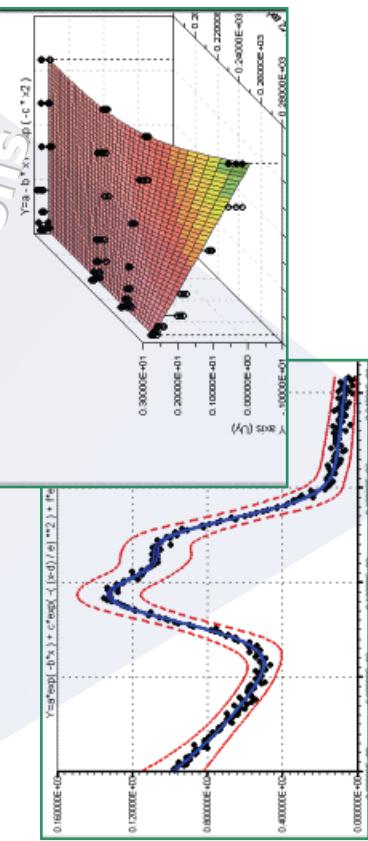
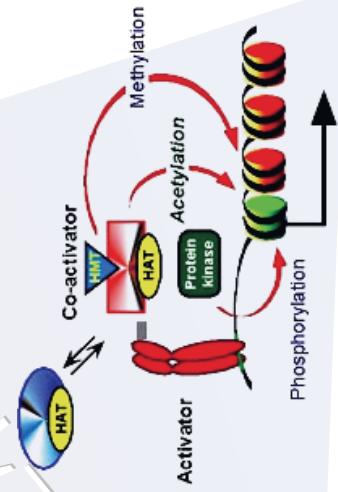
Protein Interaction

What we want to derive

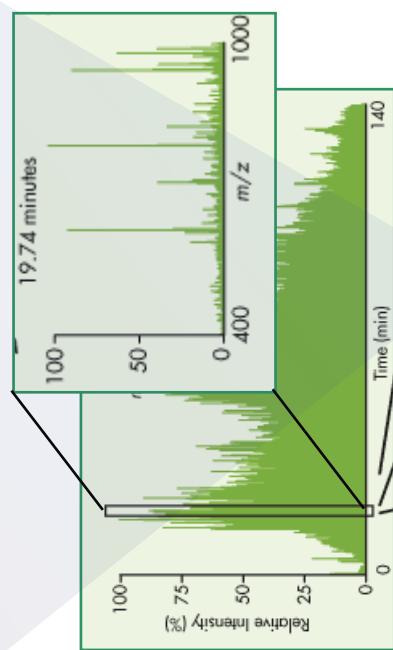
It is not just the Size – but the Complexity



Large-Scale Data



Non-linear correlations

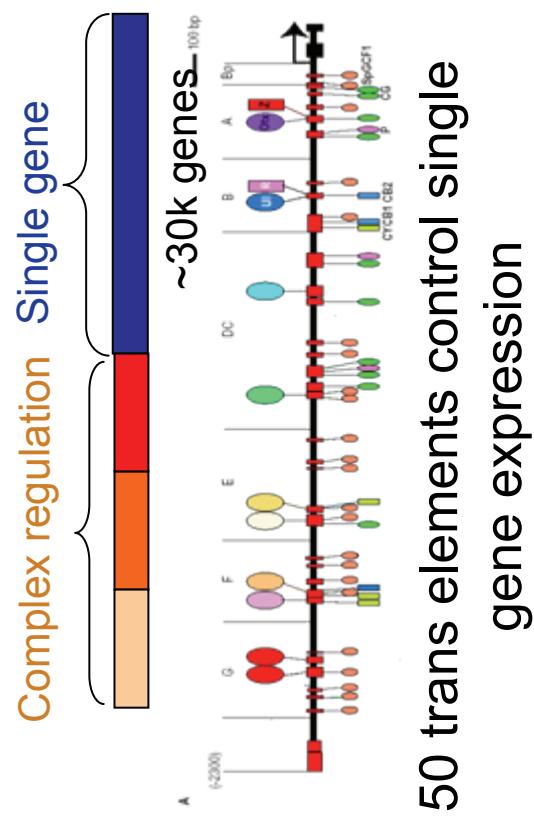


Noisy

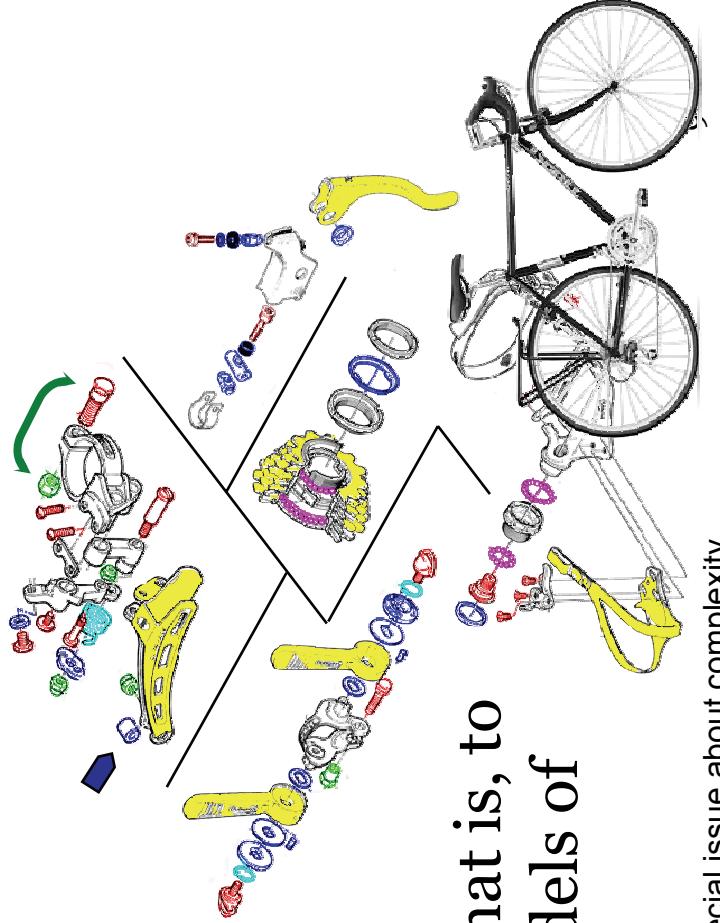
Data Describes Complex Natural Phenomena

How to untangle riddles of BioComplexity

Complexity of biological systems comes from interconnections.



Biological systems are complex because of non-linear coupling of their structural, genotypic, & phenotypic properties.



Challenge:

How to “connect the dots”, that is, to construct predictive *in silico* models of these biological systems.

How Can the Viz Community Help?

Outline

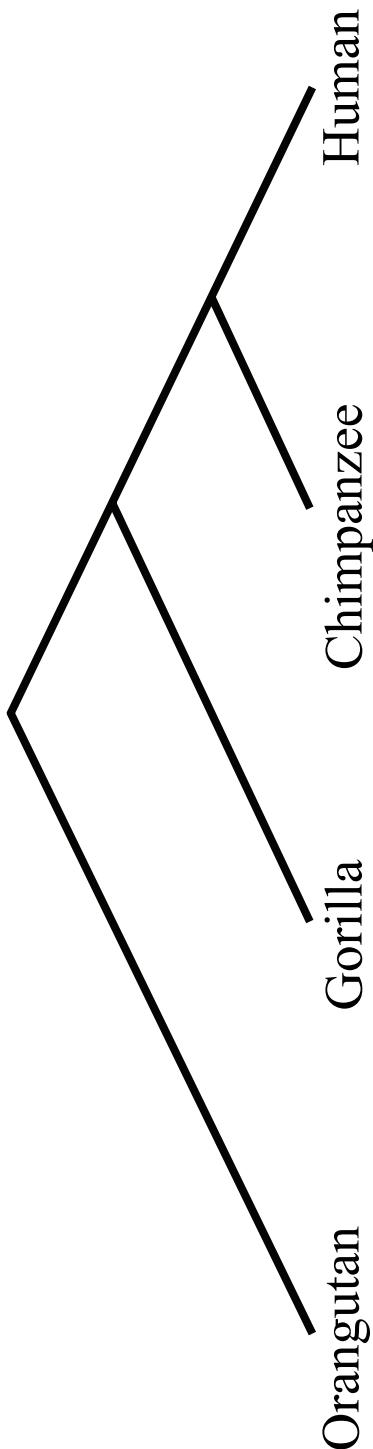
- **Biology 101**
- **Motivational Problems**
- **Cross-Cutting Ultrascale Visualization Challenges**
 - Top 6 Quests for Next Generation UV of Biological Data

Reconstruction of the Tree of Life

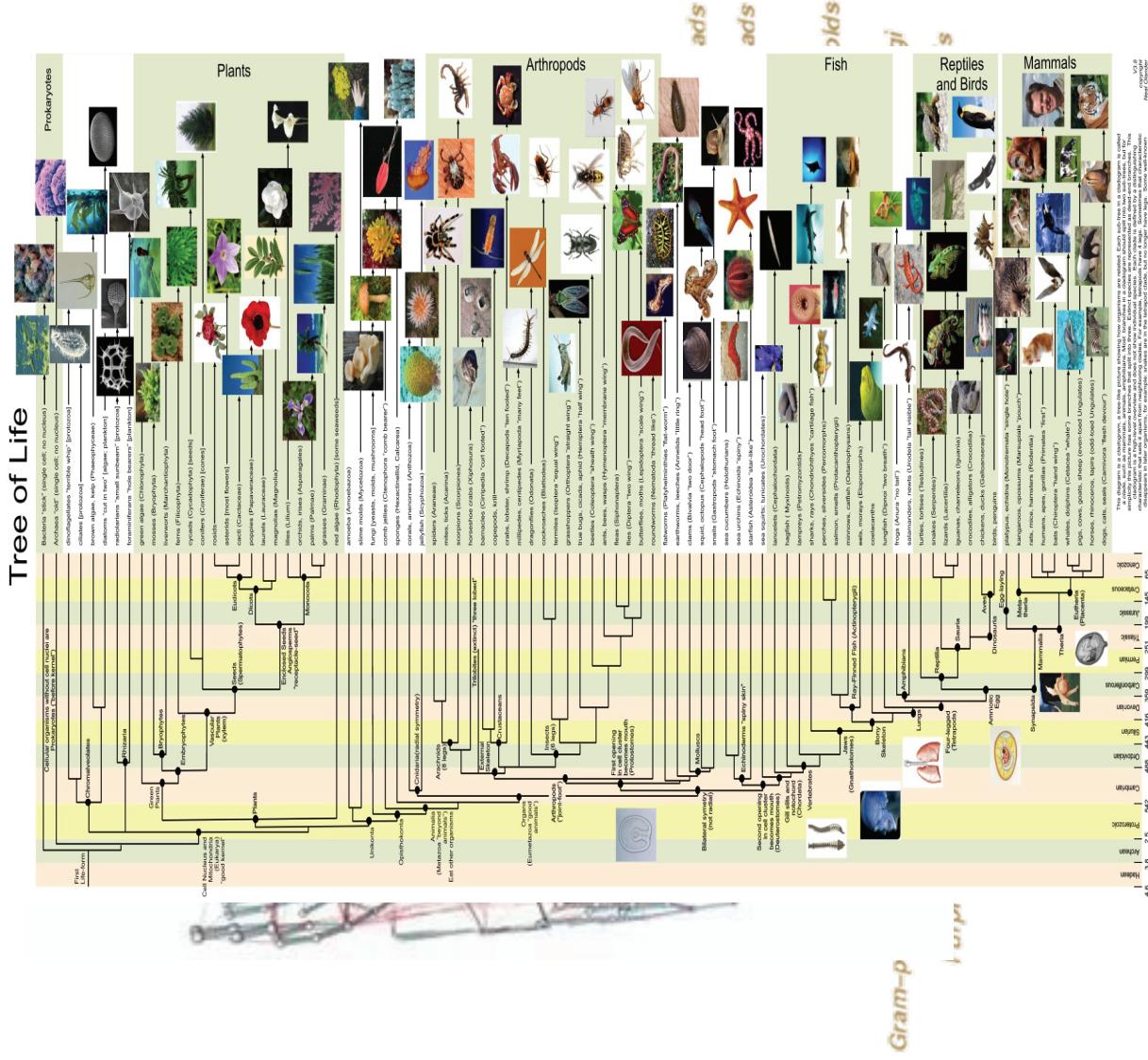
Ultrascale Visualization Problem

Phylogeny

A **phylogeny** is a tree representation for the evolutionary history relating the species of interest.



The Tree of Life for 10-100 Million Organisms



Ultrascaling Challenge

Computing the optimal phylogenetic tree based on the entire genome of 10 species will remain intractable even with peta-scale computers

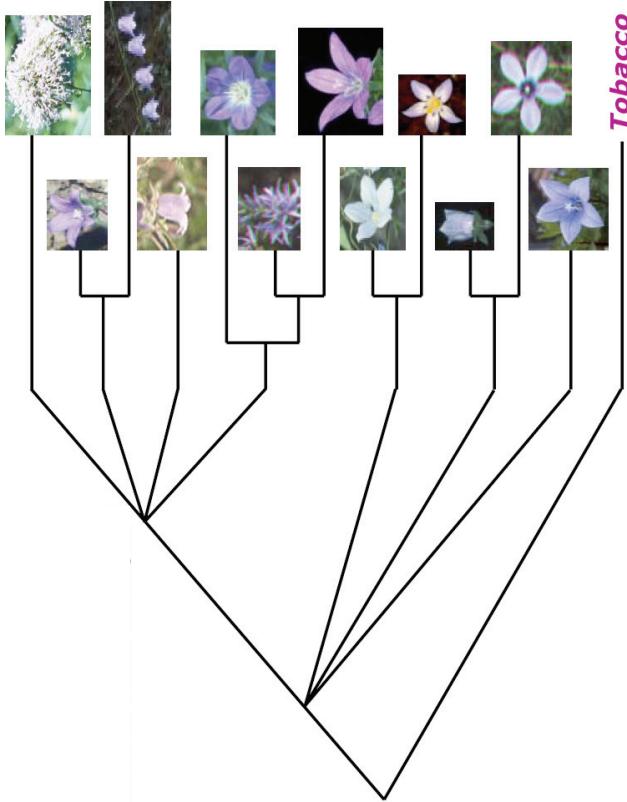
Complexity

Maximum likelihood

$C \sim n^2 * m$ where

n is number of contemporary species

m is size of the genome in question



Requirements

For 10^6 species on a gene sequence of about 1000 pairs or 1000 species on the entire genome:

Runtime: 1-3h at 1 Pflop

Memory: ~3 TB

Dealing with Computational Intractability

Computational intractability drives for various search heuristics to navigate a small fraction of exponentially-sized tree space in practical time.

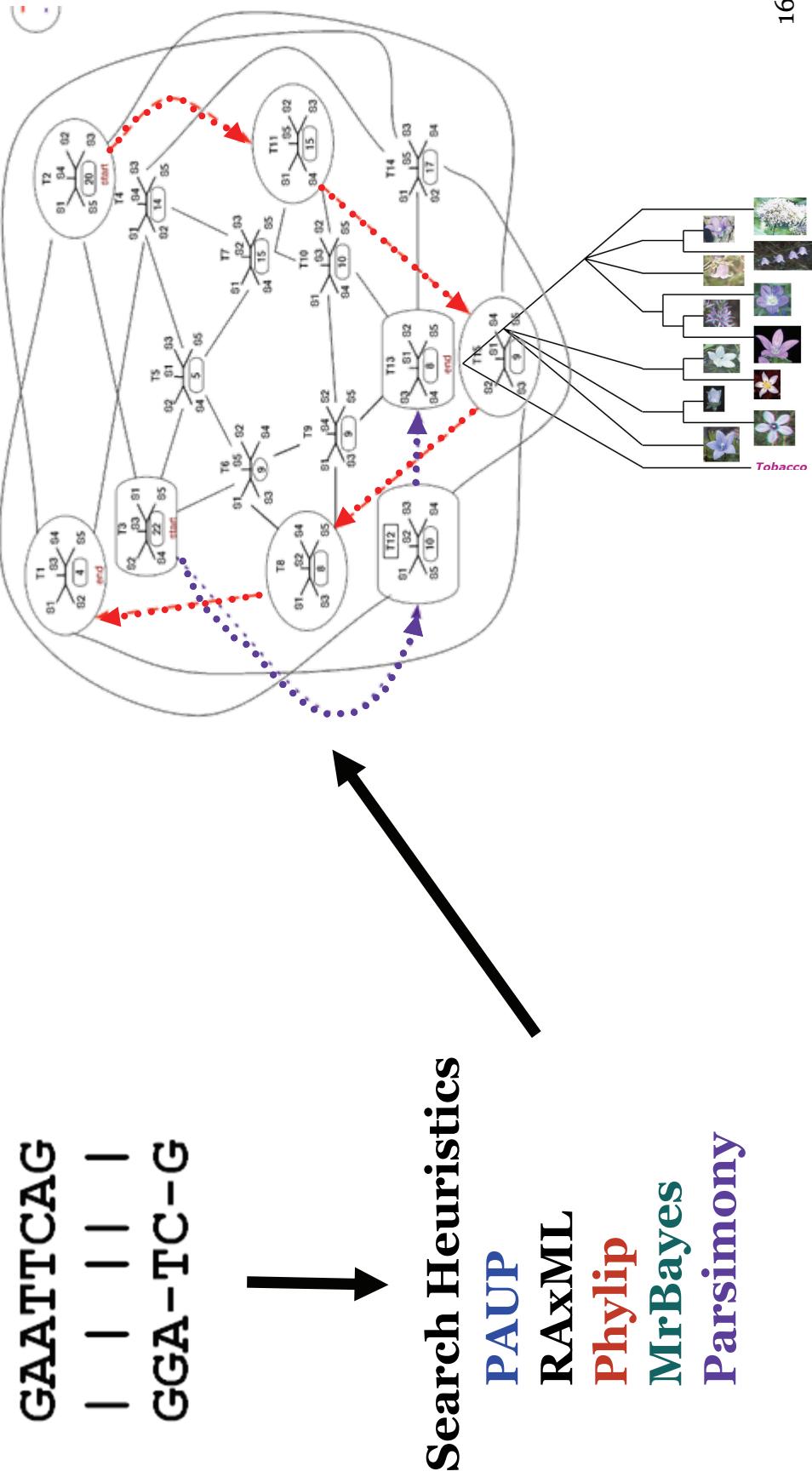
Sequence Alignment

GAATTCA

| | | | |

GGA-TC-G

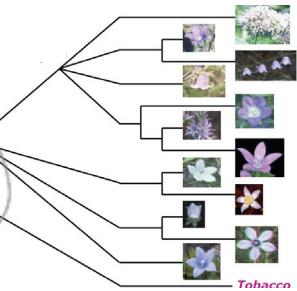
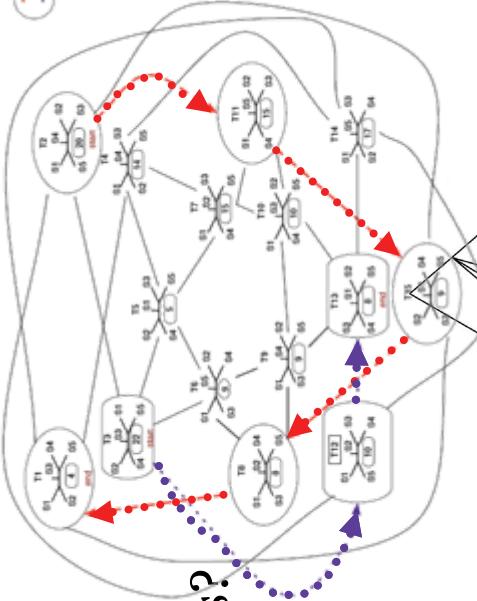
Search Histories



Visual Exploration of the Search History

Ultrascale Visualization Questions:

- How to visualize the landscape of local optimality (search histories)?
- How to visually compare the search histories from different heuristics?
- How to visually align two/many trees?
- How to visualize the hierarchical clusters of trees?



Impact

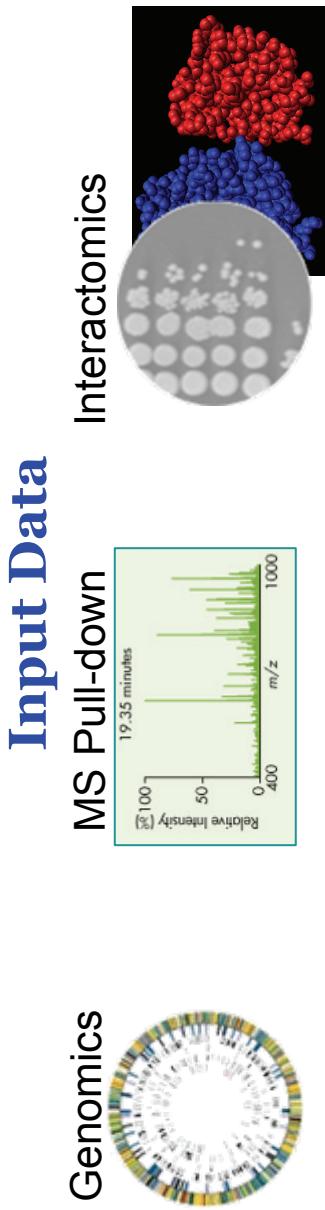
Design of better heuristics
More accurate reconstruction of phylogenies

Comparative Analysis of Networks

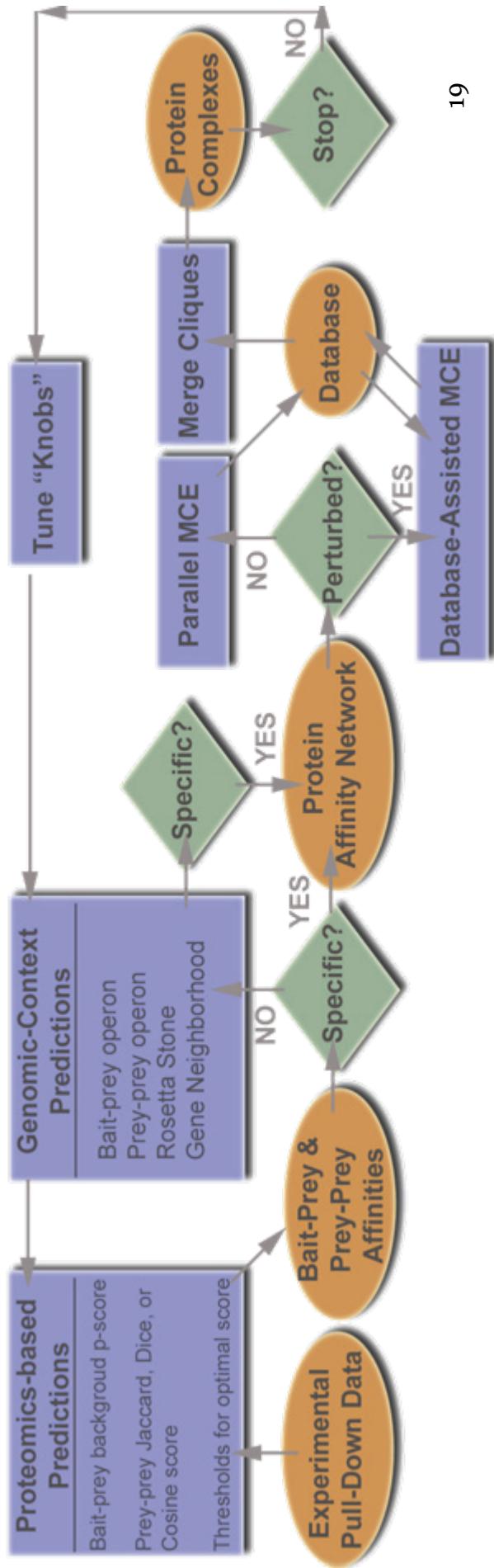
Ultrascale Visualization Problem

What is Genome-scale Interactome?

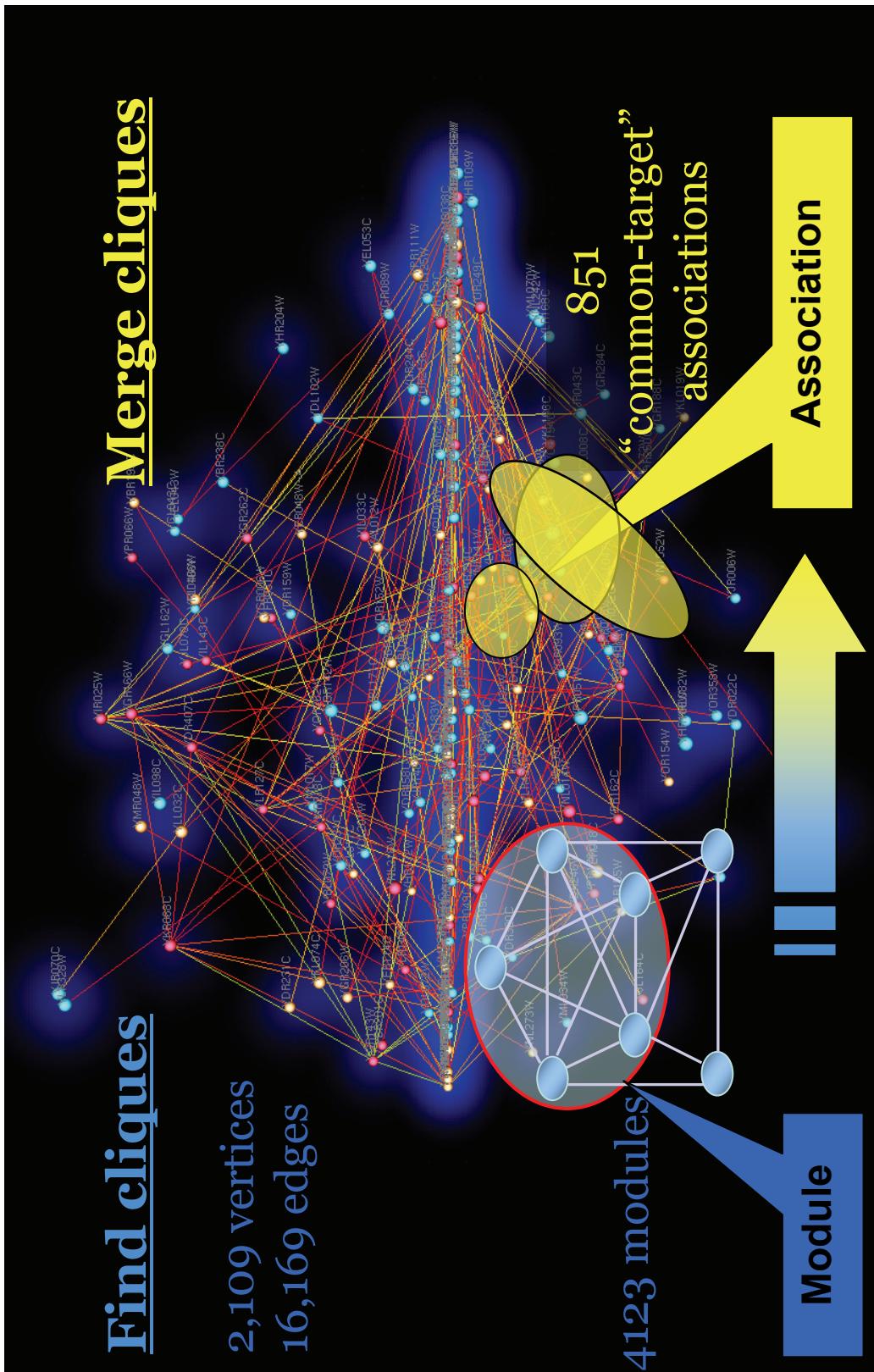
Integrated structural, genomic, and MS pull-down experimental clues for inference of genome-scale Interactome.



Network Inference Pipeline



Visual Exploration of Genome-Scale Interactome?

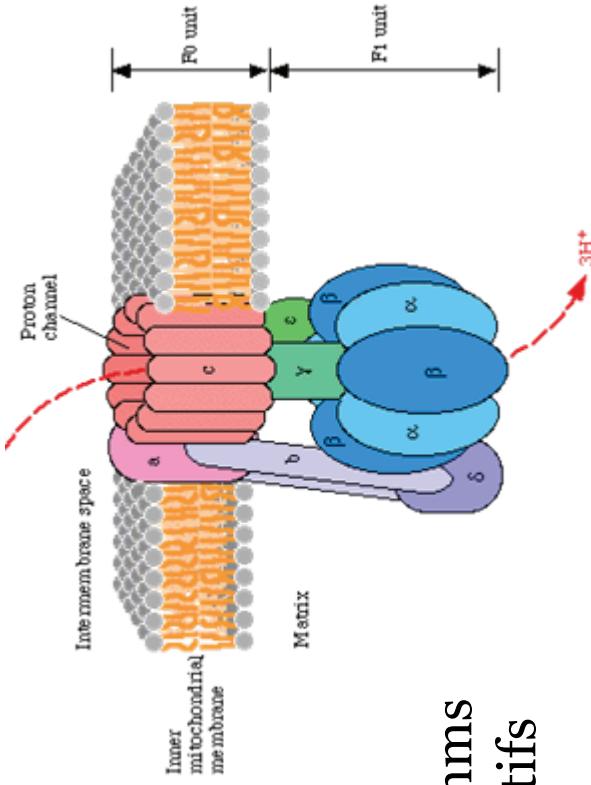


Visualization provided by Dr. Jian Huang and Mr. Joshe New, SciDAC Ultrascale Viz. Institute

Visual Exploration of Networks

Ultrascale Visualization Questions:

- What are highly connected (dense) network motifs?
- Are the network motifs **statistically significant/biologically relevant**?
- How do network motifs change for different parameters (knobs)?
- Is a network motif of interest (e.g. ATP synthase) present and what is it connected to?



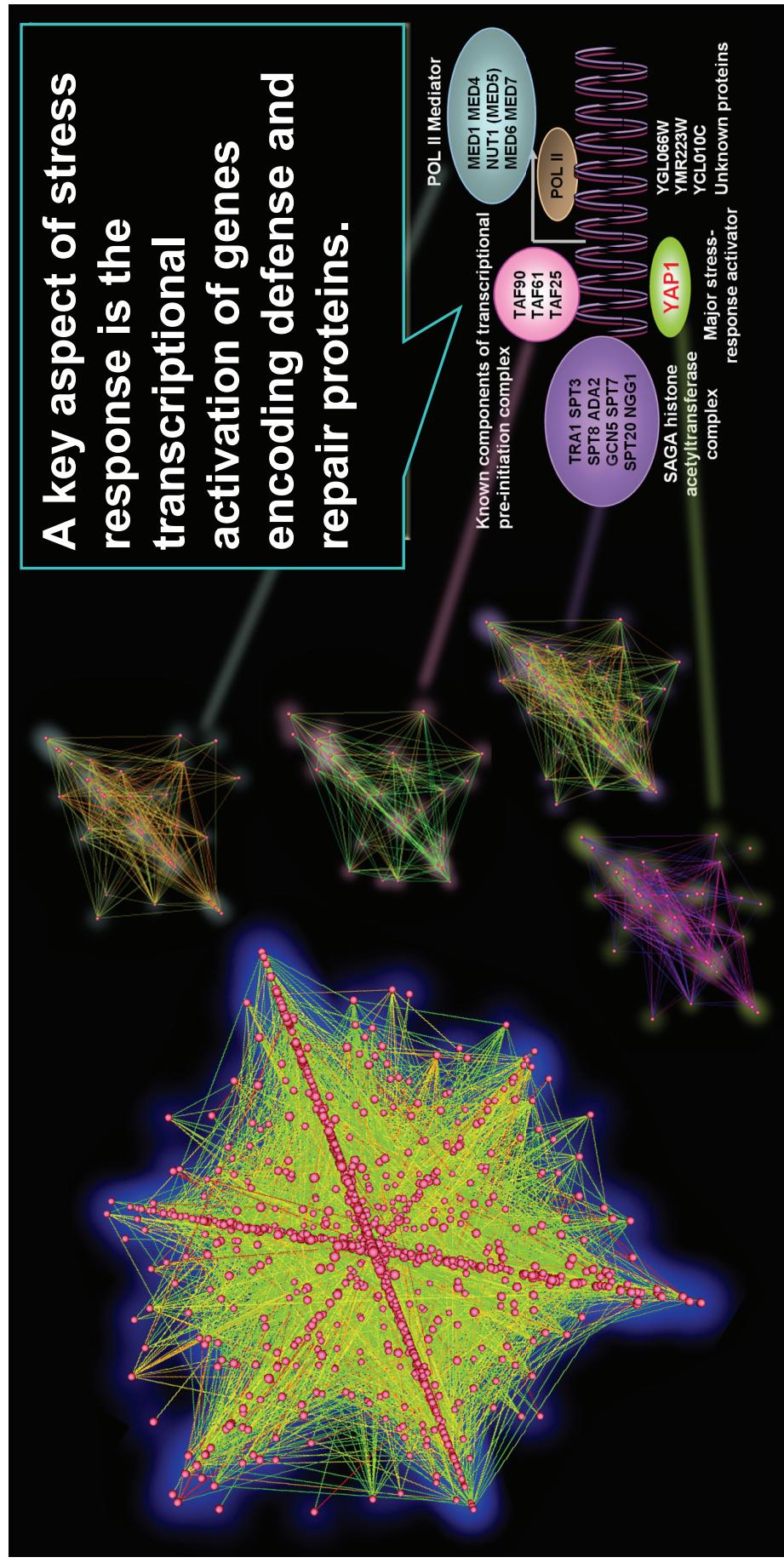
Impact

Design of better network inference algorithms
Discover biologically relevant network motifs

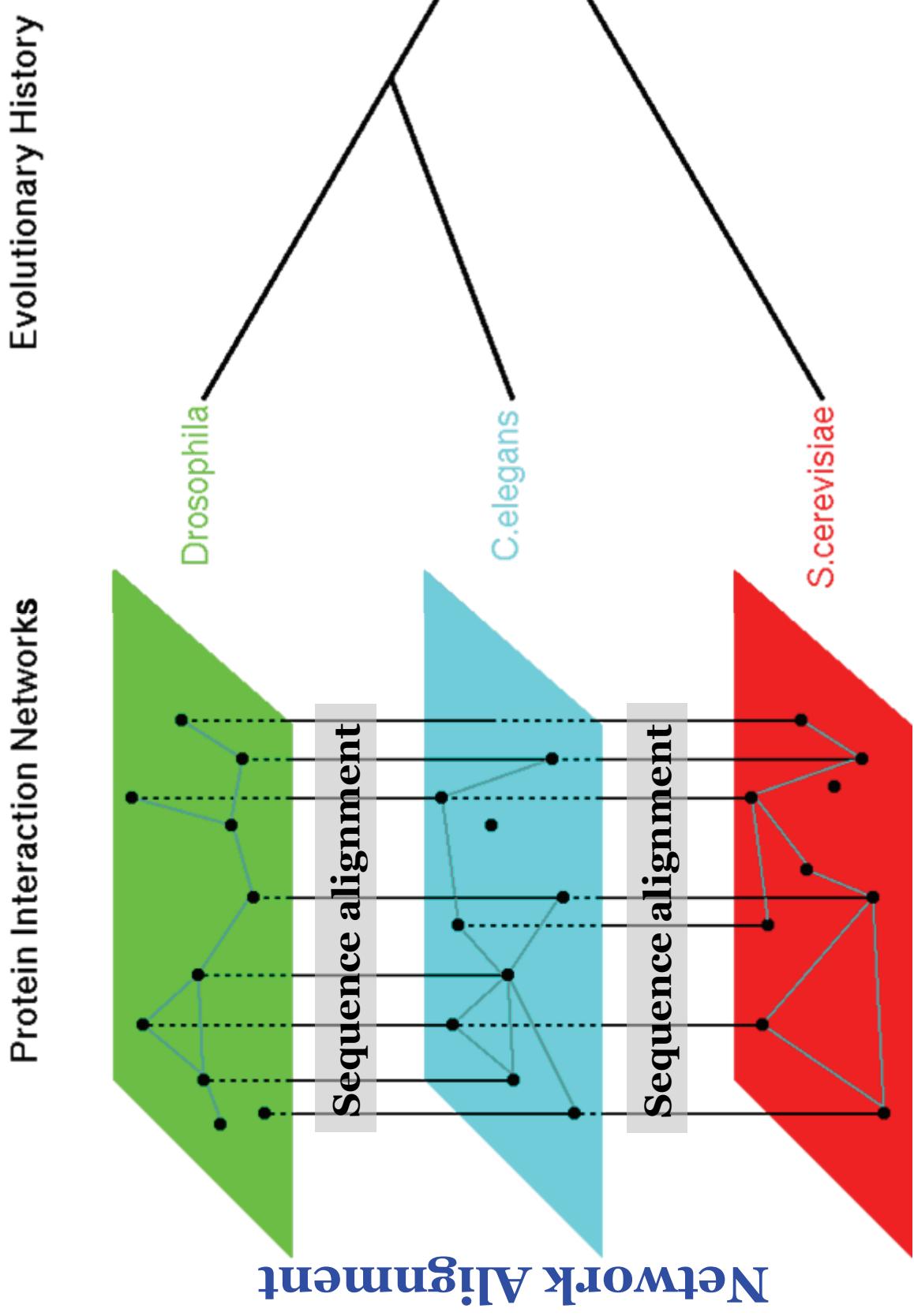
Prior Knowledge: ATP synthase

What cellular machineries are responsible for organismal resistance to stress?

Thermochemical pre-treatment of biomass and acidic compounds of sugar fermentation to ethanol require organisms resistant to stress.



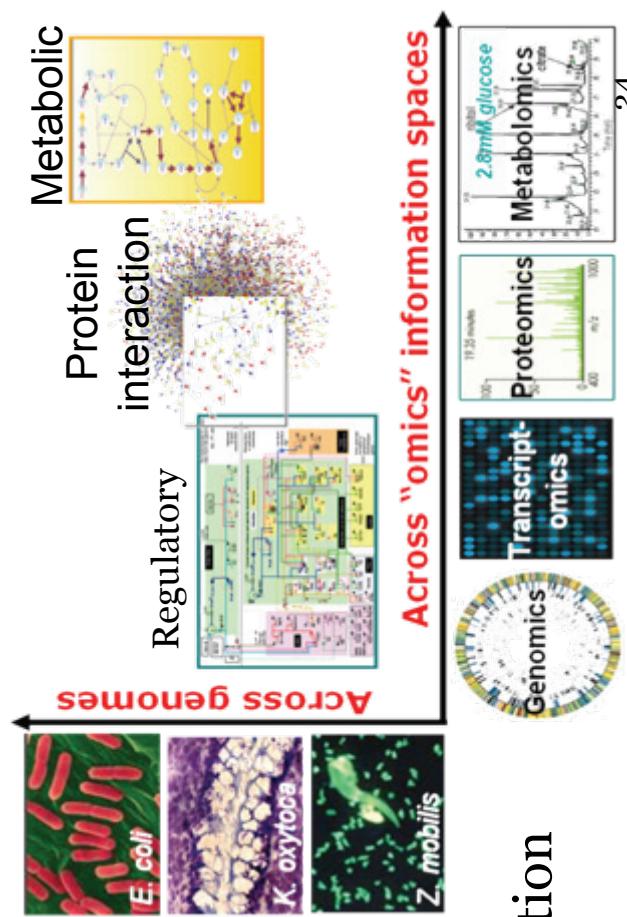
Are Network Motifs Evolutionarily Conserved?



Visual Exploration of Networks Evolution

Ultrascale Visualization Questions:

- What network motifs are evolutionary conserved?
- Is the conservation statistically significant (compared to random networks)?
- Is a network motif of interest evolutionary conserved? Across what organisms? Are these organisms evolutionary close or distant?
- How to visually compare networks across organisms and “omics” information spaces?



Impact

- Design better network analysis
- Discover novel network motifs
- Annotate proteins with unknown function

Identification of Phenotype-Specific Genes

Ultrascale Visualization Problem

What Genes are Responsible for a Phenotype?

Phenotypes:

Resistance to:

- High/low temperature
- Low/high pH
- High EtH concentration

Growth:

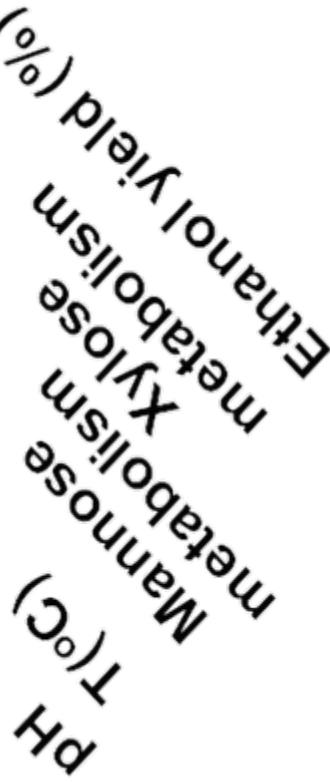
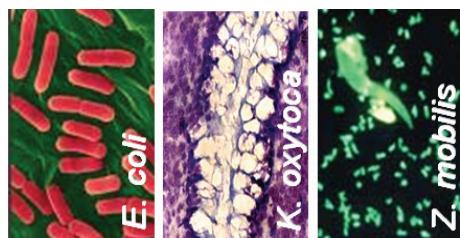
- Aerobic
- Anaerobic

Metabolism & Productivity:

- Ferment multiple sugars
- High EtH yield

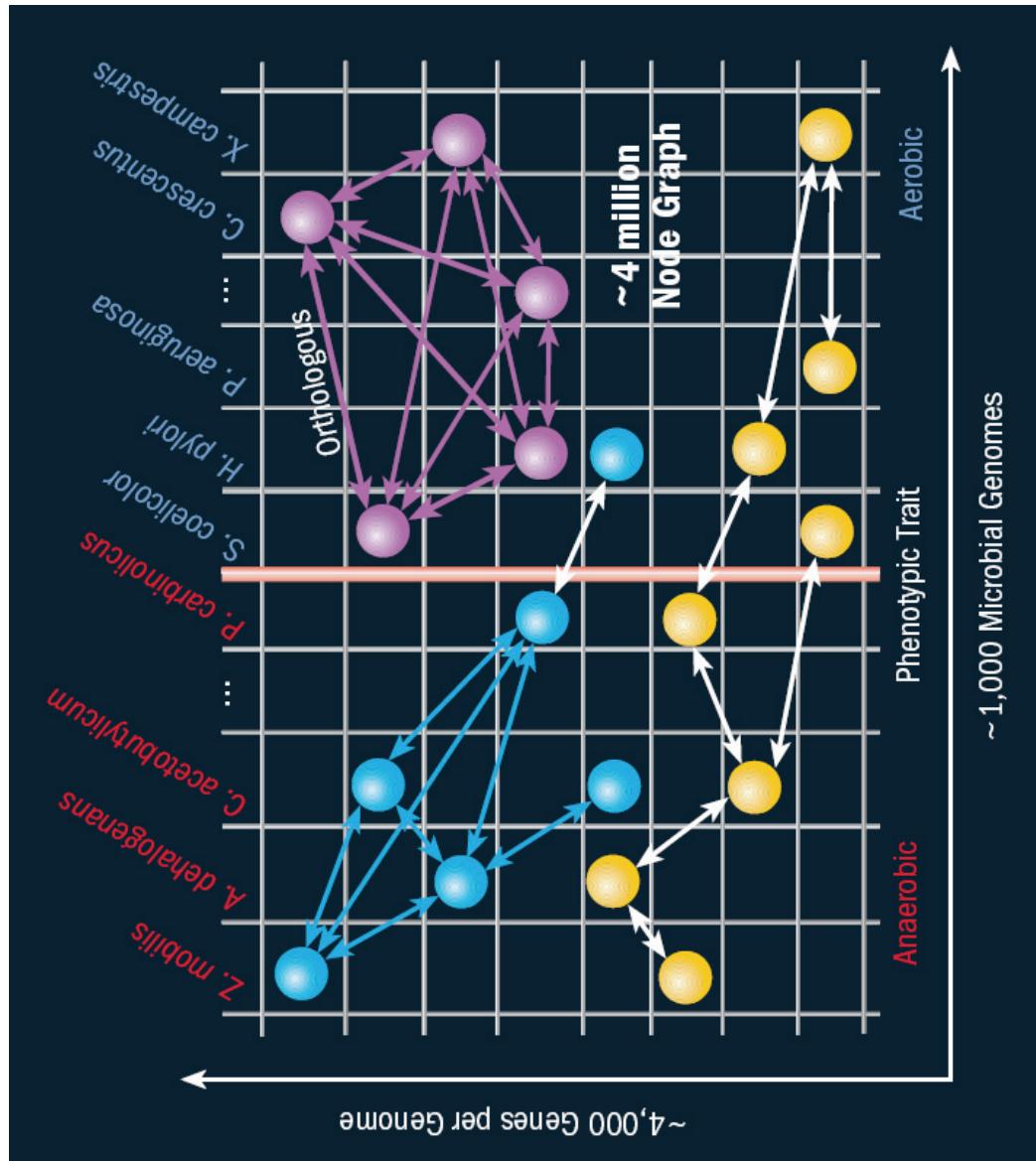
Phenotypic Traits:

	6.5	35	+	+	89
	5.5	30	+	+	95
	5.5	30	-	+	84



Identification of Phenotype-Related Genes

If a gene is responsible for a phenotypic trait, then it is **evolutionarily conserved** across several organisms.



Multi-partite
graph

Genotype-Phenotype Relationships

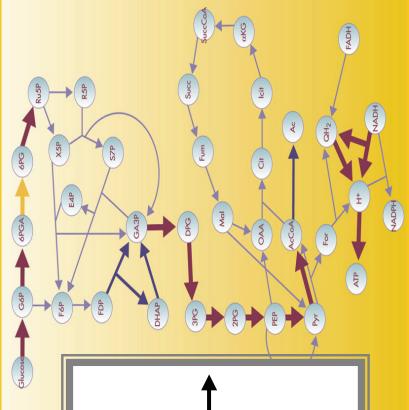
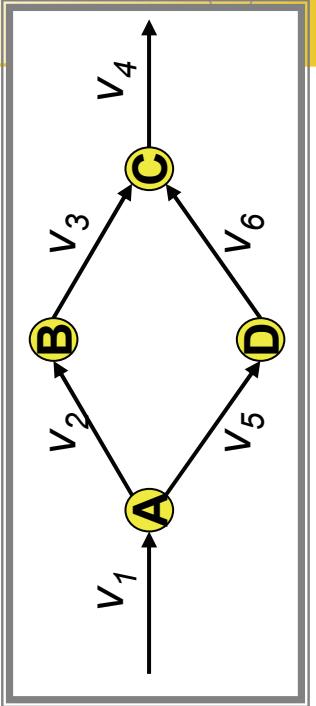
Ultrascale Visualization Problem

Mathematical Laws of Metabolism

Chemical reaction:



Reaction Network:



Stoichiometry Matrix, S :

Flux Vector, V :

		metabolites	reactions				
metabolites	reactions	A	B	C	D	E	H
V_1	V_1	a	$-c$	0	0	0	0
V_2	V_2	$-a$	b	0	0	0	0
V_3	V_3	0	$-b$	c	d	0	0
V_4	V_4	0	0	$-d$	e	f	0
V_5	V_5	0	0	0	$-f$	g	0
V_6	V_6	0	0	0	0	$-g$	h

**Steady State/
Mass Balances:**

$$\mathbf{A}: V_1 - V_2 - V_5 = 0$$

$$\mathbf{B}: V_2 - V_3 = 0$$

$$\mathbf{C}: V_3 - V_4 + V_6 = 0$$

$$\mathbf{D}: V_5 - V_6 = 0$$

**Steady state flux
vectors are in the
Null Space of S**

$$\mathbf{S} \cdot \mathbf{V} = \mathbf{0}$$

High Dimensional Space of System Phenotypes

Metabolism could be studied within **convex analysis context.**
Extreme pathways are a basis for a metabolic genotype.
Their regulation defines all possible metabolic phenotypes.

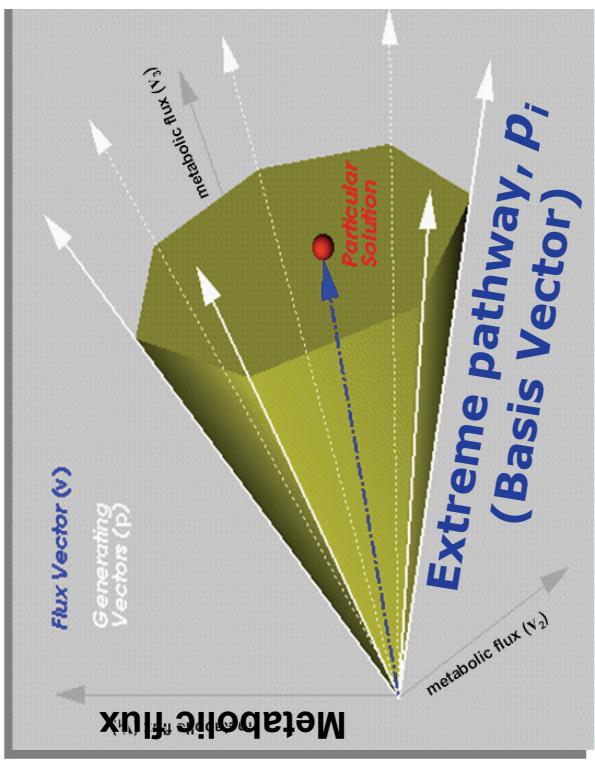
Convex polyhedral cone, C

$$C = \{ v \in R^n \mid v = \sum_{i=1}^k \alpha_i p_i, \alpha_i \geq 0 \}$$

Convex Analysis

Cellular Biology

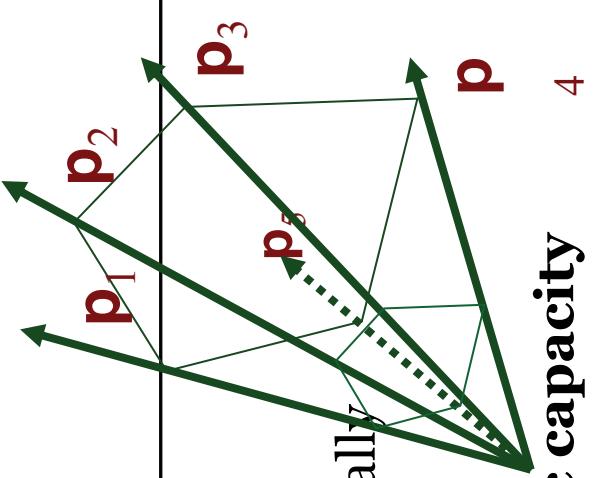
Capabilities of a Metabolic Genotype
Independent Extreme Pathways
Metabolic Phenotype
Positive Combination of Extreme Pathways



Every phenotype that the system can exhibit is a combination of these extreme pathways, which are then turned on or off.

The UV Challenges

- The number of extreme pathways can exceed the dimensions of the cone (i.e. linearly dependent):
 - How to visualize the pathways (potentially exponentially many)?



- The volume of the cone is a measure of metabolic capacity of the organism:

- Volume computation for a d -polytopes is an exponential problems, ($O(n^d)$)
- How to visually represent the volume?
- Is metabolic potential of an organism is rich or poor?

- If the objective function is known, then it is a linear or non-linear programming problem:

- Organism may pursue multiple objective functions
- How to solve the inverse problem? – How to find the objective functions based on the observed phenotypes?
- Can visualization help with this problem?

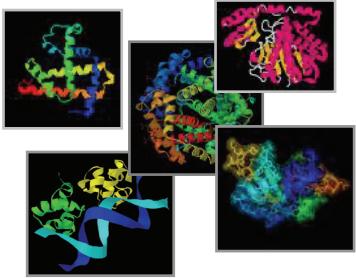
Outline

- Biology 101
- Motivational Problems
- Cross-Cutting Ultrascale Visualization Challenges
 - Top 6 Quests for Next Generation UV of Biological Data

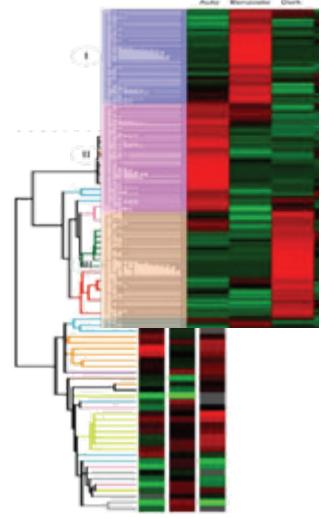
Q1: Comparative Ultrayscale “Visual-omics”

Genome
Proteome
Transcriptom
Physiome
Metabolome
Phenome
Morphome
Interactome
Glycome
Secretome
Ribonome
Orfeome
Regulome
Cellome
Operome
Transportome
Functome
Translatome
Pseudome
Foldome

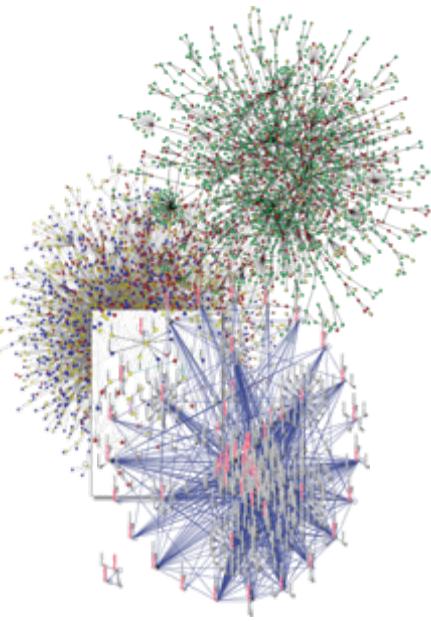
Comparative Proteomics



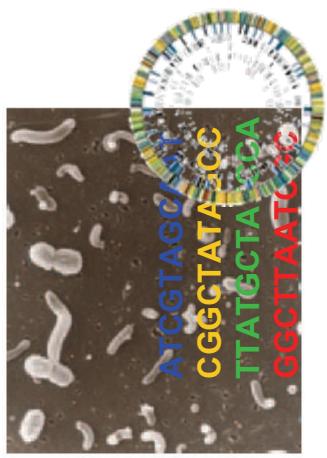
Comparative Transcriptomics



Comparative Interactomics



Comparative Metagenomics

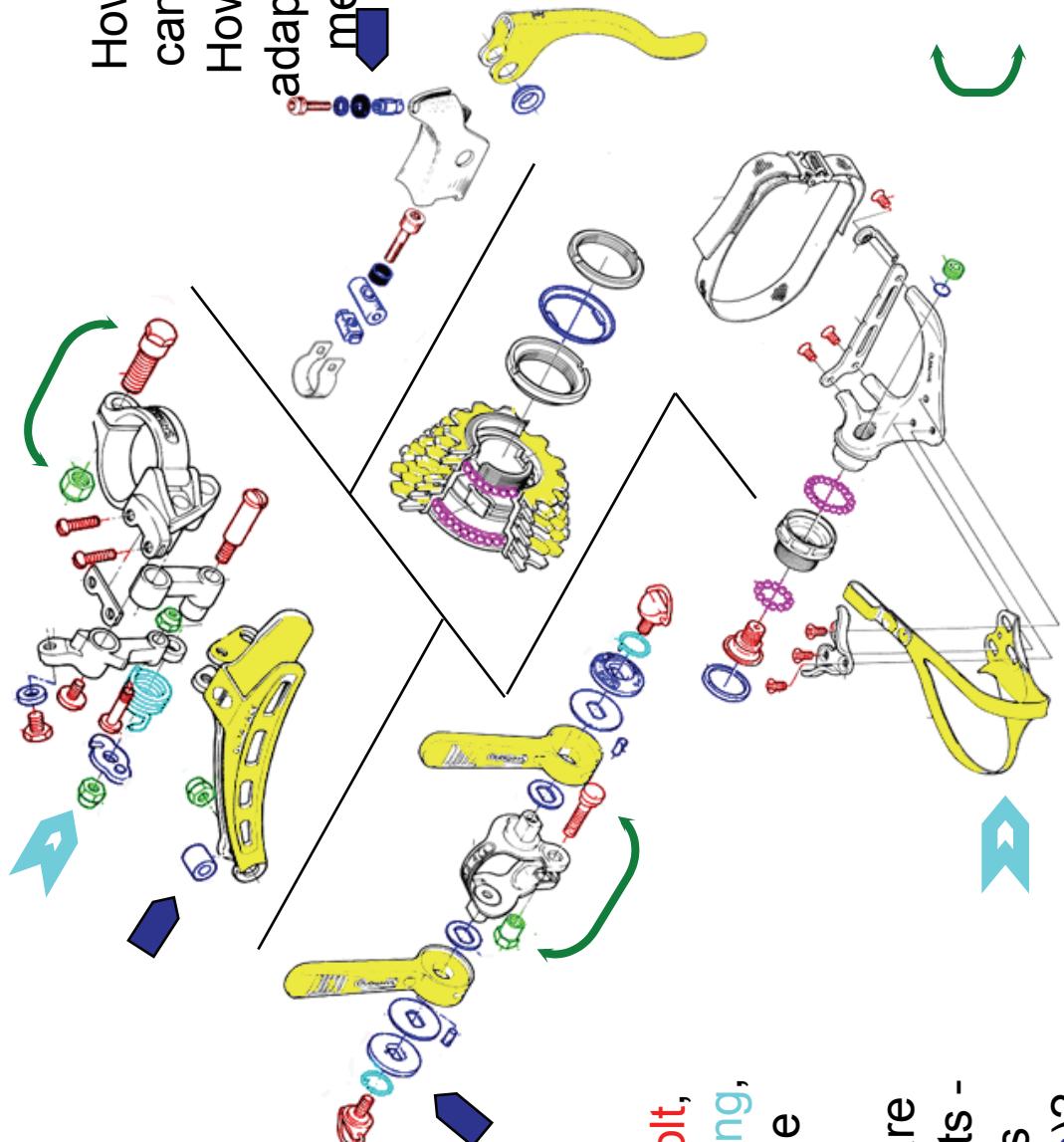


Q2: Integrative Ultrascalar “Visual-omics”

A collection of Parts



From Parts to Parts Assembly

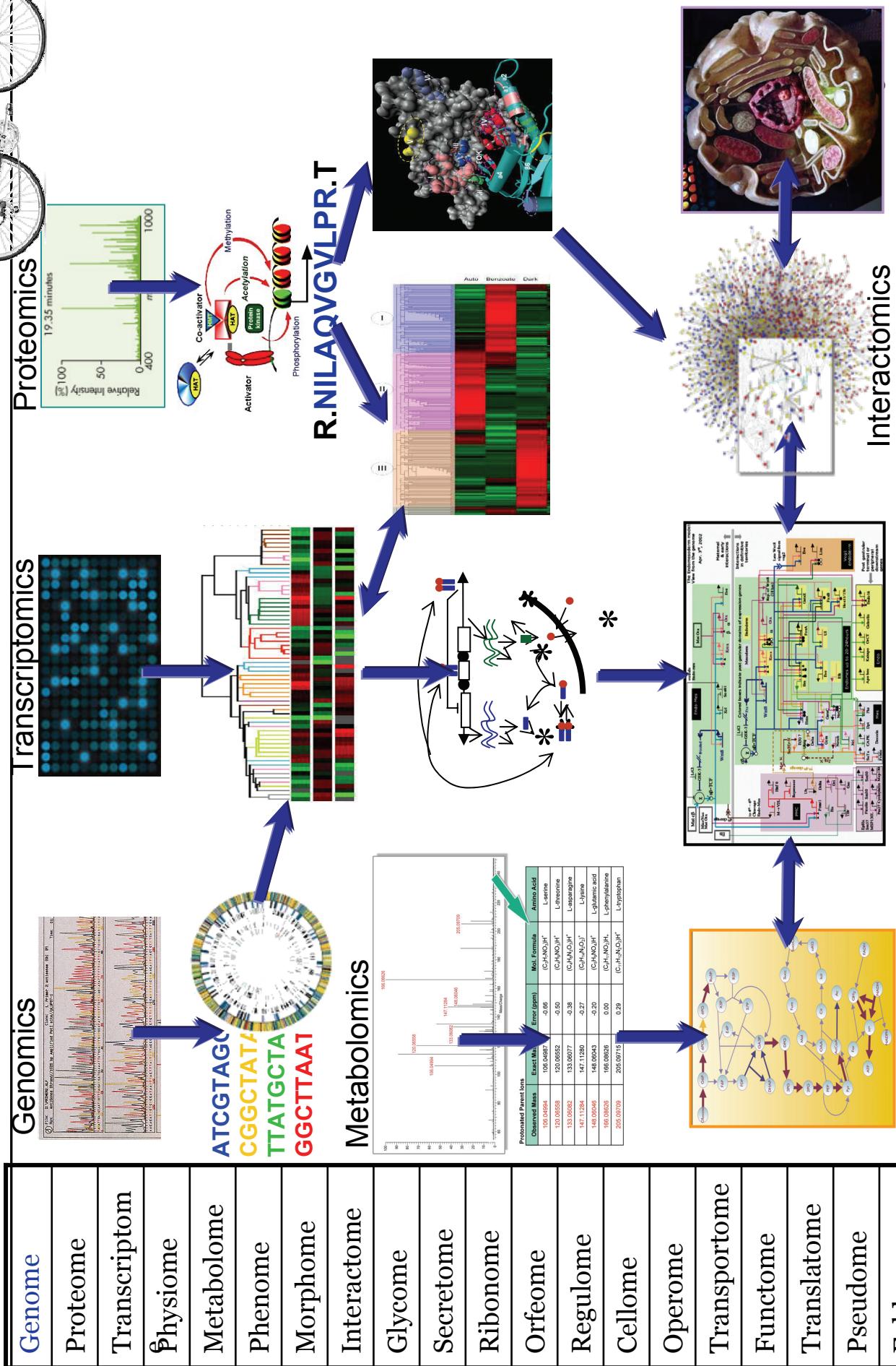


What are the shared parts (**bolt**, **nut**, **washer**, **spring**, **bearing**), unique parts (**cogs**, **levers**)? What are the common parts - types of parts (**nuts** & **washers**)?

How many roles can these play?
How flexible and adaptable are they mechanically?

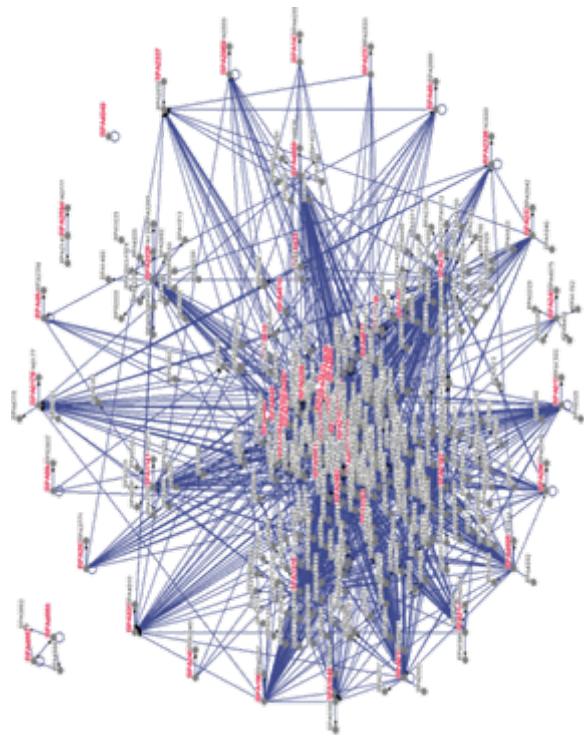
Where are the parts located?
Which parts interact?

Towards dynamic systems through Integrative “visual-omics”

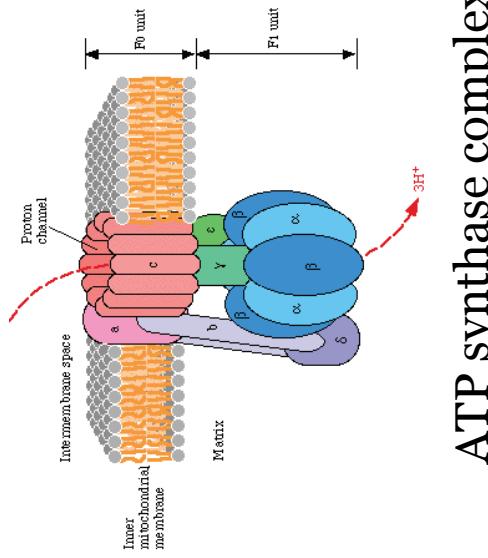


Q3: “Fuzzy” Visualization w/ Knowledge Priors

Mass spec pull-down experiments



Prior-knowledge



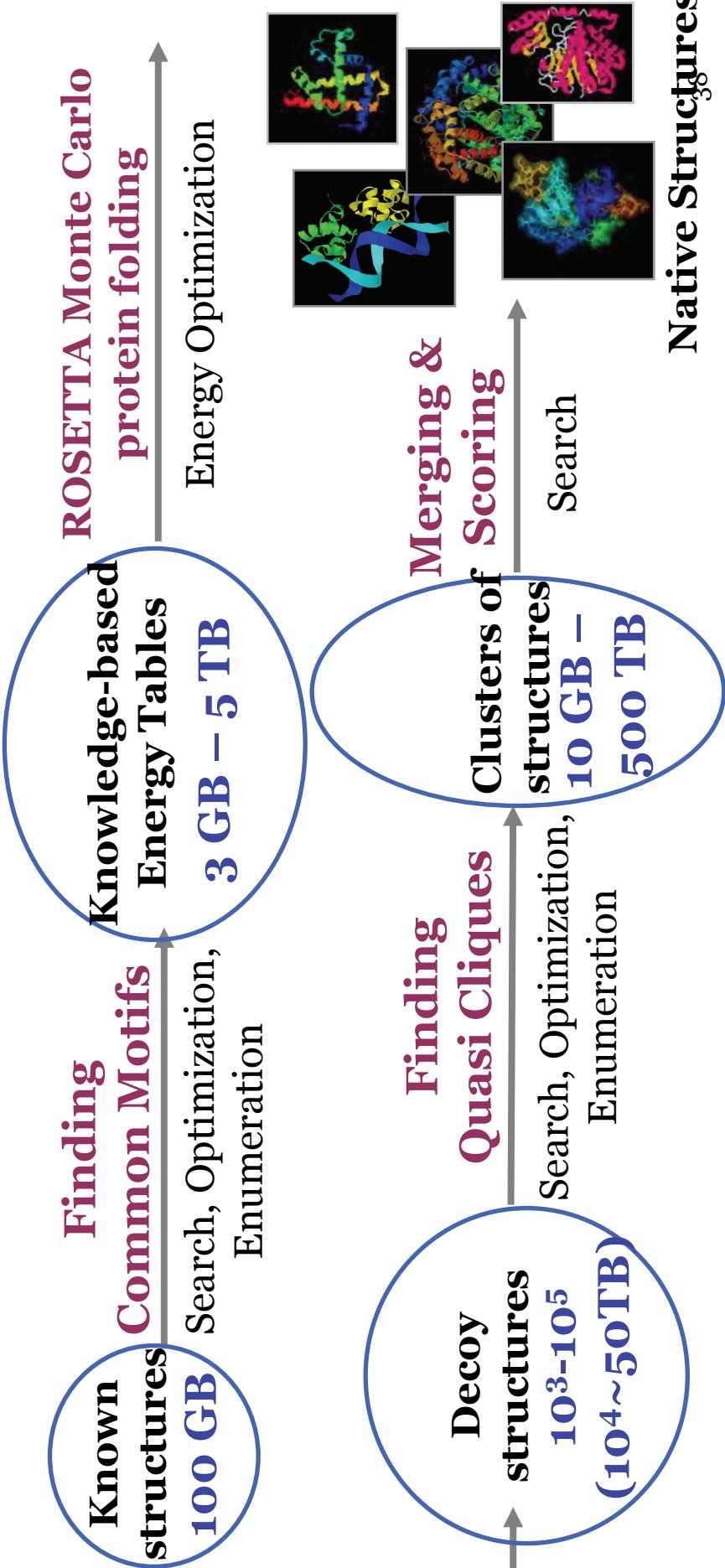
ATP synthase complex
or Constraints

- “Proactive” visualization
- Analogy: Query by example
- Analogy: Active learning
- Forms a (quasi) clique
- Contains at least 3 nodes
- Enriched by GO/KEGG terms
- Statistically significant

Q4: Visual Landscapes of Optimality

Each step is an *NP-hard combinatorial optimization problem* with different search heuristics.

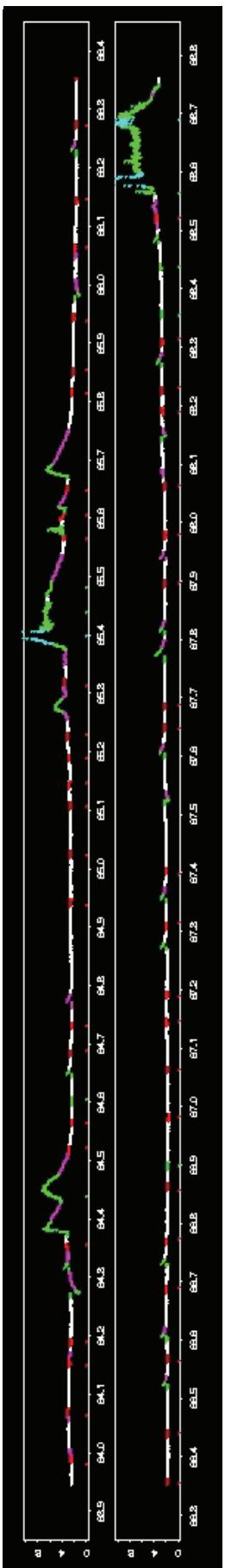
Example: *Ab Initio* Prediction of Protein 3-d Structure



Native Structures

Q5: Modeling the Usual to Discover the Unusual

To reduce data & detect extreme events in global context.



1. Segment series (100 obs)

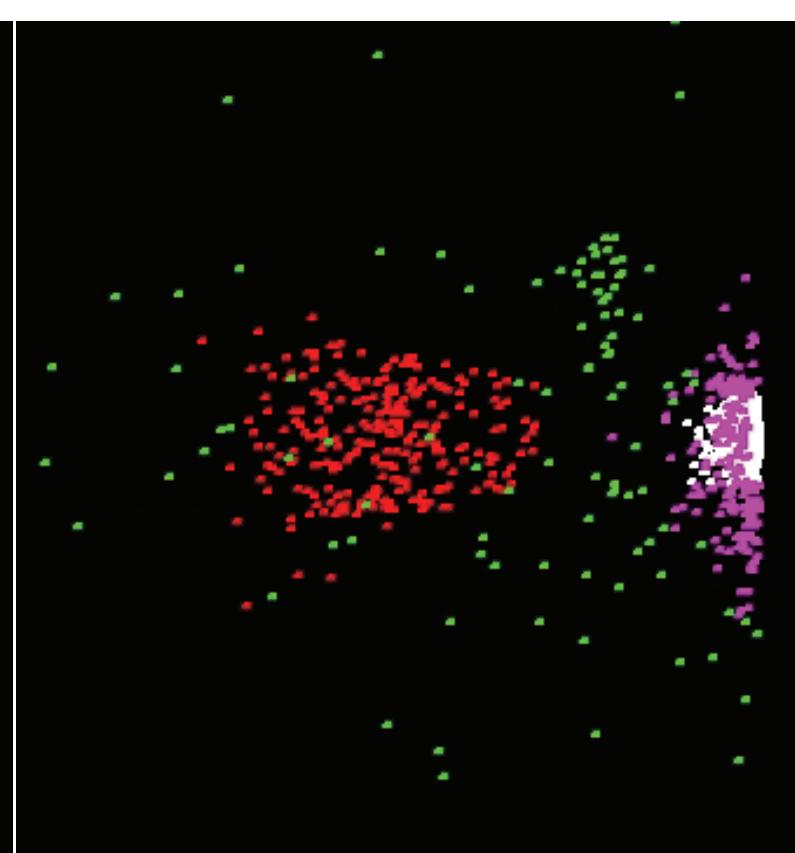
2. Fit simple local models to series
($c_0, c_1, c_2, \|e\|_\infty, \|e\|_2$)

3. Reduce data to model parameters

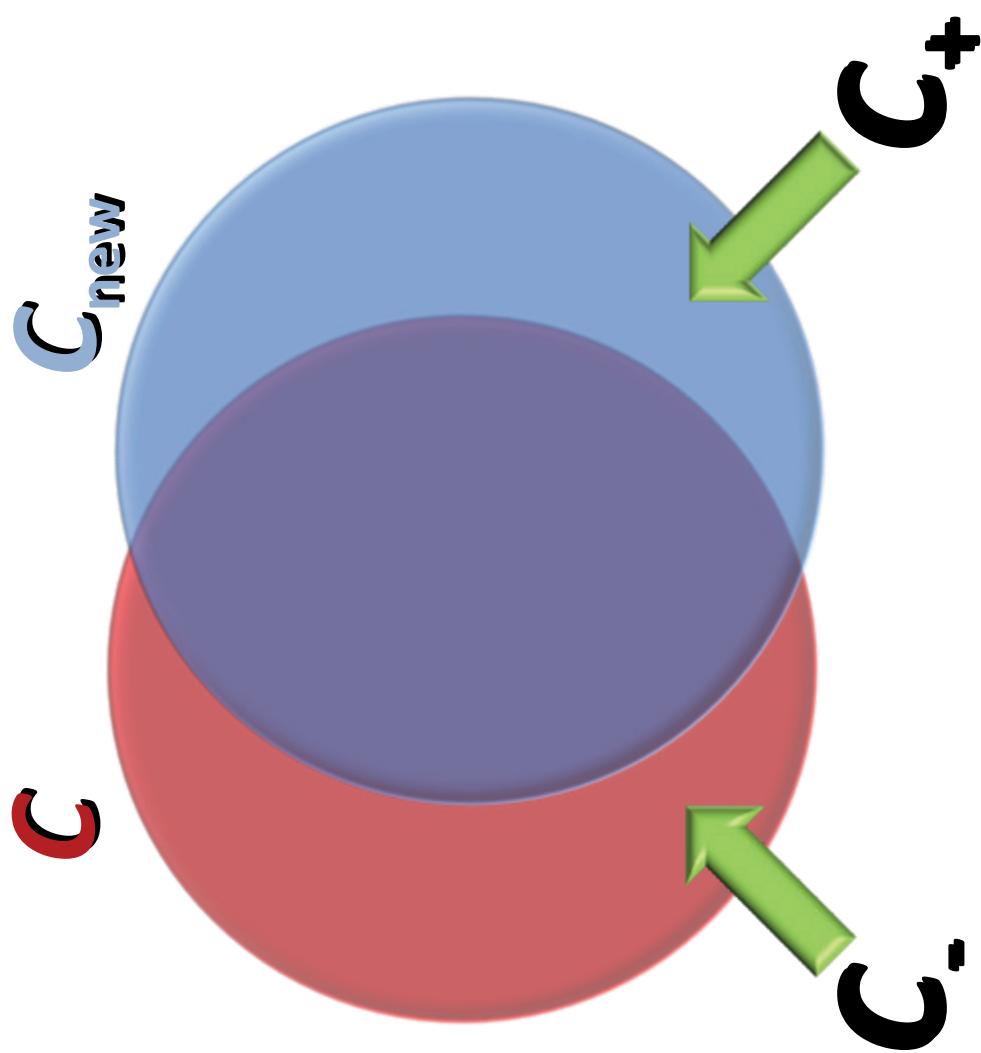
4. Select extremes for global analysis

5. Cluster the extremes (4)

6. Map back to series



Q6: “Perturbation” Visualization



Summary – Top 6 Quests for UV of Bio Data

1. Comparative “Visual-omics”
2. Integrative “Visual-omics”
3. “Fuzzy” Visualization with Knowledge Priors
4. Visual Landscapes of Optimalities
5. Modeling the Usual to Discover the Unusual
6. “Perturbation” Visualization

Acknowledgements

- Juan Huang, UTK
- DOE ASCR
- DOE BER