

HIGH PERFORMANCE  
COMPUTATIONAL ANALYSIS OF  
DNA SEQUENCES  
FROM DIFFERENT ENVIRONMENTS

Rob Edwards

Computer Science  
Biology



[edwards.sdsu.edu](http://edwards.sdsu.edu)

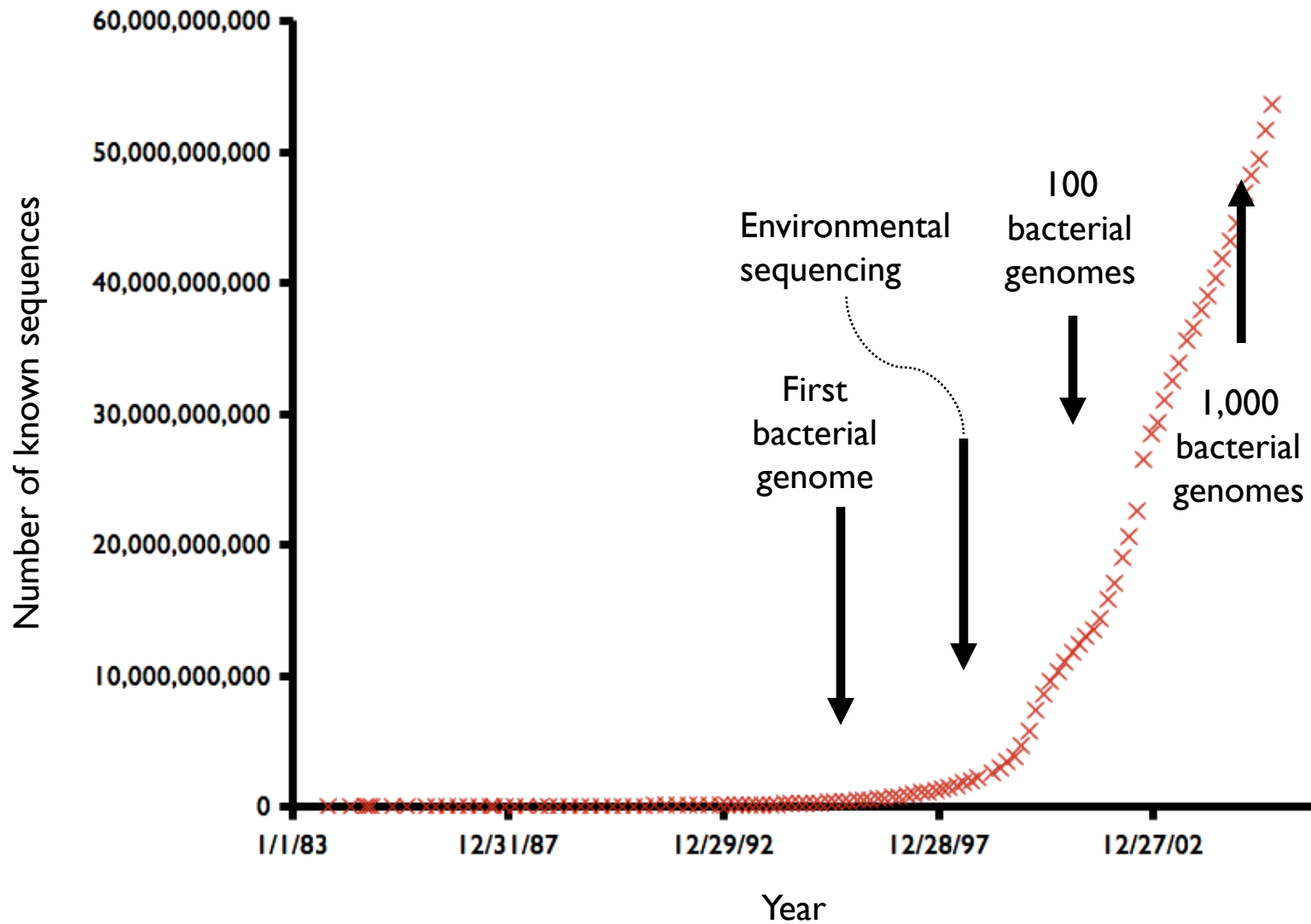
[www.theseed.org](http://www.theseed.org)



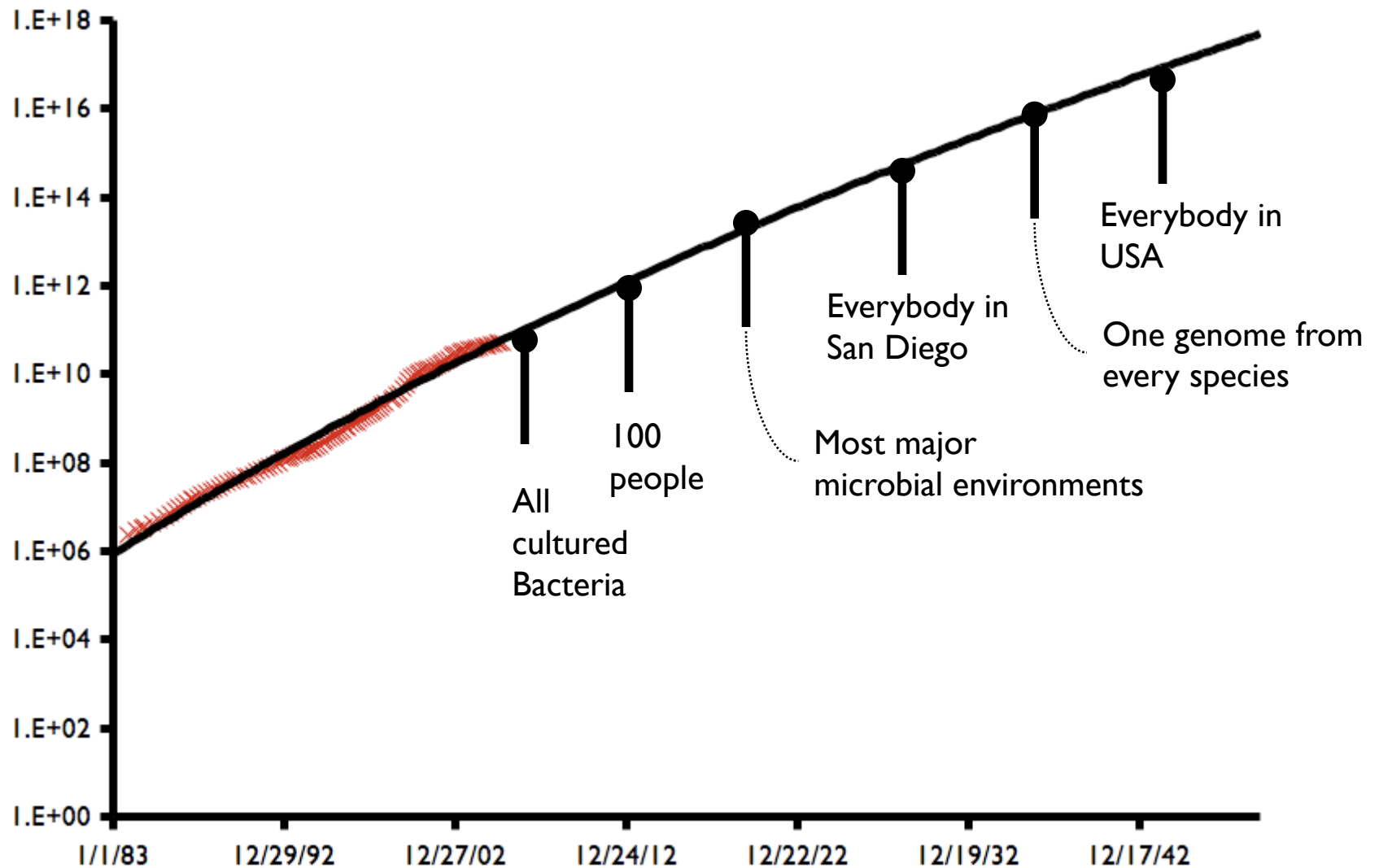
# OUTLINE

- There is a lot of sequence
- Tools for analysis
- More computers
- Can we speed analysis

# HOW MUCH HAS BEEN SEQUENCED?



# HOW MUCH WILL BE SEQUENCED?



# METAGENOMICS

(JUST SEQUENCE IT)

200 liters water  
5-500 g fresh fecal matter  
50 g soil



Concentrate and purify bacteria,  
viruses, etc



Extract nucleic acids

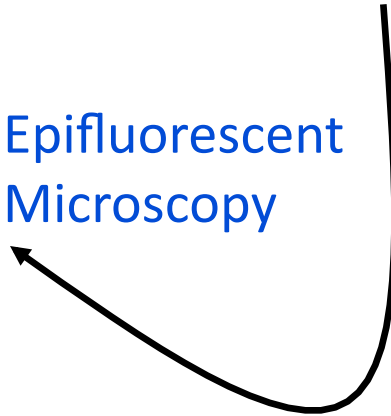
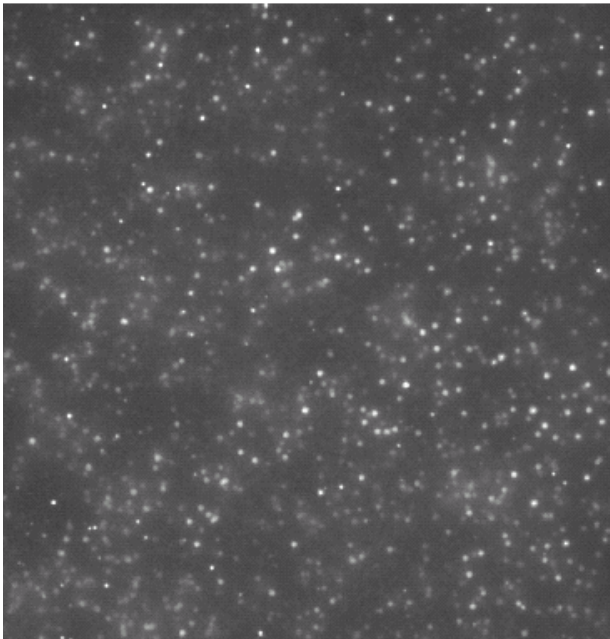


Sequence



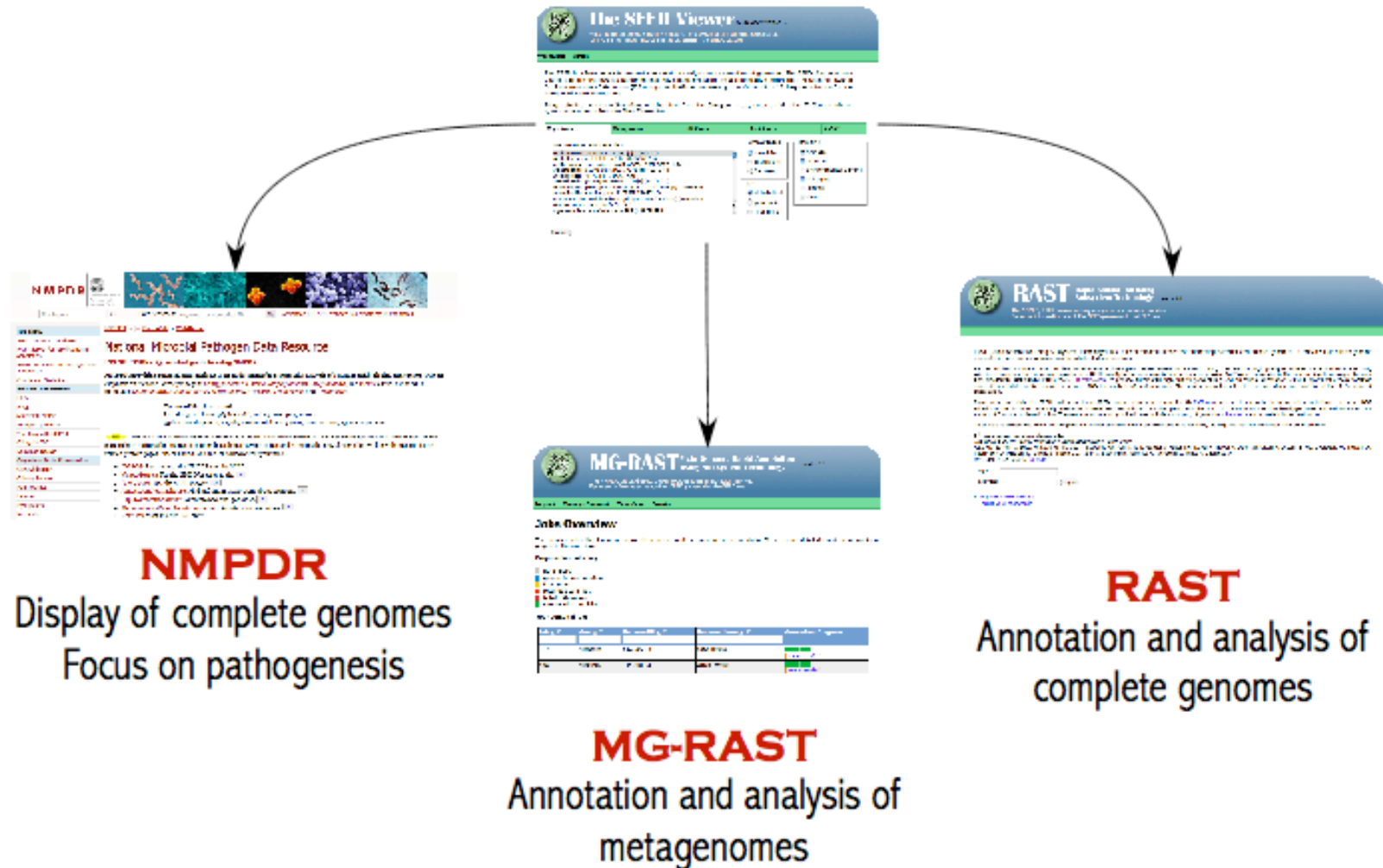
Publish papers

Epifluorescent  
Microscopy




# THE SEED FAMILY

## THE SEED Environmental, Viral, Bacterial, Archaeal, and Eukaryal Genome Interpretation



# THE METAGENOMICS RAST SERVER




## MG-RAST

Meta Genome Rapid Annotation  
using Subsystem Technology

Metagenomics SEED Viewer version 2.0

Welcome to the Metagenomics SEED Viewer.  
For more information about The SEED please visit [theSEED.org](http://theSEED.org).

[»Navigate](#) [»Help](#)

 Rob Edwards

**MG-RAST is a fully-automated service for annotating metagenome samples.**

It provides:

- **annotation** of sequence fragments,
- their **phylogenetic classification**,
- **metabolic reconstructions** and
- **comparison tools**

The service is built as a modified version of the RAST server which was originally designed to support high-quality annotation of complete or draft microbial genomes.



**If you use our service, please cite:**

*The Metagenomics RAST server - A public resource for the automatic phylogenetic and functional analysis of metagenomes* F. Meyer, D. Paarmann, M. D'Souza, R. Olson, E. M. Glass, M. Kubal, T. Paczian, R. Stevens, A. Wilke, J. Wilkening and R. A. Edwards submitted



» [Manage your uploaded data](#) » [Register a new account](#) » [Upload new metagenome to MG-RAST](#)

## You have access to the following metagenomes:

Number of publicly available metagenome: 158

start typing to narrow selection

Public: 5-Way (CG) Acid Mine Drainage Biofilm (4441137.3) from project

Public: 640F6 (4440355.3) from project Cow Rumen

Public: 710F6 (4440387.3) from project Cow Rumen

Public: 80F6 (4440356.3) from project Cow Rumen

Public: ALVINELLA (4441102.3) from project Alvinella Pompejana Epibio

Public: ArcticVir2002 (4440306.3) from project Ocean Viruses

Public: Australian Phosphorus Removing (EBPR) Sludge (4441092.3) from

Public: BBCVir96to04 (4440305.3) from project Ocean viruses

View Metagenome

# AUTOMATED PROCESSING

User uploads  
DNA sequences

Pre-processing:

✓ Genome Upload has been successfully completed.

Metagenome ID - Name:

Job:

User:

Date:

Number of uploaded sequences:

Total uploaded sequence length:

✓ Preprocessing has been successfully completed.

The following statistics are based on:

Number of sequences:

23

Total sequence length:

53295251

Average read length:

227.606717772416

Longest sequence id:

>58459919

Longest sequence length:

374

Shortest sequence id:

>58696741

Shortest sequence length:

36

✓ Similarity Computation has been successfully completed.

✓ Similarity postprocessing has been successfully completed.

✓ Final Assignment has been successfully completed.

## Jobs Details #570

» [Browse annotated metagenome in SEED Viewer](#)

» Available downloads for this job:

Genbank export



Download


» [Share this metagenome with selected users](#)

» [Make this metagenome publicly accessible](#)

» [Back to the Jobs Overview](#)




# SUMMARY VIEW

**MG-RAST** Meta Genome Rapid Annotation  
using Subsystem Technology

Metagenomics SEED Viewer version 2.0

Welcome to the Metagenomics SEED Viewer.  
For more information about The SEED please visit [theSEED.org](http://theSEED.org).

[»Navigate](#) [»Metagenome](#) [»Compare Metagenomes](#) [»Management](#) [»Help](#)

 Rob Edwards

## Metagenome Overview for MG\_Soybean\_Field\_Sample (4440698.3)

<b>Project:</b>	Fermi_metagenomics_samples
<b>Metagenome</b>	MG_Soybean_Field_Sample
<b>Metagenome ID:</b>	4440698.3
<b>Description:</b>	No description available.
<b>Uploaded on:</b>	Sun Mar 23 20:49:40 2008
<b>Total no. of sequences</b>	234,155
<b>Total sequence size</b>	53,295,251
<b>Shortest sequence length</b>	36
<b>Longest sequence length</b>	374
<b>Average sequence length</b>	227.61
<b>Average GC content</b>	not computed

Overview Metabolic Analysis Phylogenetic Analysis Compare

The metagenome overview page provides basic information and a summary regarding the selected metagenome. Information includes project name, project description, metagenome name and unique id as well as sequence length and percent GC statistics. Histograms of sequence length and GC content is also provided. In order to provide a brief overview of the taxonomic distribution, a table is provided with domain distribution for RNA and protein based analysis.

The Overview is accessible through the menu via  
» Metagenome » [Overview](#)

## Summary and Statistics

The MG\_Soybean\_Field\_Sample data set contains 234,155 contigs totaling 53,295,251 basepairs with an average fragment length of 227.61 (you can [download](#) the entire data set). A total of 61,041 sequences (26.07%) could be matched to proteins in [SEED subsystems](#) (using an e-value cut-off of 1e-5), you can explore metabolic reconstructions based on different parameters on the [Metabolic Reconstruction Page](#). Based on 94,012 hits against the SEED protein non-redundant database (40.15 % of the fragments) and on the 127 hits against the ribosomal RNA database [Greengenes](#) (0.05%) we computed the following table (using an e-value cut-off of 1e-5 and a minimum alignment length of 50bp).

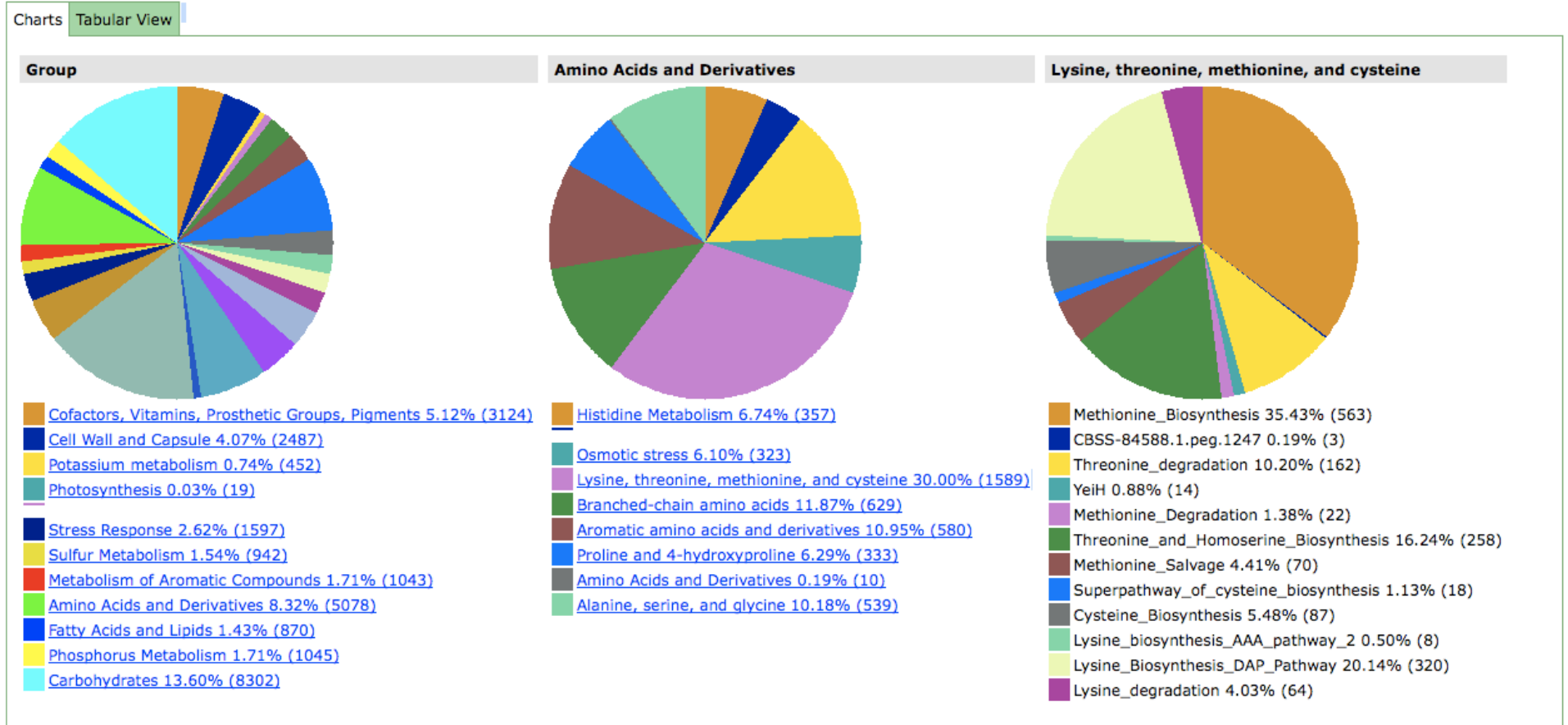
The [Phylogenetic Reconstruction](#) page will allow you to view taxonomic distributions in greater detail, change parameters and incorporate additional databases into your analysis.

The [MG-RAST manual](#) has more pointers for working with the system.

	Protein based	16s based
<b>Archaea</b>	2.07% (1946)	0.00% (0)
<b>Bacteria</b>	84.76% (79682)	90.55% (115)
<b>Eukaryota</b>	1.04% (980)	0.00% (0)
<b>Virus</b>	0.00% (0)	0.00% (0)
<b>Other</b>	12.13% (11404)	9.45% (12)

# METAGENOMICS TOOLS

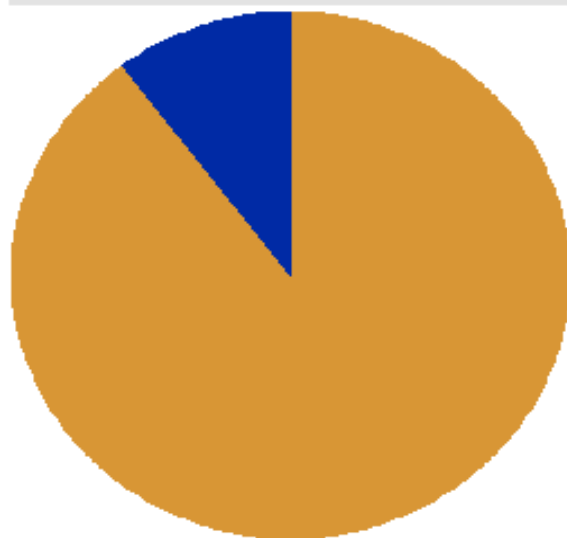
## ANNOTATION & SUBSYSTEMS



# METAGENOMICS TOOLS

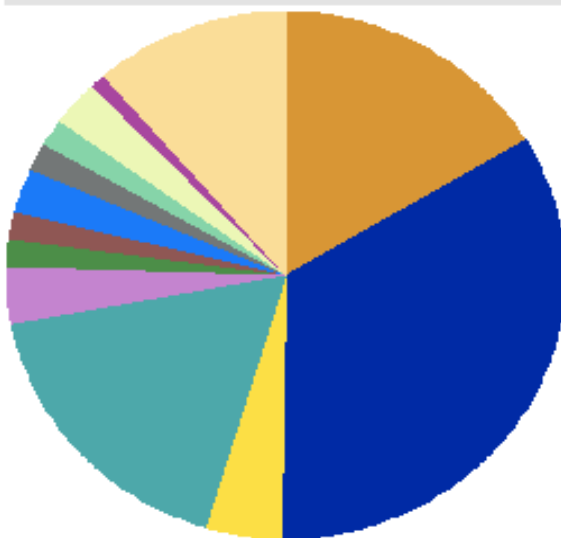
## PHYLOGENETIC RECONSTRUCTION

**Domain**



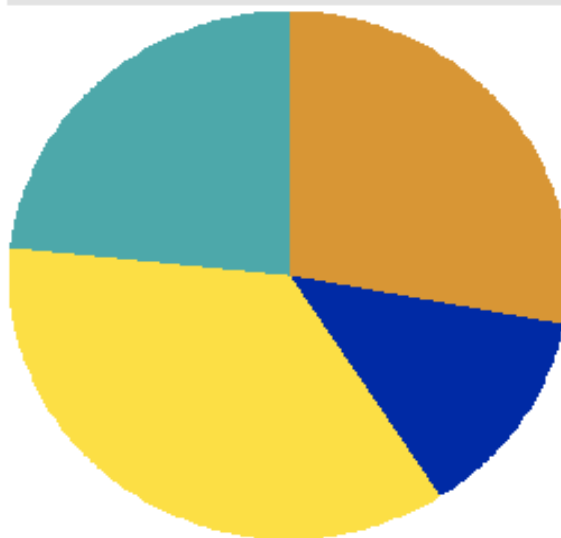
[Bacteria 89.84% \(115\)](#)  
[Unclassified 10.16% \(13\)](#)

**Bacteria**



[Actinobacteria 16.52% \(19\)](#)  
[Proteobacteria 33.91% \(39\)](#)  
[Bacteroidetes 4.35% \(5\)](#)  
[Unclassified 17.39% \(20\)](#)  
[Planctomycetes 3.48% \(4\)](#)  
[Genera incertae sedis OP10 1.74% \(2\)](#)  
[Genera incertae sedis WS3 1.74% \(2\)](#)  
[Firmicutes 2.61% \(3\)](#)  
[Chloroflexi 1.74% \(2\)](#)  
[Acidobacteria 1.74% \(2\)](#)  
[Nitrospira 2.61% \(3\)](#)  
[Verrucomicrobia 0.87% \(1\)](#)  
[Gemmatimonadetes 11.30% \(13\)](#)

**Proteobacteria**



[Alphaproteobacteria 28.21% \(11\)](#)  
[Gammaproteobacteria 12.82% \(5\)](#)  
[Betaproteobacteria 35.90% \(14\)](#)  
[Deltaproteobacteria 23.08% \(9\)](#)

# METAGENOMICS TOOLS

## COMPARATIVE TOOLS

Subsystem Hierarchy 1 ▲▼ all	4440690.3 ▲▼	4440740.3 ▲▼	4440739.3 ▲▼	4440698.3 ▲▼
Carbohydrates	<a href="#">0.0059</a>	<a href="#">0.0074</a>	<a href="#">0.0084</a>	<a href="#">0.0106</a>
Clustering-based subsystems	<a href="#">0.0035</a>	<a href="#">0.0045</a>	<a href="#">0.0051</a>	<a href="#">0.0066</a>
Amino Acids and Derivatives	<a href="#">0.0030</a>	<a href="#">0.0038</a>	<a href="#">0.0042</a>	<a href="#">0.0053</a>
Virulence	<a href="#">0.0030</a>	<a href="#">0.0036</a>	<a href="#">0.0039</a>	<a href="#">0.0048</a>
Cofactors, Vitamins, Prosthetic Groups, Pigments	<a href="#">0.0024</a>	<a href="#">0.0029</a>	<a href="#">0.0032</a>	<a href="#">0.0042</a>
Respiration	<a href="#">0.0021</a>	<a href="#">0.0025</a>	<a href="#">0.0027</a>	<a href="#">0.0035</a>
Protein Metabolism	<a href="#">0.0020</a>	<a href="#">0.0024</a>	<a href="#">0.0028</a>	<a href="#">0.0035</a>
Cell Wall and Capsule	<a href="#">0.0018</a>	<a href="#">0.0022</a>	<a href="#">0.0026</a>	<a href="#">0.0032</a>
Unclassified	<a href="#">0.0018</a>	<a href="#">0.0022</a>	<a href="#">0.0024</a>	<a href="#">0.0030</a>
Metabolism of Aromatic Compounds	<a href="#">0.0015</a>	<a href="#">0.0018</a>	<a href="#">0.0020</a>	<a href="#">0.0025</a>
RNA Metabolism	<a href="#">0.0015</a>	<a href="#">0.0018</a>	<a href="#">0.0020</a>	<a href="#">0.0024</a>
Stress Response	<a href="#">0.0013</a>	<a href="#">0.0017</a>	<a href="#">0.0018</a>	<a href="#">0.0022</a>
Membrane Transport	<a href="#">0.0010</a>	<a href="#">0.0012</a>	<a href="#">0.0014</a>	<a href="#">0.0017</a>
DNA Metabolism	<a href="#">0.0009</a>	<a href="#">0.0012</a>	<a href="#">0.0013</a>	<a href="#">0.0017</a>
Regulation and Cell signaling	<a href="#">0.0009</a>	<a href="#">0.0012</a>	<a href="#">0.0011</a>	<a href="#">0.0016</a>
Nucleosides and Nucleotides	<a href="#">0.0008</a>	<a href="#">0.0010</a>	<a href="#">0.0012</a>	<a href="#">0.0014</a>

# OUTLINE

- There is a lot of sequence
- Tools for analysis
- More computers

## HOW MUCH DATA SO FAR

986 metagenomes

~300 GS20

79,417,238 sequences

~300 FLX

~300 Sanger

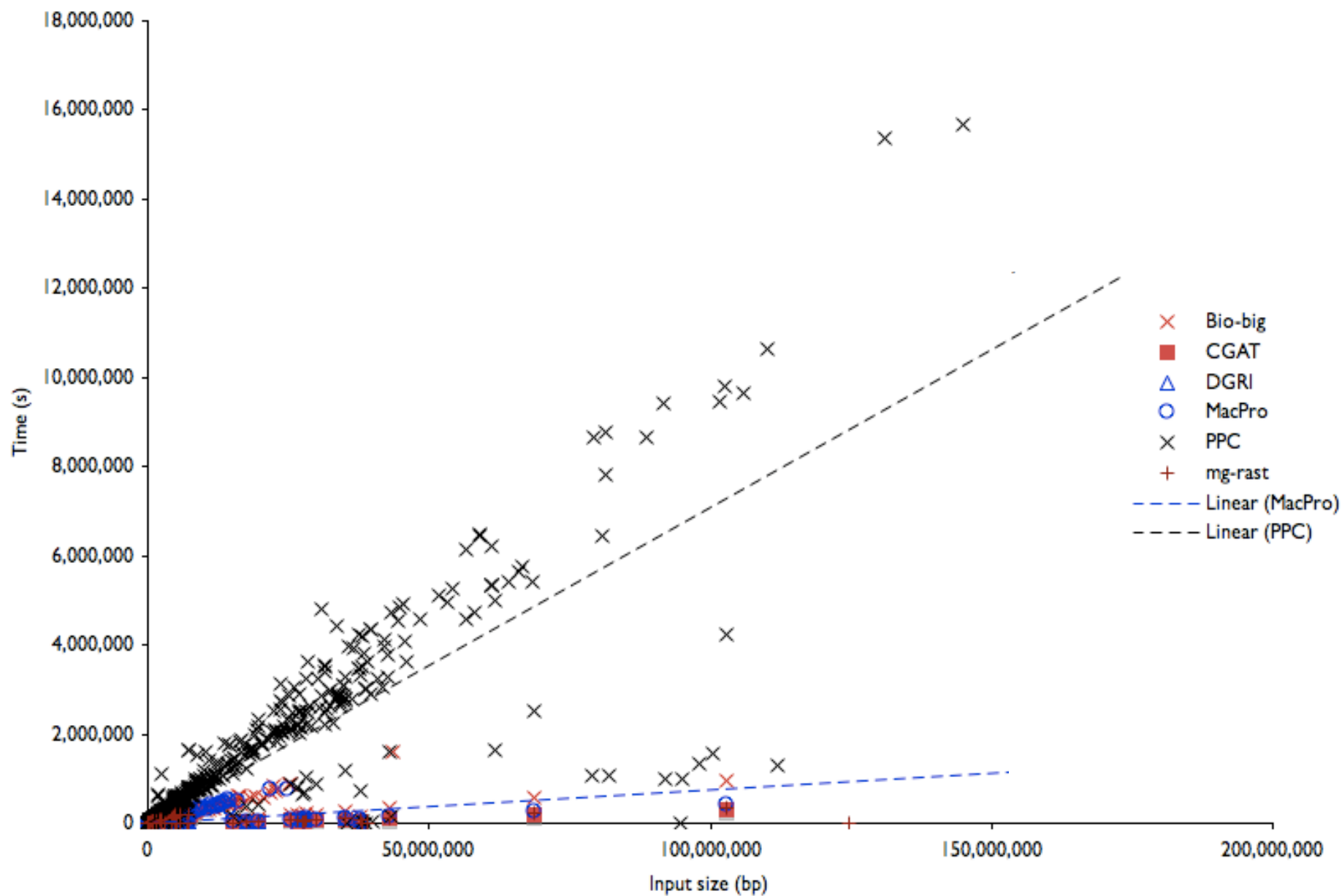
17,306,834,870 bp (17 Gbp)

Average: ~15-20 M bp per genome

## COMPUTES

<i><b>Processor</b></i>	<i><b>Memory</b></i>	<i><b>Number of units</b></i>
16 Intel Xeon CPU X7350 @ 2.93GHz	123,823 MB	1
8 Intel Xeon CPU X5365 @ 3.00GHz	16,436 MB	1
8 Intel Xeon CPU E5335 @ 2.00GHz	16,440 MB	1
1 Intel Pentium 4 CPU 3.00GHz	3,636 MB	2
8 Intel Xeon CPU X5365 @ 3.00GHz	16,387 MB	2
8 Intel Xeon CPU E5450 @ 3.00GHz	16,436 MB	1
8 Intel Xeon CPU X5365 @ 3.00GHz	16,436 MB	1
2 PPC970FX, altivec supported	4,042 MB	45

# LINEAR COMPUTE COMPLEXITY



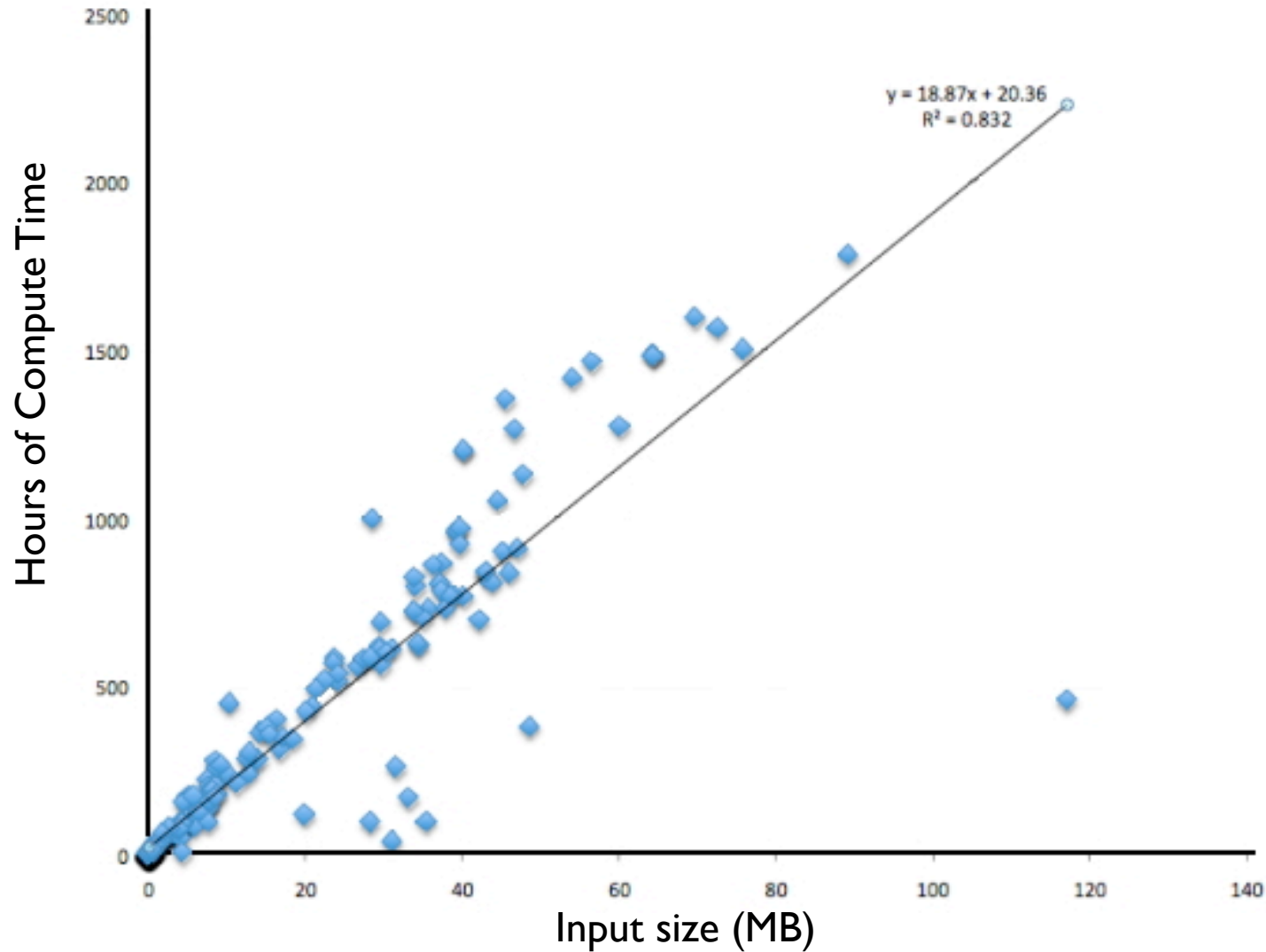


JUST WAITING ...



## OVERALL COMPUTE TIME

~19 hours of compute per input megabyte



## HOW MUCH SO FAR

986 metagenomes

~300 GS20

79,417,238 sequences

~300 FLX

~300 Sanger

17,306,834,870 bp (17 Gbp)

Average: ~15-20 M bp per genome

Compute time (on a single CPU):

328,814 hours = 13,700 days = 38 years

# OUTLINE

- There is a lot of sequence
- Tools for analysis
- More computers
- Can we speed analysis

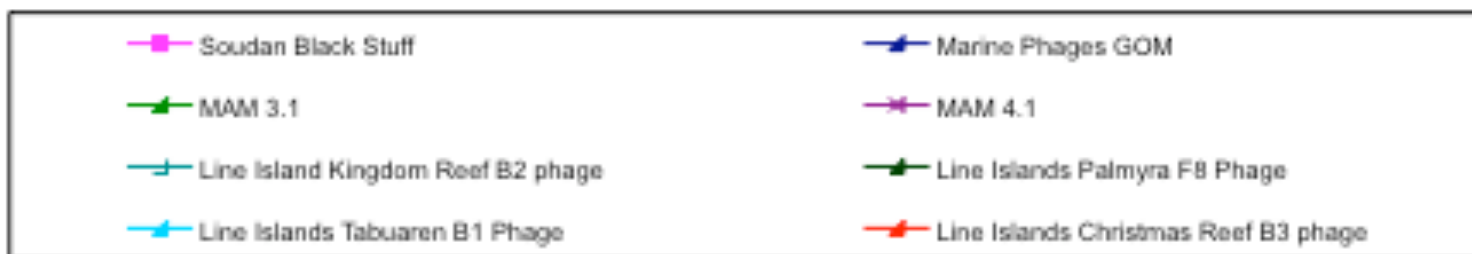
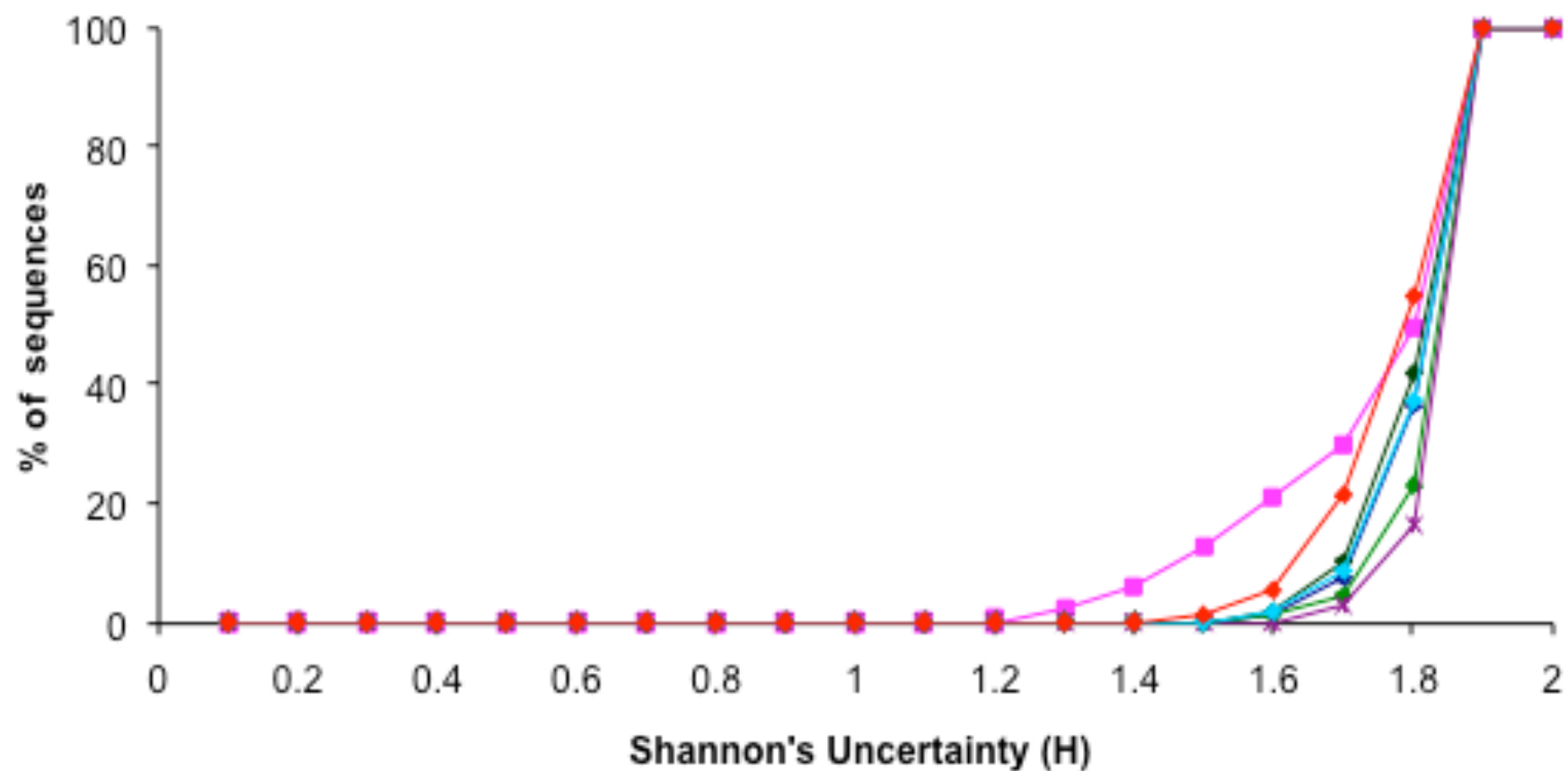
## SHANNON'S UNCERTAINTY

- Shannon's Uncertainty – Peter's surprisal

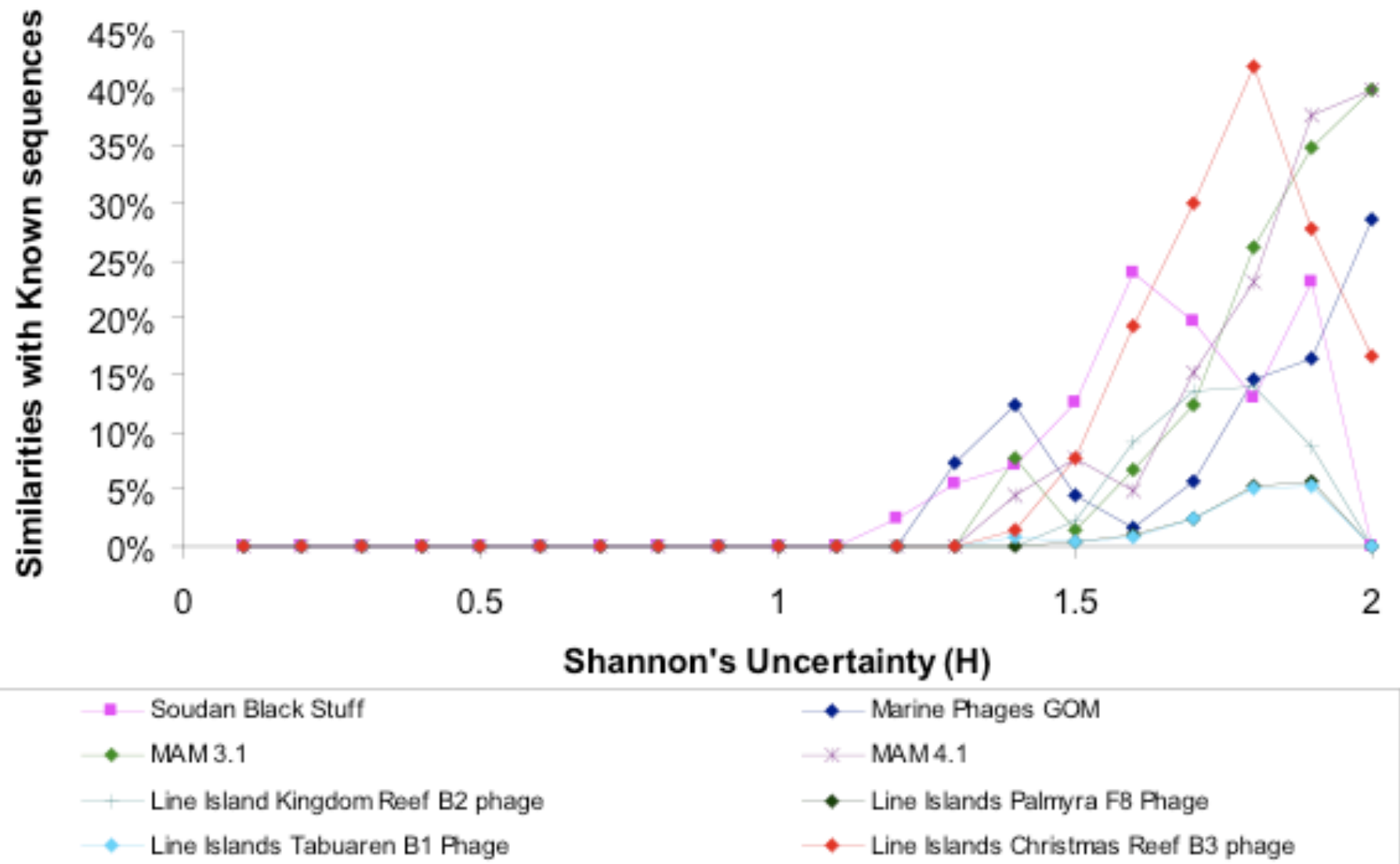
$$H(X) := - \sum_{i=1}^n p(x_i) \log_b p(x_i)$$

$p(x_i)$  is the probability of the occurrence of each base or string

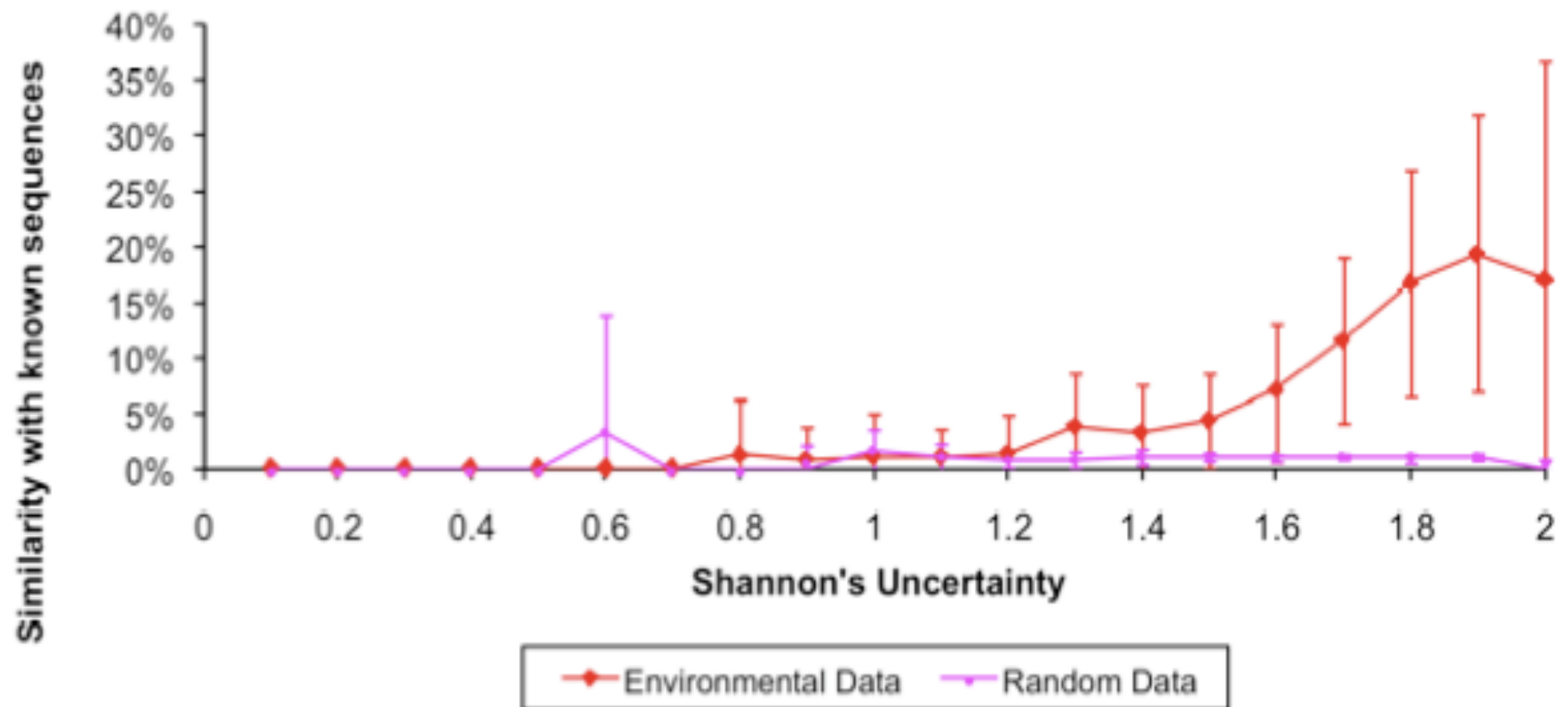
# SURPRISAL IN SEQUENCES



# UNCERTAINTY CORRELATES WITH SIMILARITY

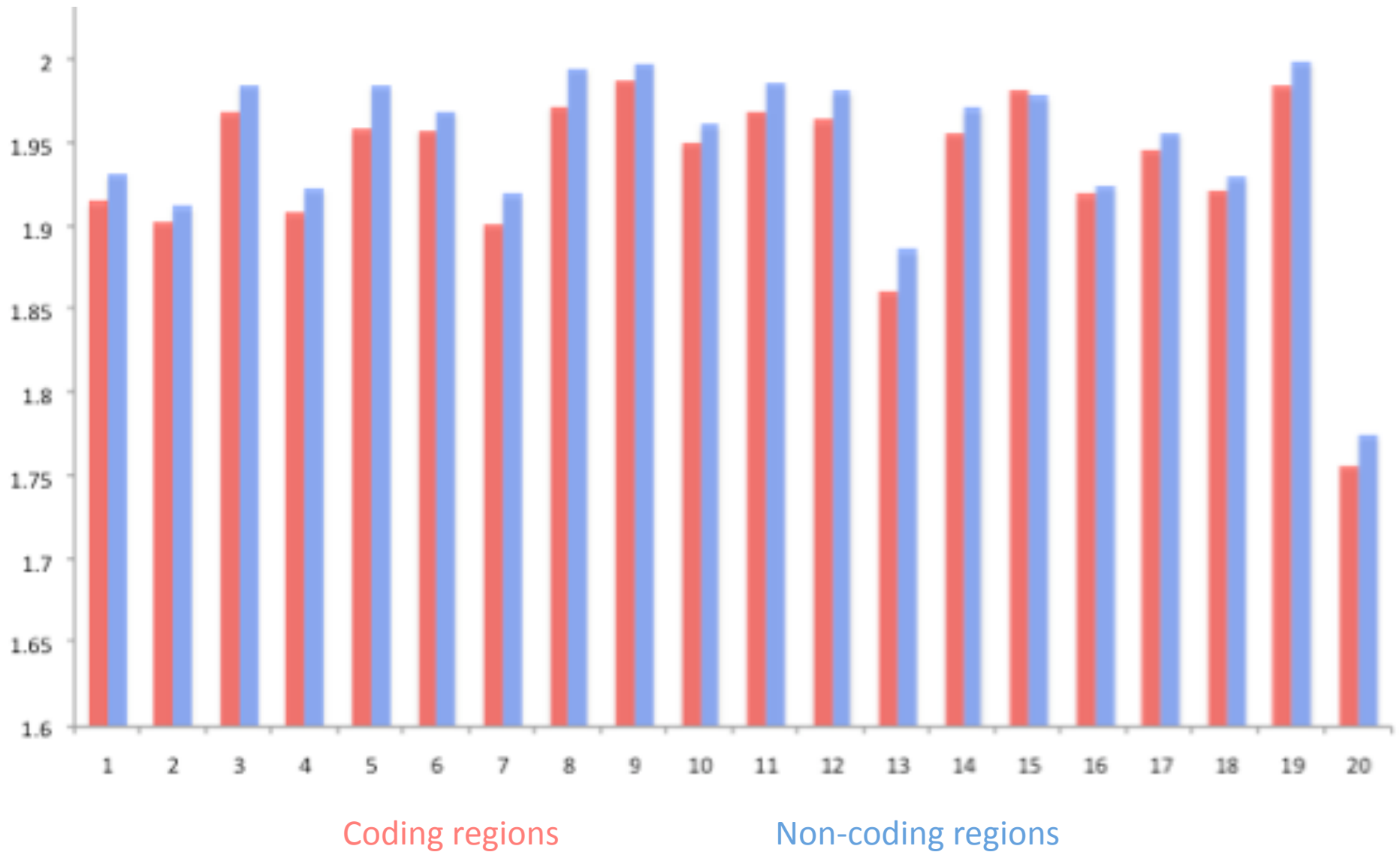


BUT IT'S NOT JUST RANDOMNESS...

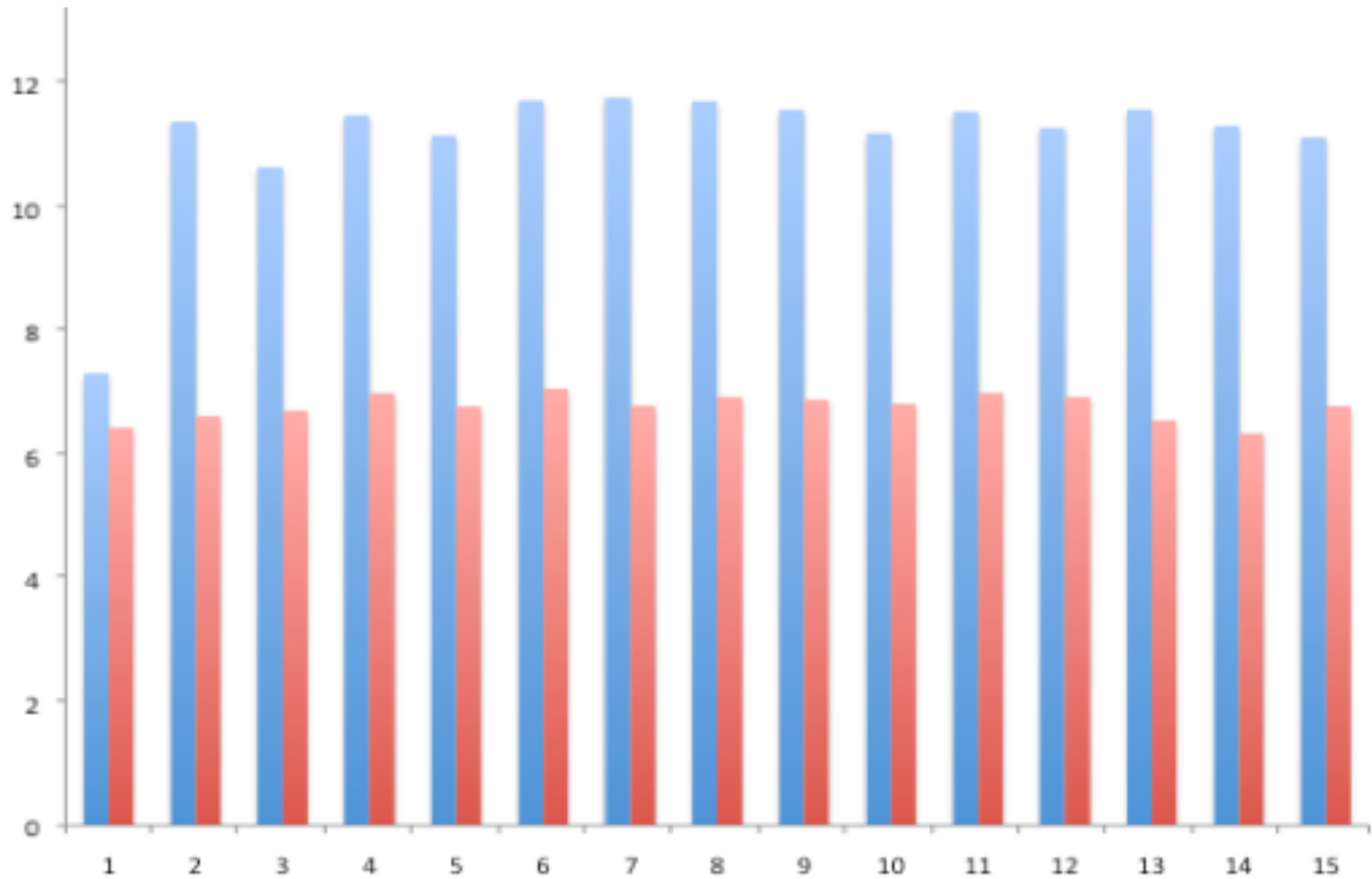




# UNCERTAINTY IN COMPLETE GENOMES



## MORE EXTREME DIFFERENCES WITH 6-MERS



Coding regions

Non-coding regions

## CAN WE PREDICT PROTEINS

- Short sequences of 100 bp
- Translate into 30-35 amino acids
- Can we predict which are real and could be doing something?
- Test with bacterial proteins

# KULLBACK-LEIBLER DIVERGENCE

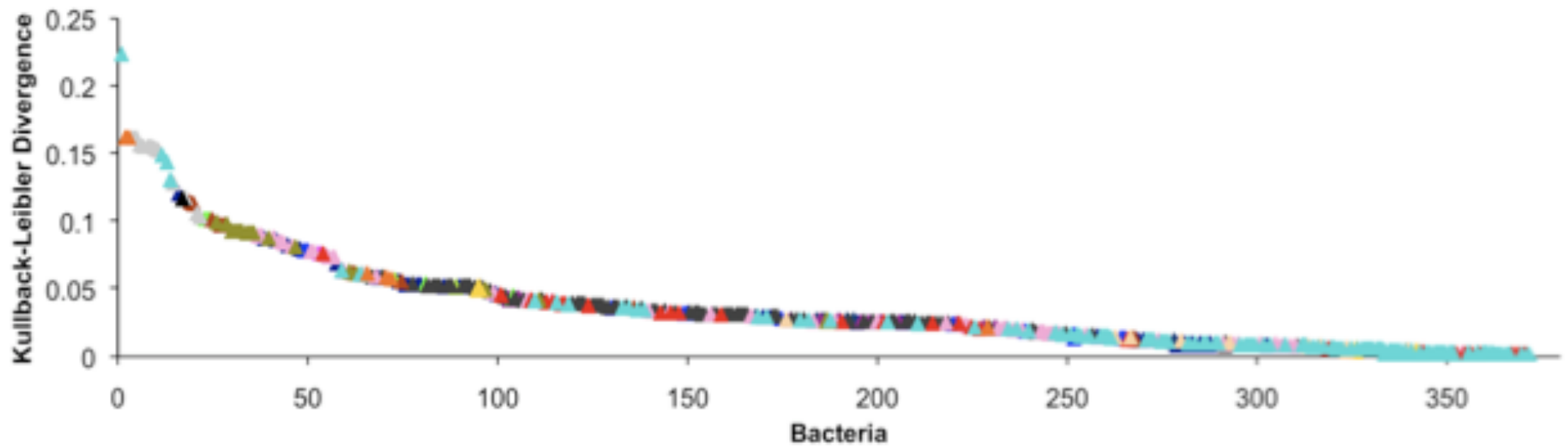
Difference between two probability distributions

$$D_{\text{KL}}(P\|Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

Difference between amino acid composition and average amino acid composition

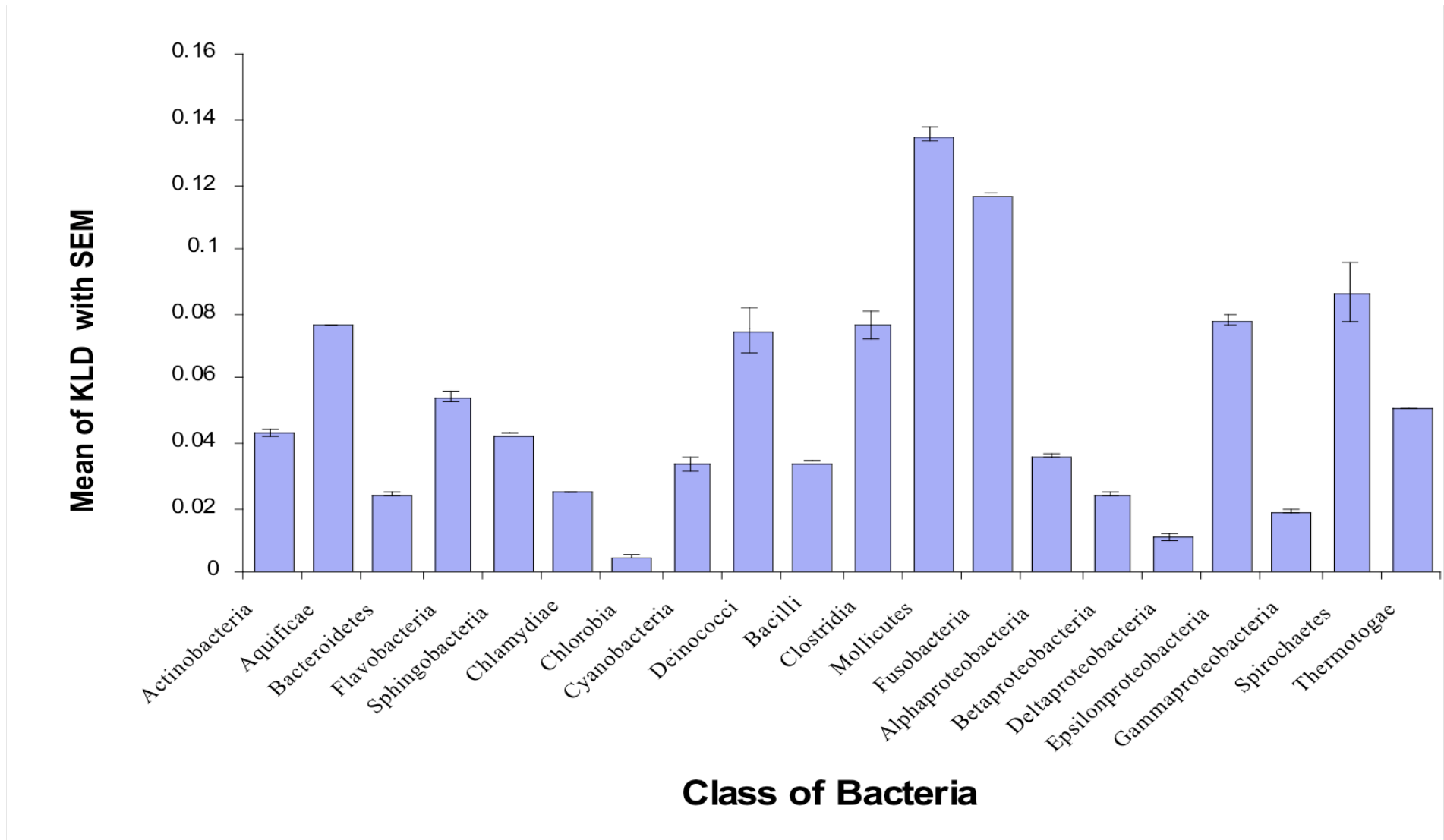
Calculate KLD for 372 bacterial genomes

# KLD VARIES BY BACTERIA



Colored by taxonomy of the bacteria

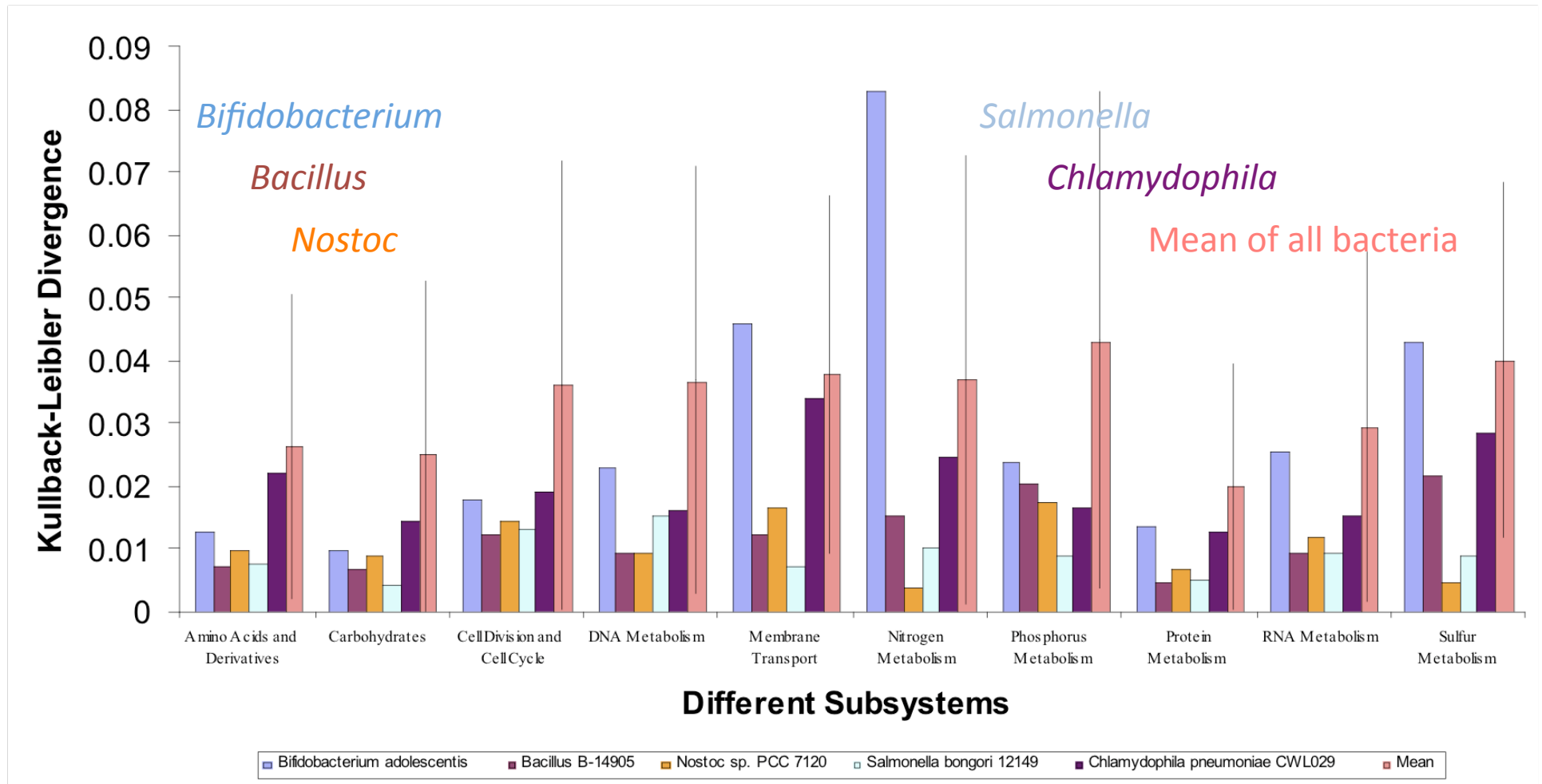
# KLD VARIES BY BACTERIA



## MOST DIVERGENT GENOMES

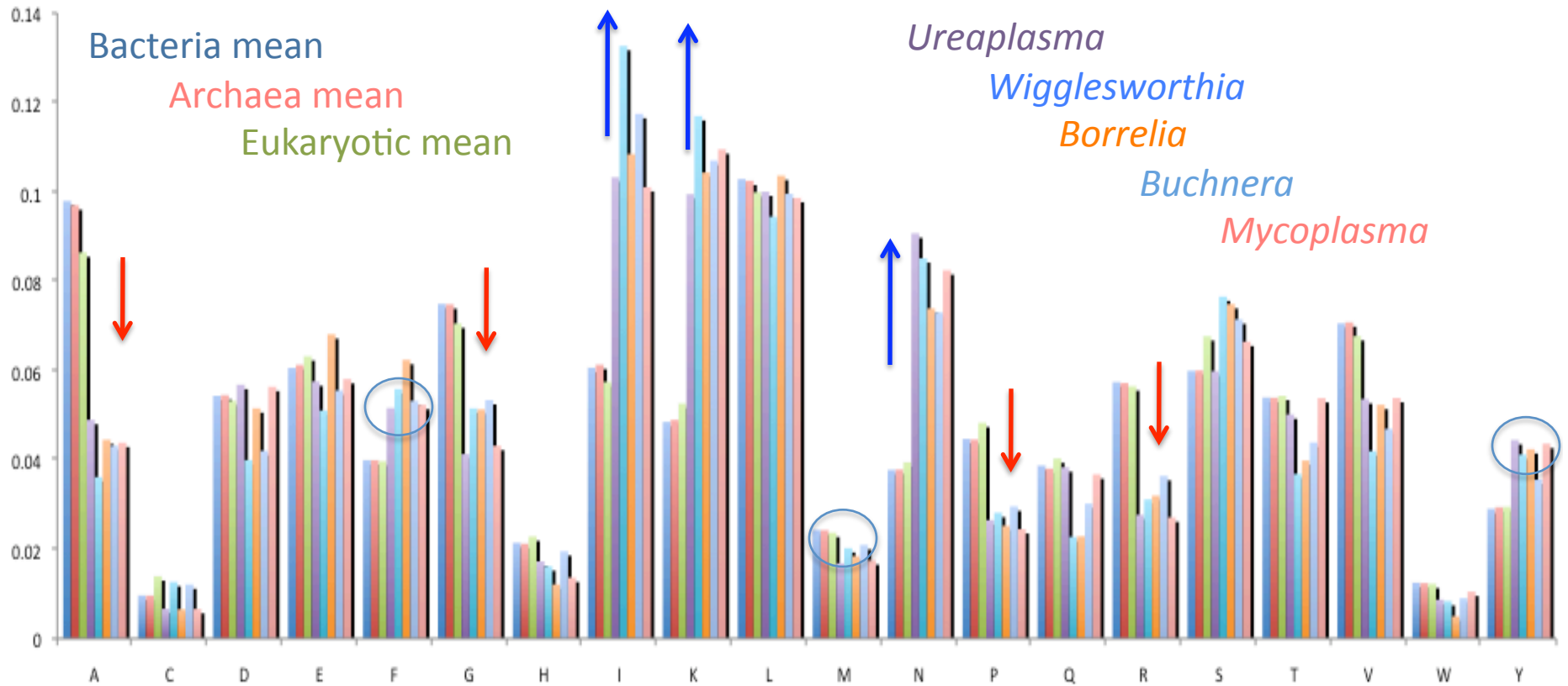
- *Borrelia garinii* – Spirochaetes
- *Mycoplasma mycoides* – Mollicutes
- *Ureaplasma parvum* – Mollicutes
- *Buchnera aphidicola* – Gammaproteobacteria
- *Wigglesworthia glossinidia* – Gammaproteobacteria

# DIVERGENCE AND METABOLISM

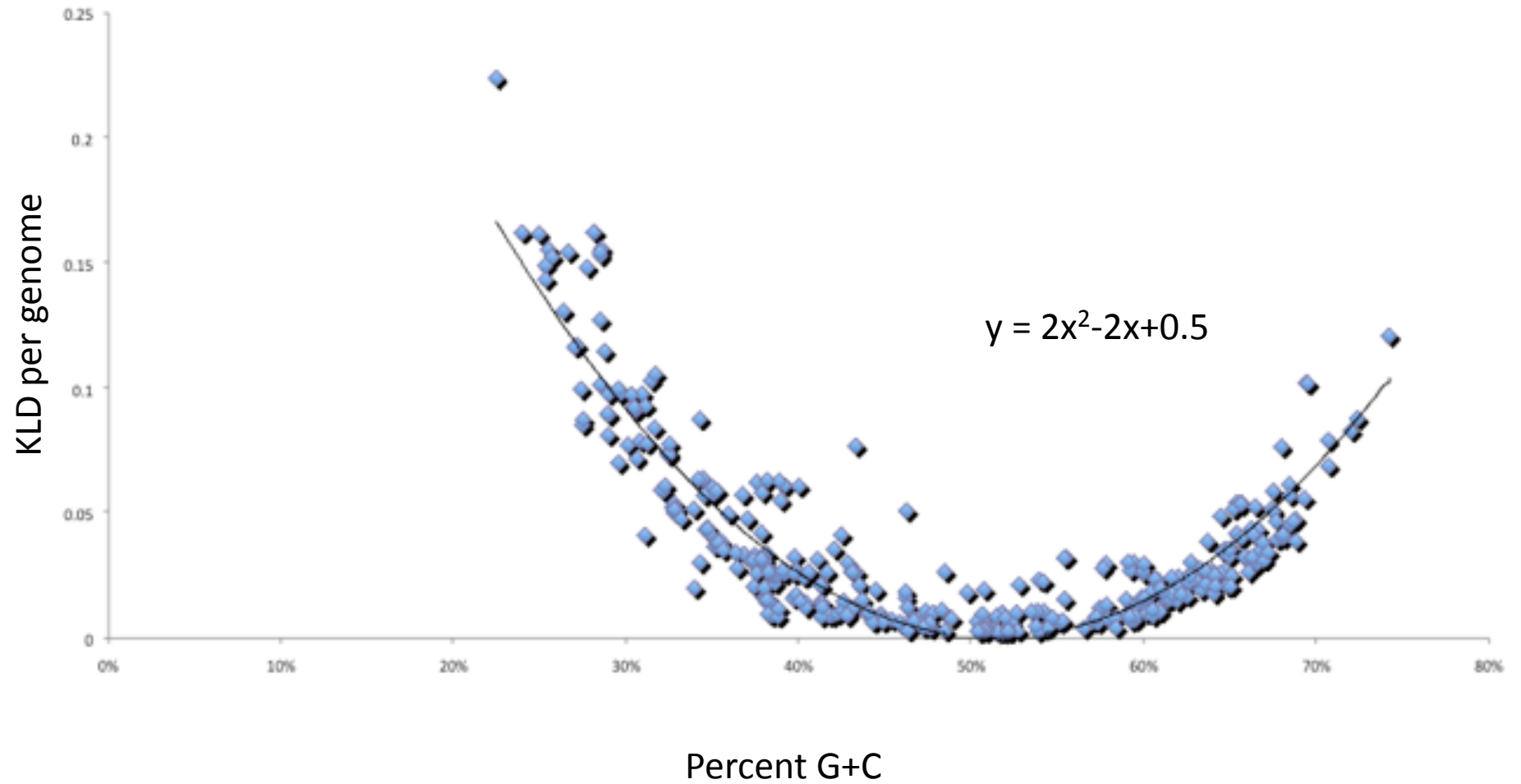




# DIVERGENCE AND AMINO ACIDS



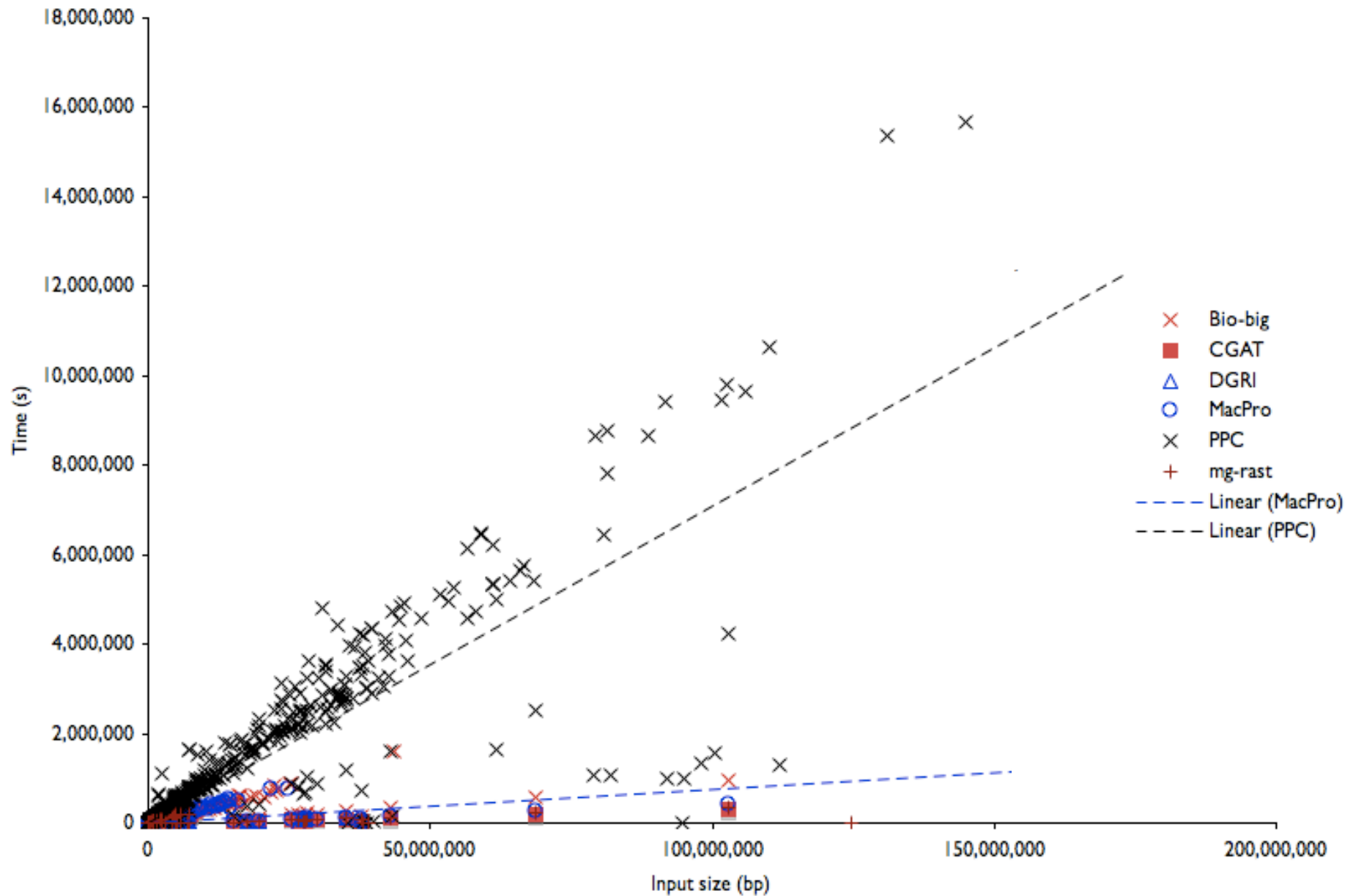
# PREDICTING KLD



## SUMMARY

- Shannon's uncertainty could predict useful sequences
- KLD varies too much to be useful and is driven by %G+C content

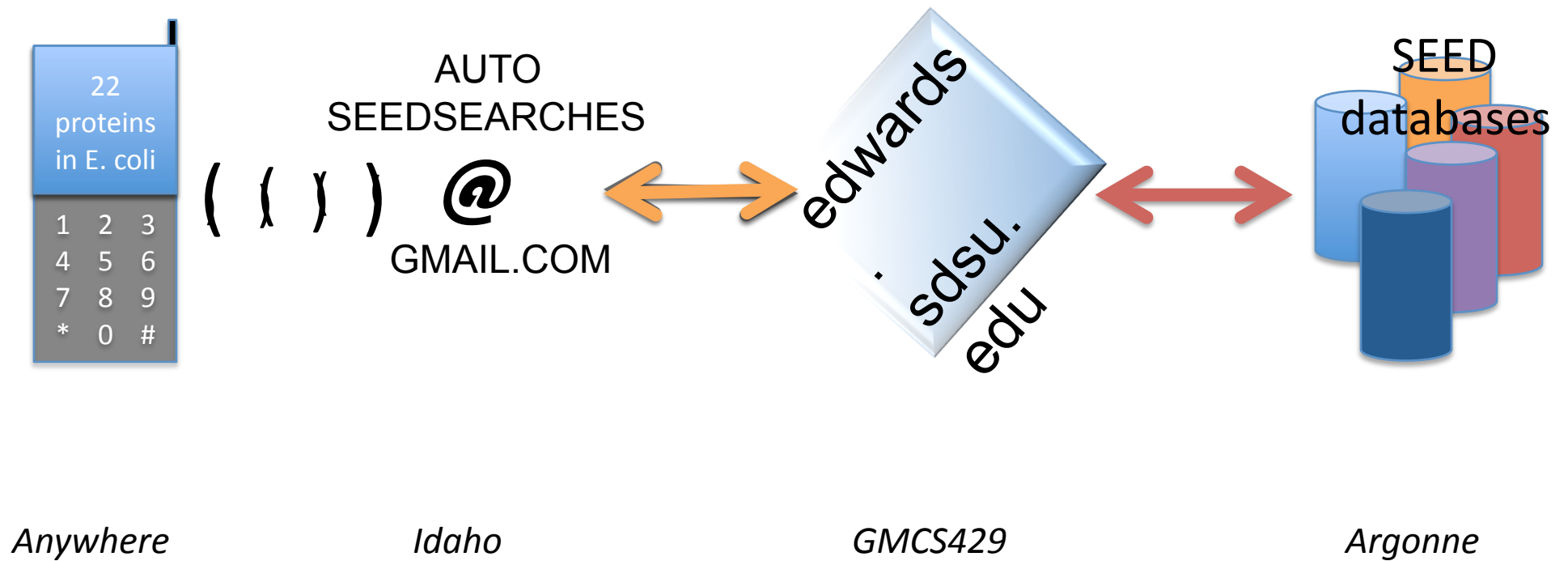
# NEW SOLUTIONS FOR OLD PROBLEMS?



# XEN AND THE ART OF IMAGERY

# THE CELL PHONE PROBLEM

# SEARCHING THE SEED BY SMS



## CHALLENGES

- Too much data
- Not easy to prioritize
- New models for HPC needed
- New interfaces to look at data



## ACKNOWLEDGEMENTS

- Sajia Akhter
- Rob Schmieder
- Nick Celms
- Sheridan Wright
- Ramy Aziz
- FIG
- The mg-RAST team
- Rick Stevens
- Peter Salamon
- Barb Bailey
- Forest Rohwer
- Anca Segall