# Metabolic Reconstructions from Global Ocean Sampling (GOS) Marine Metagenome

Mathangi Thiagarajan

J. Craig Venter Institute

Pathways Tools Workshop 2010

**J. Craig Venter**
**I N S T I T U T E**

- Metagenomics
- The Global Ocean Sampling (GOS) Project
- GOS - Community Makeup
- High Throughput Data Processing
- Metabolic Reconstruction – Mapping to MetaCyc and KEGG
- Metarep (Visualization) – Integrating with MetaCyc and KEGG
- Pathways Tools for GOS & metagenomic projects
- Conclusion
- Acknowledgements

J. Craig Venter
I N S T I T U T E

# Metagenomics

- Examining genomic content of organisms in community/environment to better understand
  - Diversity of organisms
  - Their roles and interactions in the ecosystem

- Cultivation independent approach to study microbial communities
  - DNA directly isolated from environmental sample and sequenced

J. Craig Venter
I N S T I T U T E
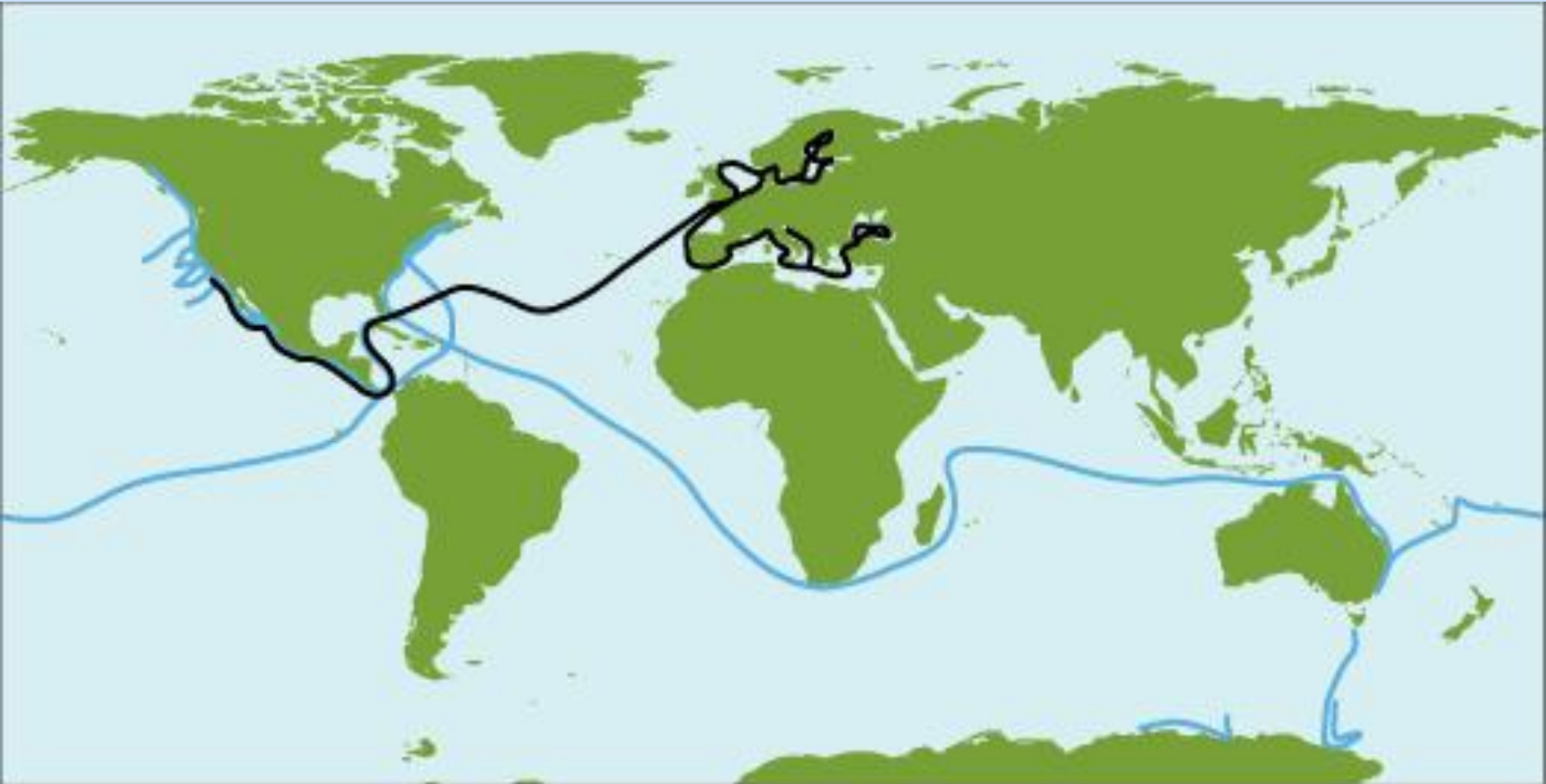
# Global Ocean Sampling Expedition

Investigate the fundamental microbial contributions from the Ocean waters to energy and nutrient cycling by analyzing its

a) biogeochemical cycling
b) community structure and function
c) microbial diversity
d) adaptation and evolution

GOS  Phase I  - Published in PLOS Biology 2007

GOS  Circumnavigation - Analysis Phase

# Global Ocean Sampling Expedition Route
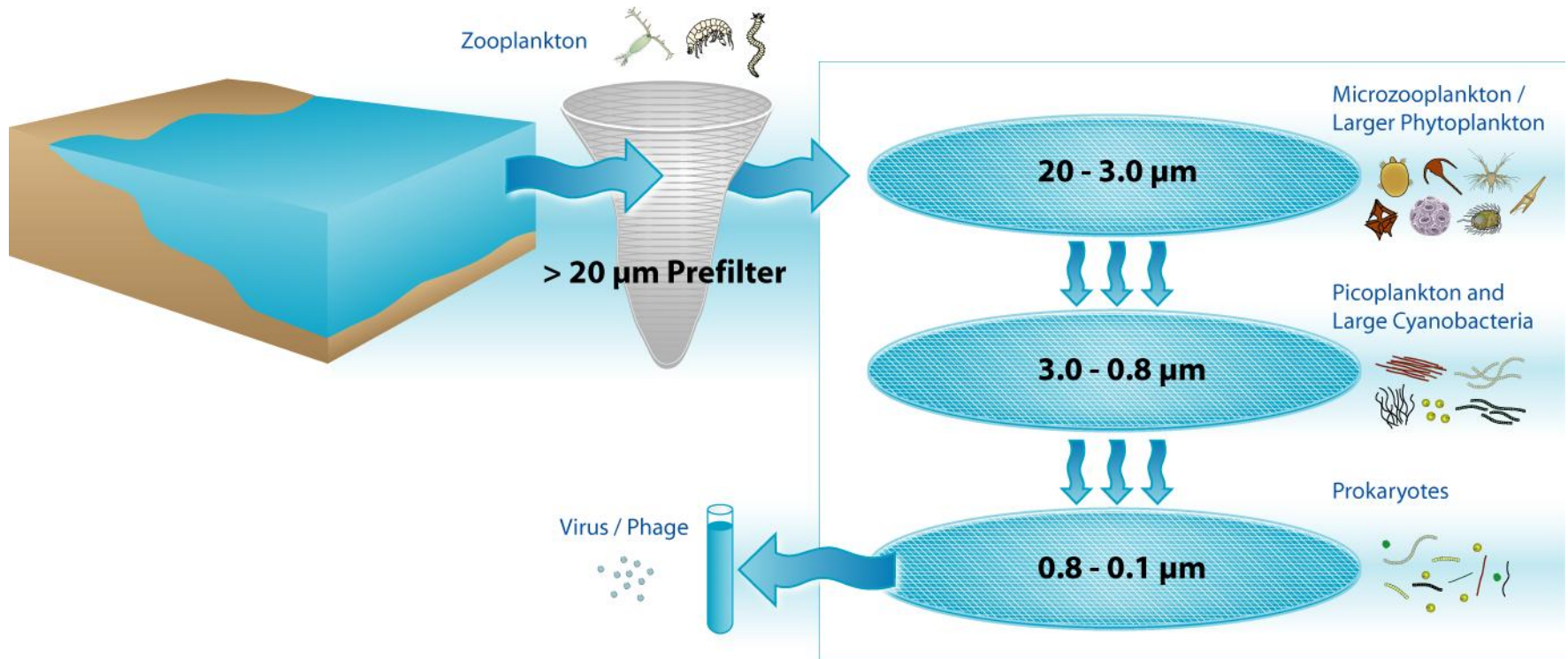


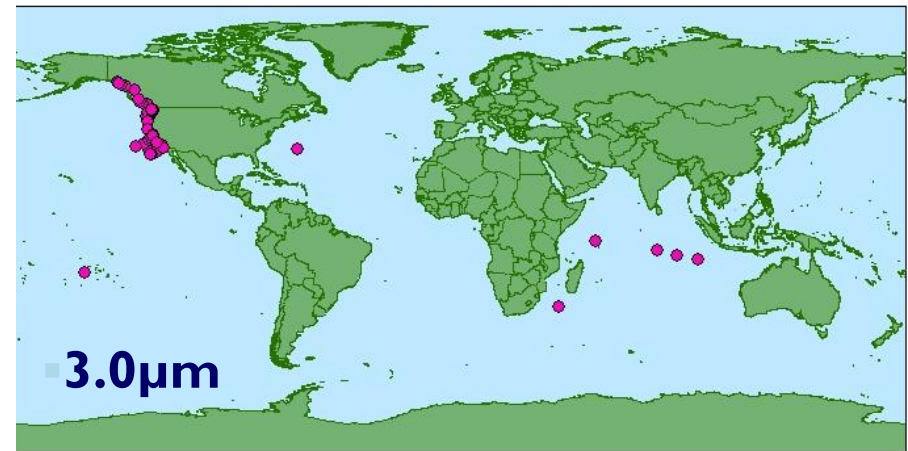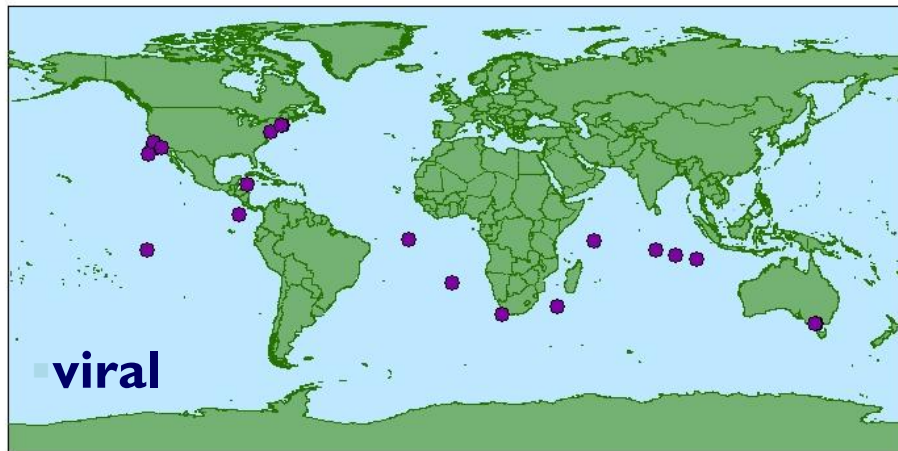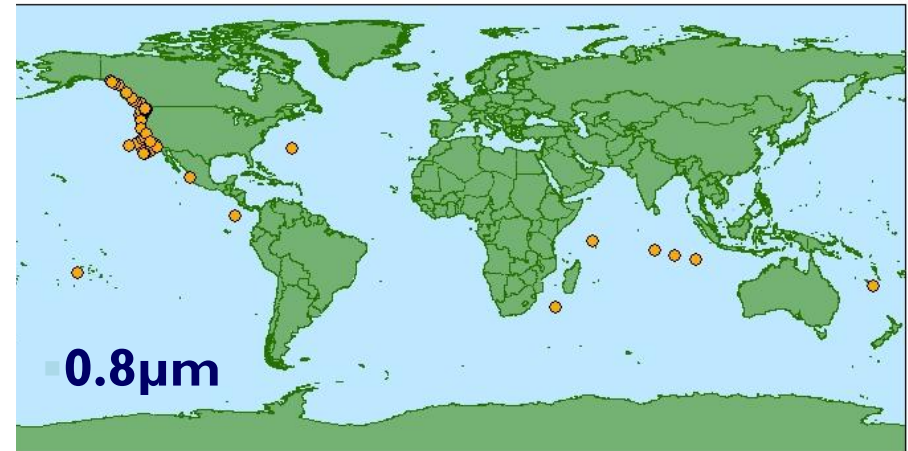2003 – 2008 Routes    2009 – 2010 Route

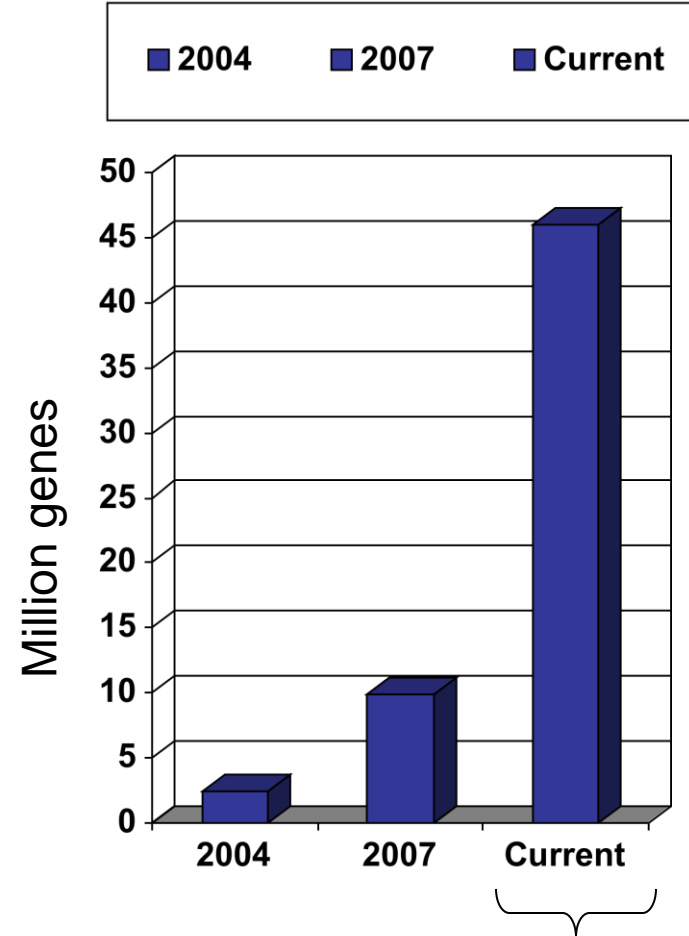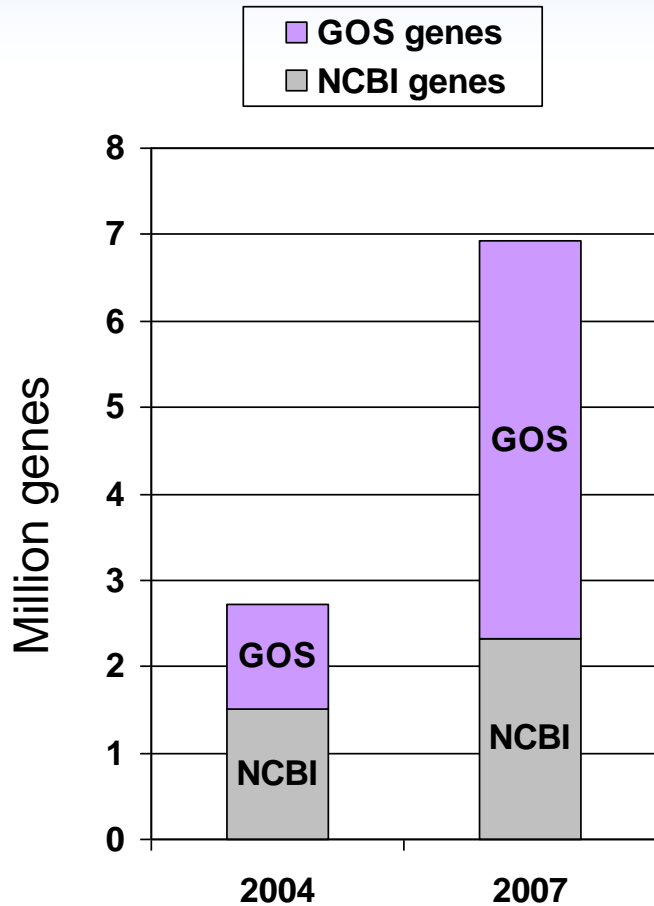J. Craig Venter
INSTITUTE

# Sample Filtration

# GOS circumnavigation data
# 229 stations and 291 samples



0.1μm

0.8μm

viral

3.0μm

# GOS data

| | Reads | Proteins | Sequencing Technology |
|---|---|---|---|
| Phase I | 7.6 Million | 9.8 Million | Sanger |
| Circumnavigation | 48 Million | ~53Million | Sanger + 454 |

J. Craig Venter

I N S T I T U T E

# GOS dataset is expanding the protein universe



Extrapolation based on amount of GOS sequence data currently available but not yet released to public domain

J. Craig Venter
INSTITUTE

# Community makeup

# Taxonomic makeup of GOS samples based on 16S data from shotgun sequencing

| Phylum or Class | Fraction[a] |
|---|---|
| Alpha *Proteobacteria* | 0.32 |
| Unclassified *Proteobacteria* | 0.155 |
| Gamma *Proteobacteria* | 0.132 |
| *Bacteroidetes* | 0.13 |
| *Cyanobacteria* | 0.079 |
| *Firmicutes* | 0.075 |
| *Actinobacteria* | 0.046 |
| Marine Group A | 0.022 |
| Beta *Proteobacteria* | 0.017 |
| OP11 | 0.008 |
| Unclassified *Bacteria* | 0.008 |
| Delta *Proteobacteria* | 0.005 |
| *Planctomycetes* | 0.002 |
| Epsilon *Proteobacteria* | 0.001 |

[a]Values shown are averages over all samples.

J. Craig Venter

I N S T I T U T E

# Phylogenetic Distribution in the Indian Ocean across size-classes
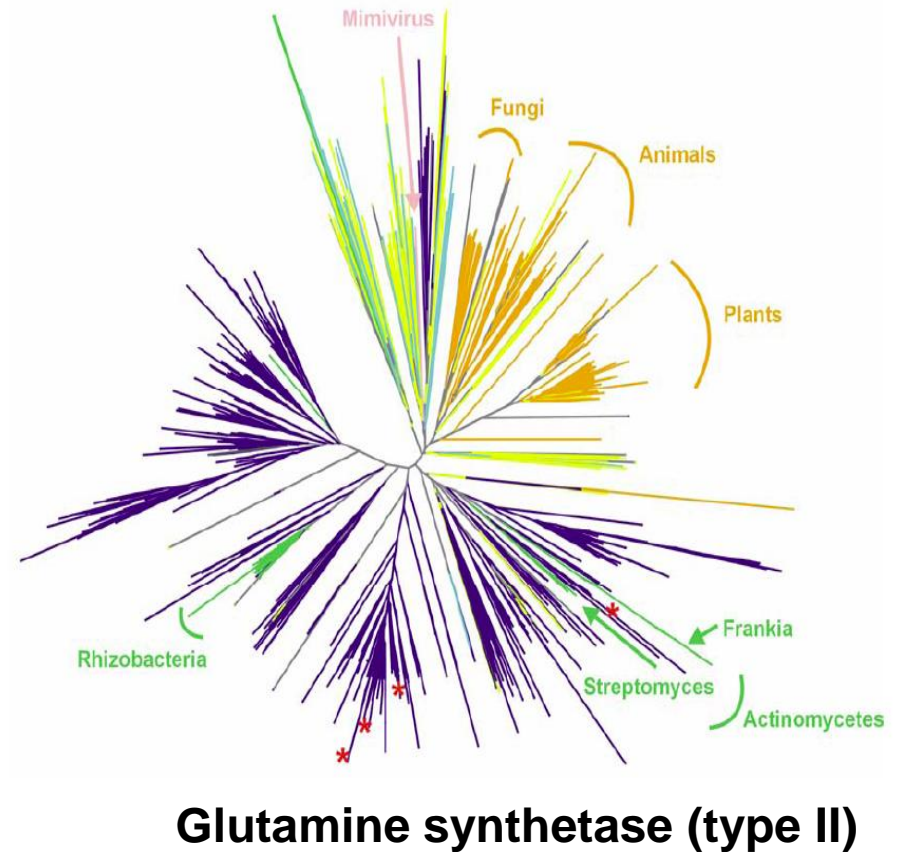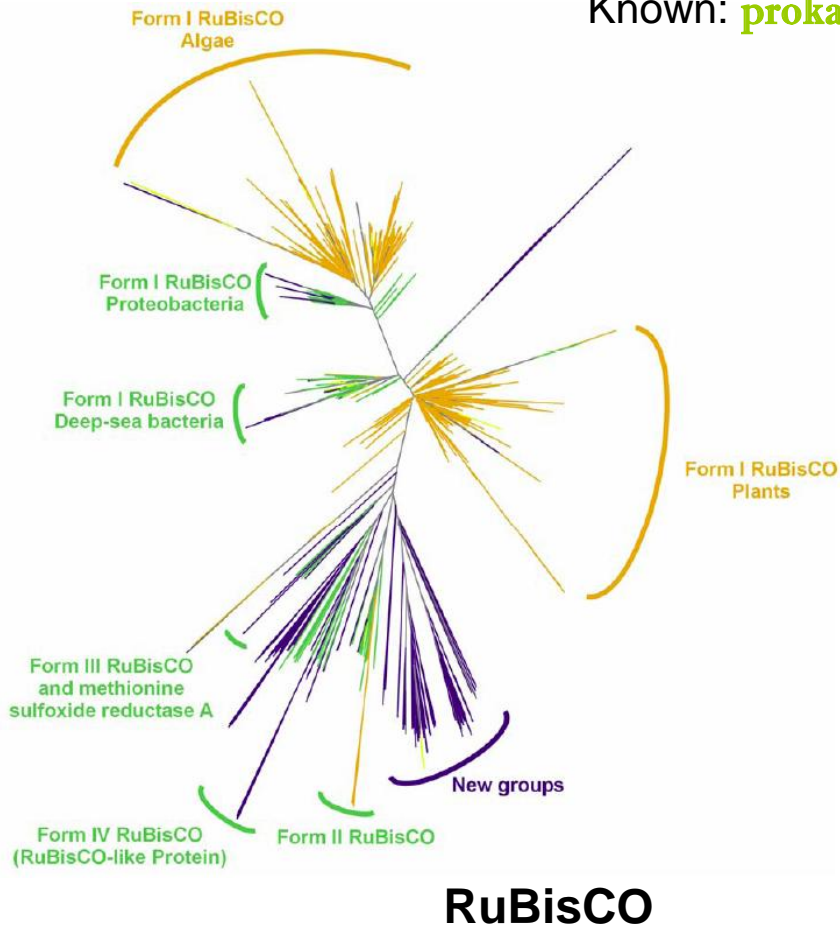
# GOS increases size and diversity of known protein families



GOS: prokaryotes, eukaryotes
Known: prokaryotes, eukaryotes

**RuBisCO**

**Glutamine synthetase (type II)**
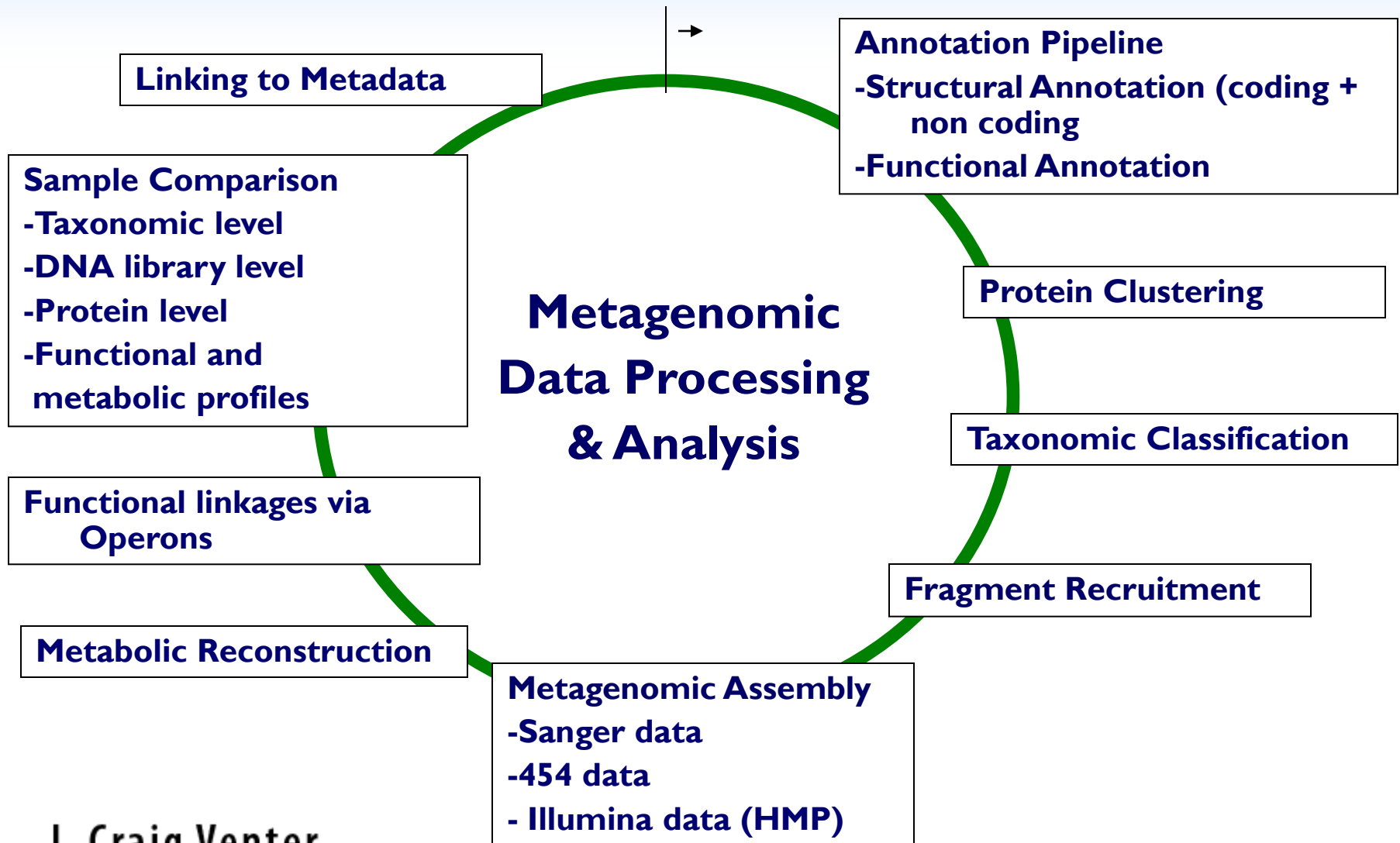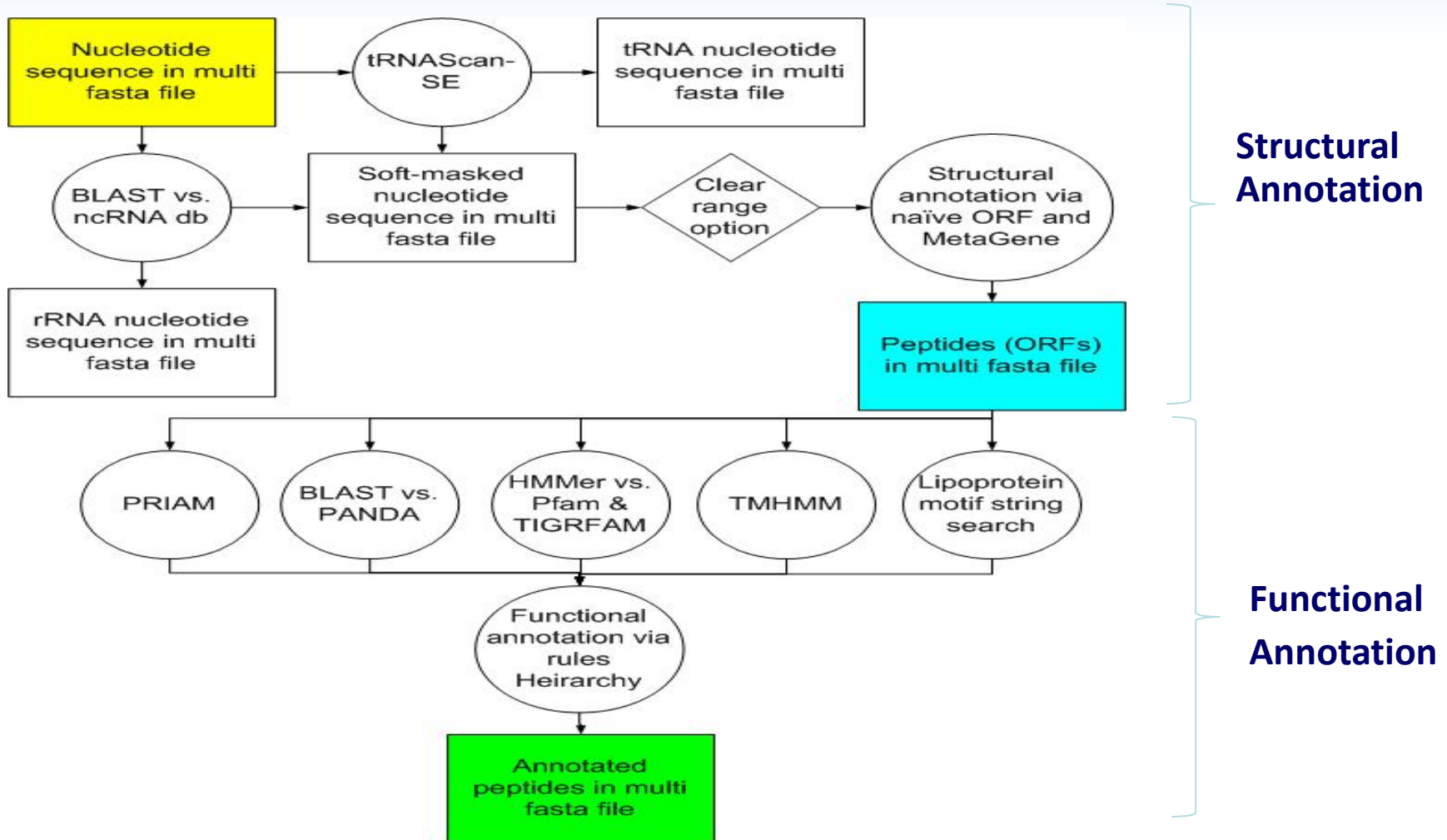
J. Craig Venter
INSTITUTE

# Viruses in the Marine Environment

- Abundant: ~$10^7$ /ml$^{-1}$ of surface seawater

- Diverse: VBR $\cong$ 10 ; ~ 10-fold greater diversity than microbial hosts

- Influence  microbial diversity through infection and host cell lysis

- Mediators of horizontal gene transfer

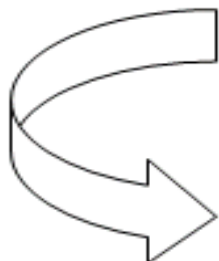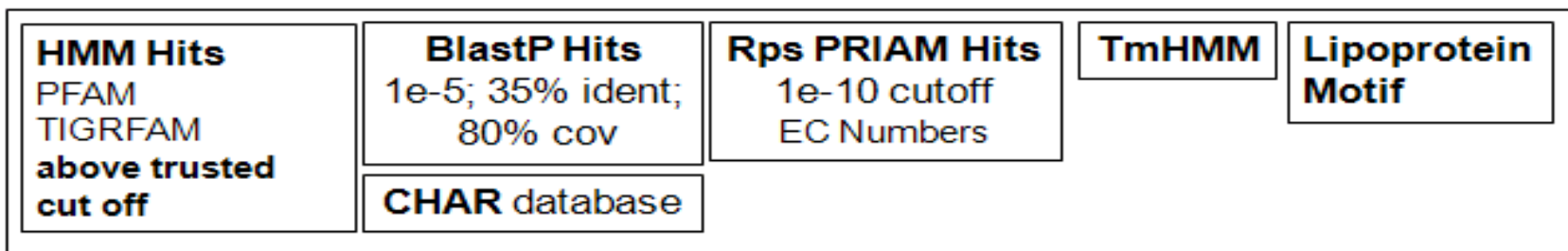- Influence biogeochemical cycling, particularly carbon

J. Craig Venter

I N S T I T U T E

# High-throughput Metagenomic Data Analysis

**Linking to Metadata**

**Annotation Pipeline**

-**Structural Annotation (coding + non coding**

-**Functional Annotation**

**Sample Comparison**

-**Taxonomic level**

-**DNA library level**

-**Protein level**

-**Functional and metabolic profiles**

## Metagenomic Data Processing & Analysis

**Protein Clustering**

**Taxonomic Classification**

**Functional linkages via Operons**

**Fragment Recruitment**

**Metabolic Reconstruction**

**Metagenomic Assembly**

-**Sanger data**

-**454 data**

- **Illumina data (HMP)**

J. Craig Venter

I N S T I T U T E

# Metagenomic Data Processing - Annotation pipeline



**Structural Annotation**

**Functional Annotation**

**Published in SIGS**

J. Craig Venter

INSTITUTE

# Annotation Rules Hierarchy

**Evidences**

| HMM Hits<br>PFAM<br>TIGRFAM<br>**above trusted**<br>**cut off** | **BlastP Hits**<br>1e-5; 35% ident;<br>80% cov<br><br>**CHAR** database | **Rps PRIAM Hits**<br>1e-10 cutoff<br>EC Numbers | **TmHMM** | **Lipoprotein**<br>**Motif** |
|---|---|---|---|---|

**Annotation Rules**

1. TIGRFAM/PFAM             (Equivalog)
2. Characterized (CHAR) BlastP Hit
3. TIGRFAM/PFAM             (Non-Equivalog)
4. CDP (conserved domain protein) blastp hit
5. TmHMM hit: "membrane protein"
6. Lipoprotein motif: "lipoprotein"
7. "hypthetical protein"

**Common Names, Gene Symbols, EC Numbers, GO Terms, TIGR Role ids**

# Viral Metagenomic (functional)Pipeline

J. Craig Venter
I N S T I T U T E

# Annotation Rules Hierarchy (Viral)

- PFAM/TIGRFAM_HMM, equivalog above trusted cutoff
- ACLAME_PEP, %id>= 50, coverage >= 80, e-value <= $10^{-10}$
- ALLGROUP_PEP, %id>= 50, coverage >= 80, e-value <= $10^{-10}$
- ACCLAME_HMM matches, > 90% coverage, e-value < $10^{-5}$
- PFAM/TIGRFAM_HMM, non-equivalog above trusted cutoff
- CDD_RPS, %id>= 35%, coverage >= 90% of CDD-domain, e-value <= $1e^{-10}$
- FRAG_HMM, e-value < $1e^{-5}$
- ACLAME_PEP, %id >= 30%, coverage >= 70%, e-value <= $1e^{-5}$
- ALLGROUP_PEP, %id >= 30%, coverage >= 70%, e-value <= $1e^{-5}$
- No evidence -> hypothetical protein

J. Craig Venter

I N S T I T U T E

# Metagenomic Assembly

## Advantages

- Provides genomic context
- Reduces redundancy and complexity
- Improves annotation
- Mechanism to isolate environment specific gene regions

## Challenges

- Coverage dependent
- Variation can limit the length of assemblies
- Can mask diversity

- Celera Hybrid Assembler has been updated to work with 454 Titanium reads
- Will further optimize assembly process to capture environmental diversity

# Metagenomic Data Processing - Continued

- **Protein Clustering** : JCVI's Protein clustering (S. Yooseph)
- **Taxonomic Classification** : APIS (J. Badger)
- **Fragment Recruitment** :Advanced Reference Viewer (D. Rusch)
- **Metagenomic Assembly** : Celera Assembler (G. Sutton & J. Miller)
- **Sample Comparison**

Making sense of everything in the context of **METADATA**

J. Craig Venter
I N S T I T U T E

# General Questions

- Who are they?

  Species , Taxonomic distribution…

- How many?

  Distribution across sites and filters

- What are they doing?

  Functional profiles

  **Metabolic profiles**

J. Craig Venter

I N S T I T U T E

# MR Specific Questions

- Metabolic profiles across sites and filters

- Pathways coverage and abundance

- What known characterized pathways and how many?

- What novel pathways are there?

- Metabolic network

J. Craig Venter

I N S T I T U T E

# Metabolic Reconstruction

- From the Annotation Pipeline (orf based)

Proteins → EC assignment → Pathways prediction
(EC to MetaCyc/Kegg mapping)

**Sources for EC :** TIGRFAM
PFAM
High confidence blast hit to Uniref100/Panda
RPSblast to EC profiles from PRIAM

- From BlastX to a Functional database (read based)

Reads → Blastx Metacyc/Kegg → Pathways prediction

J. Craig Venter
I N S T I T U T E

# Browse/analyze/compare pathways across datasets in the context of annotation and Metadata

**METAREP**

JCVI Metagenomics Reports

website **www.jcvi.org/metarep**
source code **http://github.com/jcvi/METAREP**
blog **http://blogs.jcvi.org/tag/metarep**
contact **metarep-support@jcvi.org**

**METAREP is a web interface designed to help scientists to view, query and compare annotation data derived from proteins called on metagenomics reads**

**Developer : Johannes Goll**
**Published in Bioinformatics**

**www.jcvi.org/metarep**

J. Craig Venter
INSTITUTE

# Browse pathways

METAREP
JCVI Metagenomics Reports

QUICK NAVIGATION | SEARCH | LIST PROJECTS | LIST POPULATIONS | PIPELINE LOG | DASHBOARD | LOG OUT

List Projects | View Project | Browse Dataset

## Browse MetaCyc Pathways gos-phase-I-sanger (GOS Phase I)

### Filter

Help  `*.*`    Filter

### Browse MetaCyc Pathways

- Biosynthesis (level 1) [7,924,211 hits]
  - Amino Acids Biosynthesis (level 2) [1,528,954 hits]
  - Aminoacyl-tRNA Charging (level 2) [386,014 hits]
  - Aromatic Compounds Biosynthesis (level 2) [141,820 hits]
    - 3-Dehydroquinate Biosynthesis (level 2) [11,472 hits]
    - 4-Hydroxybenzoate Biosynthesis (level 2) [24,582 hits]
    - Chorismate Biosynthesis (level 2) [35,857 hits]
    - 1,3,5-trimethoxybenzene biosynthesis (level 2) [0 hits]
    - oligomeric urushiol biosynthesis (pathway) [965 hits]
    - 4-hydroxyphenylpyruvate biosynthesis (pathway) [27 hits]
    - 2,3-dihydroxybenzoate biosynthesis (pathway) [1,094 hits]
    - phaselate biosynthesis (pathway) [12,041 hits]
    - benzoylanthranilate biosynthesis (pathway) [0 hits]
    - salicylate biosynthesis I (pathway) [8,780 hits]
    - <i>trans</i>-cinnamoyl-CoA biosynthesis (pathway) [59,838 hits]
    - benzoyl-CoA biosynthesis (level 2) [0 hits]
    - petivericin biosynthesis (level 2) [0 hits]
    - salicylate biosynthesis II (pathway) [11,461 hits]
  - Carbohydrates Biosynthesis (level 2) [728,419 hits]
  - Cell Structures Biosynthesis (level 2) [1,213,328 hits]
  - Cofactors, Prosthetic Groups, Electron Carriers Biosynthesis (level 2) [1,897,167 hits]
  - Hormones Biosynthesis (level 2) [2,113,243 hits]
  - Fatty Acids and Lipids Biosynthesis (level 2) [3,883,182 hits]
  - Metabolic Regulators Biosynthesis (level 2) [28,729 hits]
  - Nucleosides and Nucleotides Biosynthesis (level 2) [411,479 hits]
  - Other Biosynthesis (level 2) [287,519 hits]
  - Amines and Polyamines Biosynthesis (level 2) [384,516 hits]
  - Secondary Metabolites Biosynthesis (level 2) [3,369,338 hits]
  - Siderophore Biosynthesis (level 2) [9,041 hits]
  - <i>Methanobacterium thermoautotrophicum</i> biosynthetic metabolism (pathway) [220,03
- Degradation/Utilization/Assimilation (level 1) [4,350,781 hits]
- Detoxification (level 1) [575,806 hits]
- Generation of precursor metabolites and energy (level 1) [2,156,128 hits]
- Signal transduction pathways (level 1) [0 hits]
- Superpathways (level 1) [6,376,473 hits]
- Transport (level 1) [198,282 hits]
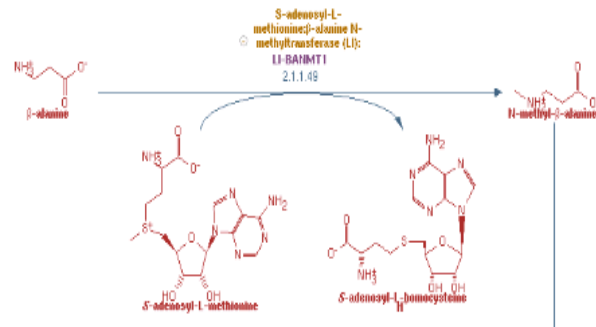- tRNA processing pathway (pathway) [0 hits]

### Pathway Classification



Aromatic Compounds Biosynthesis

MetaCyc Pathway: β-alanine betaine biosynthesis

Less Detail    Species Comparison

J. Craig Venter
INSTITUTE

# Compare pathways across datasets

# Pathways Tools for GOS

- Metagenomic specific predictions - Incorporate taxonomic resolution when predicting pathways

- Confidence Scores for the pathways

- Incorporate more annotation evidence types in predictions other than EC

- Ability to overlay and visualize expression data

- Full integration of pathways tools into Metarep

- Performance enhancements to handle metagenomic data volume

# Conclusion

- Who are they?

    Species , Taxonomic distribution…

- How many?

    Distribution across sites and filters

- What are they doing?

    Functional profiles

    **Metabolic profiles**

# Acknowledgements

J. Craig Venter

I N S T I T U T E

# Questions


## Thank You

J. Craig Venter
I N S T I T U T E