# A computationally intensive approach to discover novel adaptation genes in extreme environments

Genomics of Gene Expression Lab

Tatyana Zamkovaya[1], Jamie S. Foster[2], Valérie de Crecy-Lágard[1,3], Ana Conesa[1,3]

[1] Microbiology and Cell Science Department, Institute for Food and Agricultural Sciences, University of Florida, Gainesville, FL, USA
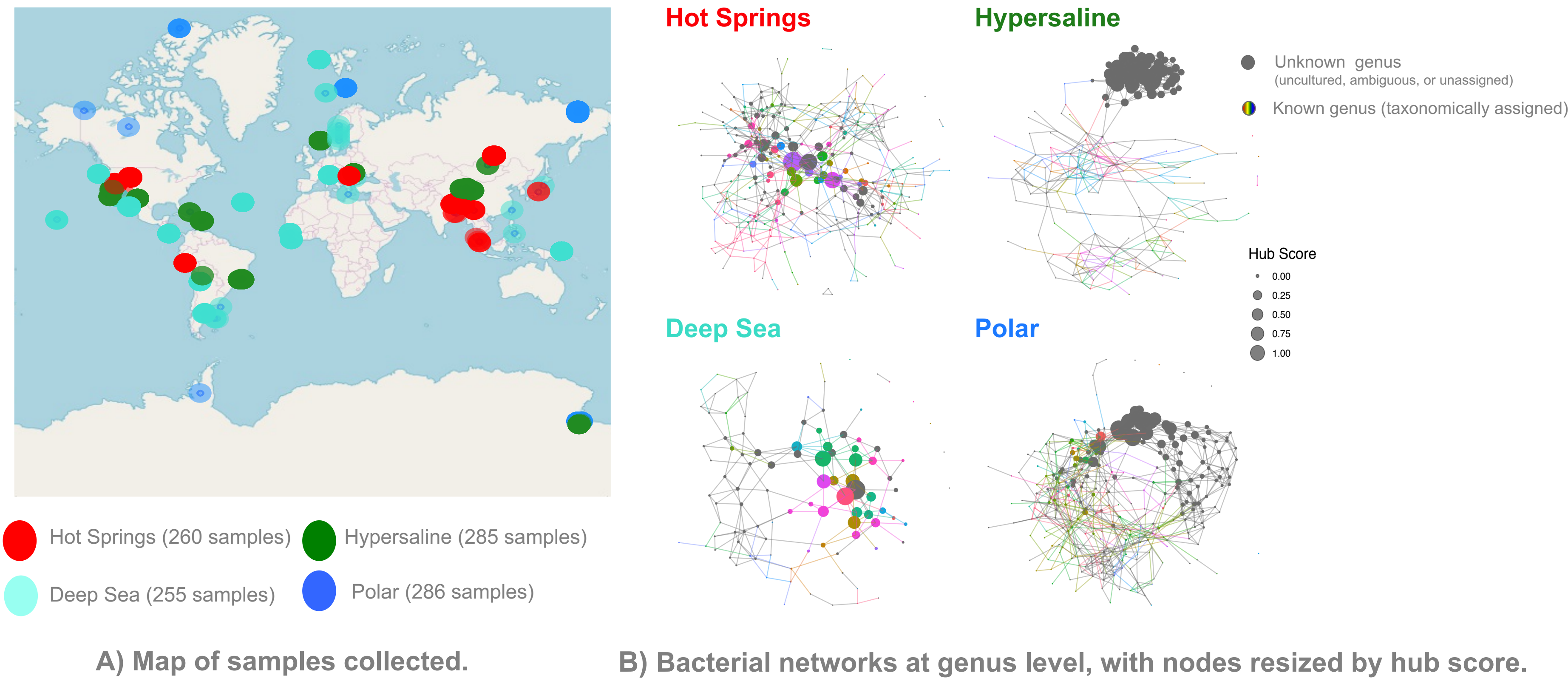[2] Department of Microbiology and Cell Science, Space Life Sciences Lab, Merritt Island, FL, USA
[3] Genetics and Genomics Institute, University of Florida, Gainesville, FL, USA

UF UNIVERSITY of FLORIDA

## Introduction

**Unknown taxa act as top hubs in extreme environmental networks, yet have unknown functional roles**



A) Map of samples collected.

- Hot Springs (260 samples)
- Hypersaline (285 samples)
- Deep Sea (255 samples)
- Polar (286 samples)

Hot Springs | Hypersaline | Deep Sea | Polar

- Unknown genus (uncultured, ambiguous, or unassigned)
- Known genus (taxonomically assigned)

Hub Score: 0.00, 0.25, 0.50, 0.75, 1.00

B) Bacterial networks at genus level, with nodes resized by hub score.
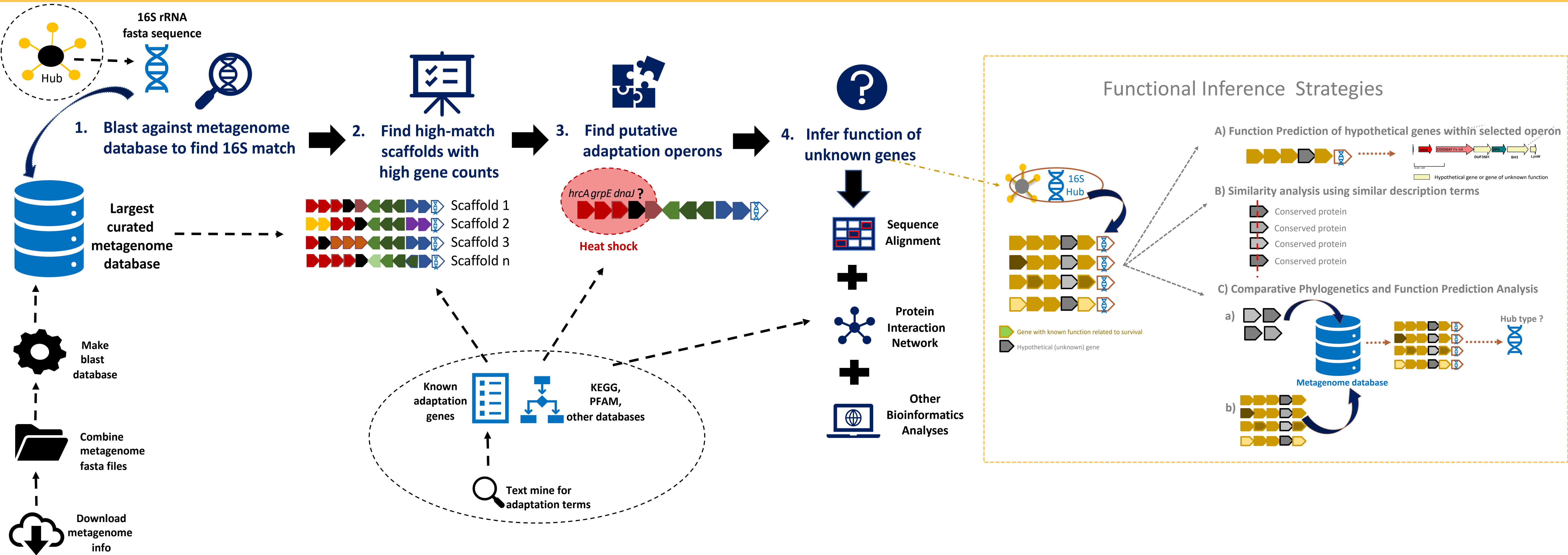
## Hypothesis and Aim

The key to understand how life evolves and adapts to extreme conditions, within and beyond Earth, lies in the uncharacterized, unknown taxa, called Microbial Dark Matter (MDM), that dominate even the most extreme of environments. We were able to identify unknown hubs of extreme environments using a network-based approach on amplicon data. We found that unknowns are highly prevalent and connected and wanted to know more about how and why these unknowns were able to adapt to such harsh conditions. Building on these results, we now propose a computational approach to identify adaptation-related genes in these unknown organisms.

Due to their integral network position, these extreme environmental hubs may possess optimized adaptation strategies and therefore are ideal candidates to search for novel biological pathways of survival.

Lack of well-annotated reference genomes prevents existing tools from predicting gene function of unknown taxa.

Here, we circumvent this problem by using **16S hubs as probes** to **find novel genes** within **functionally annotated metagenome scaffolds**.

## Hub Blast Pipeline



1. Blast against metagenome database to find 16S match
2. Find high-match scaffolds with high gene counts
3. Find putative adaptation operons
4. Infer function of unknown genes

Largest curated metagenome database

Make blast database
Combine metagenome fasta files
Download metagenome info

Scaffold 1 / Scaffold 2 / Scaffold 3 / Scaffold n

*hrcA grpE dnaJ* ? — Heat shock

Known adaptation genes — KEGG, PFAM, other databases
Text mine for adaptation terms

Sequence Alignment + Protein Interaction Network + Other Bioinformatics Analyses

Functional Inference Strategies
A) Function Prediction of hypothetical genes within selected operon
B) Similarity analysis using similar description terms
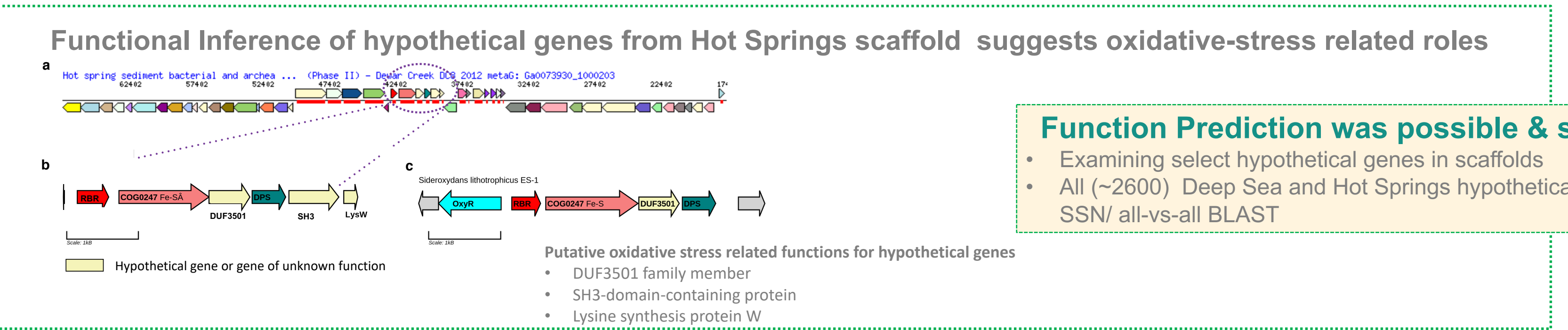C) Comparative Phylogenetics and Function Prediction Analysis

## Results

**Approach successfully returns high-confidence gene-rich scaffolds for known and unknown hubs belonging to different environments and enables prediction of novel adaptation genes**
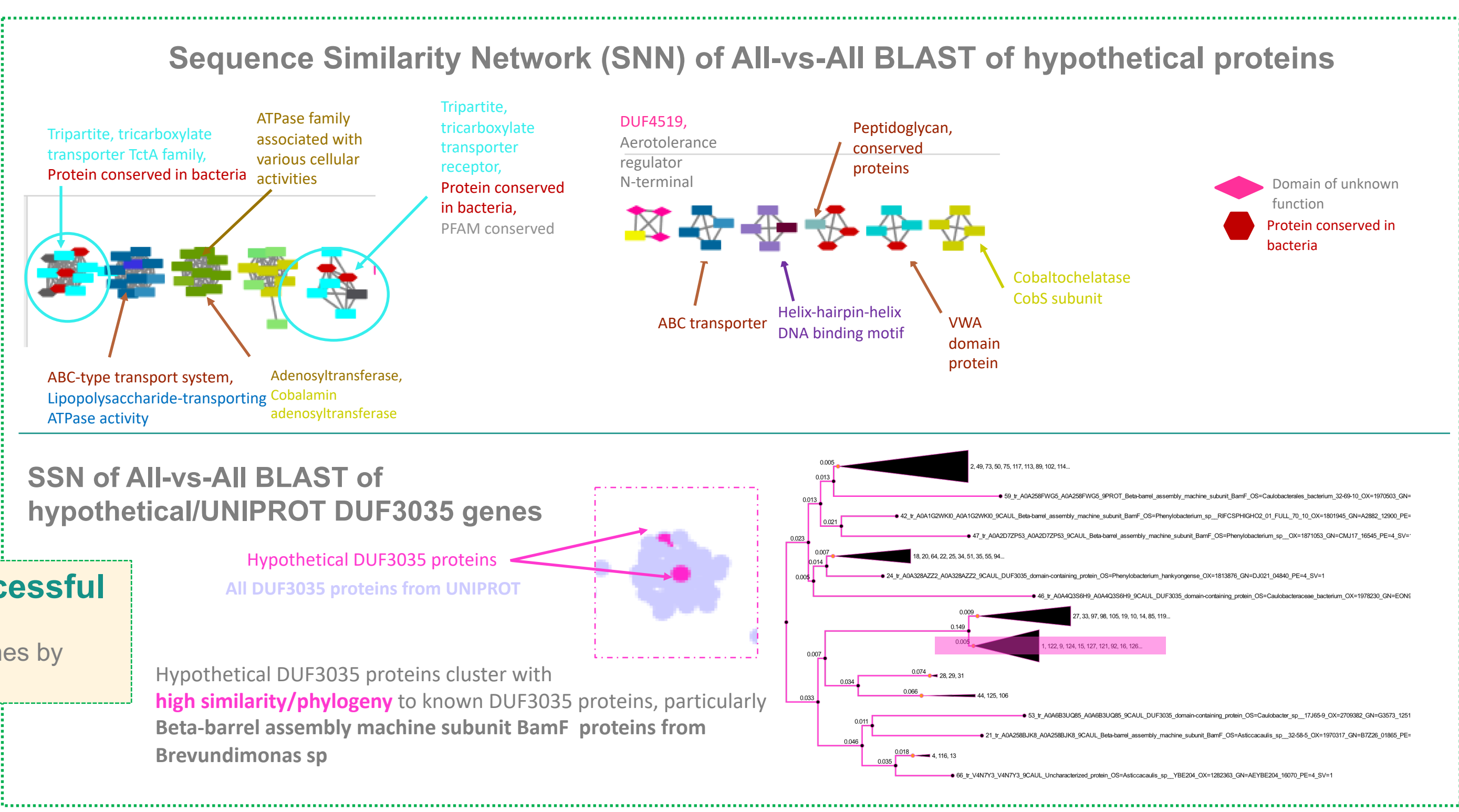
### Total Gene, Hypothetical Gene, and Adaptation Operon Counts

| | met_hub_type | avg_gene_count | avg_percent_hyp_genes | avg_percent_adapt_genes | total_gene_count | total_hyp_count | total_adapt_count | sum_operon_counts | mean_operon_counts |
|---|---|---|---|---|---|---|---|---|---|
| 1 | HS_Known | 238.1429 | 17.78571 | 8.585714 | 3334 | 683 | 254 | 75 | 2.678571 |
| 2 | HS_Unknown | 345.6596 | 17.09362 | 8.212766 | 16246 | 2985 | 1396 | 344 | 3.659574 |
| 3 | HY_Known | 102.9524 | 15.41429 | 8.552381 | 2162 | 335 | 201 | 52 | 1.238095 |
| 4 | HY_Unknown | 34.0000 | 11.80000 | 8.800000 | 34 | 4 | 3 | 2 | 1.000000 |
| 5 | DS_Known | 236.8462 | 19.77308 | 7.980769 | 6158 | 1191 | 520 | 174 | 3.346154 |
| 6 | DS_Unknown | 144.8750 | 15.92500 | 8.318750 | 2318 | 406 | 190 | 72 | 2.250000 |
| 7 | PO_Known | 198.6000 | 19.85000 | 8.050000 | 1986 | 367 | 170 | 56 | 2.800000 |
| 8 | PO_Unknown | 70.5000 | 28.80000 | 7.100000 | 141 | 41 | 8 | 2 | 0.500000 |

**High # of hypothetical genes, adaptation genes, and adaptation-related operons were found**

### Functional Inference of hypothetical genes from Hot Springs scaffold suggests oxidative-stress related roles



**Function Prediction was possible & successful**
- Examining select hypothetical genes in scaffolds
- All (~2600) Deep Sea and Hot Springs hypothetical genes by SSN/ all-vs-all BLAST

Putative oxidative stress related functions for hypothetical genes
- DUF3501 family member
- SH3-domain-containing protein
- Lysine synthesis protein W

- Hypothetical gene or gene of unknown function

### Sequence Similarity Network (SNN) of All-vs-All BLAST of hypothetical proteins



- Domain of unknown function
- Protein conserved in bacteria

### SSN of All-vs-All BLAST of hypothetical/UNIPROT DUF3035 genes



- Hypothetical DUF3035 proteins
- All DUF3035 proteins from UNIPROT

Hypothetical DUF3035 proteins cluster with **high similarity/phylogeny** to known DUF3035 proteins, particularly Beta-barrel assembly machine subunit BamF proteins from *Brevundimonas sp*

## Conclusions

- **Our approach enables identification and prediction of novel genes**
- Our pipeline **accommodates the computational power and memory required of large-scale datasets**
- Function prediction of select operons shows conserved roles and similar functional categories among neighboring genes
- Large-scale function prediction of hypothetical genes via SSN shows high similarity and potential conservation among genes with similar Gene Ontology annotations, even among genes originating from different environmental metagenomes.

## Future Work

Similarity analysis of hypothetical genes and adaptation operons may help identify novel but well-conserved survival/ stress genes

## Acknowledgements

## References

Chen, I.-M. A., Chu, K., Palaniappan, K., Pillay, M., Ratner, A., Huang, J., ... Kyrpides, N. C. (2019). IMG/M v.5.0: an integrated data management and comparative analysis system for microbial genomes and microbiomes. *Nucleic Acids Research*, 47(D1), D666–D677. https://doi.org/10.1093/nar/gky901
Conesa, A., Gotz, S., Garcia-Gomez, J. M., Terol, J., Talon, M., & Robles, M. (2005). Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics*, 21(18), 3674–3676. https://doi.org/10.1093/bioinformatics/bti610