

# Peptides associated with scaffold-derived metaproteomic proteins from samples taken during R/V Atlantic Explorer cruise AE1913 from the Sargasso Sea to Northeast US shelf waters in June of 2019

**Website:** <https://www.bco-dmo.org/dataset/934718>

**Data Type:** Cruise Results

**Version:** 1

**Version Date:** 2024-08-01

## Project

» [Collaborative Research: Direct Characterization of Adaptive Nutrient Stress Responses in the Sargasso Sea using Protein Biomarkers and a Biogeochemical AUV](#) (Nutrient Stress Responses and AUV Clio)

Contributors	Affiliation	Role
<a href="#">Saito, Mak A.</a>	Woods Hole Oceanographic Institution (WHOI)	Principal Investigator
<a href="#">Cohen, Natalie</a>	Woods Hole Oceanographic Institution (WHOI)	Scientist
<a href="#">York, Amber D.</a>	Woods Hole Oceanographic Institution (WHOI BCO-DMO)	BCO-DMO Data Manager

## Abstract

These are individual peptides associated with scaffold-derived metaproteomic proteins (includes total spectral counts for each peptide). Samples were taken during R/V Atlantic Explorer cruise AE1913 in Subtropical North Atlantic, beginning at the Bermuda Atlantic Time-series Station (BATS) of the Sargasso Sea and ending in coastal Northeast US shelf waters in June of 2019.

## Table of Contents

- [Coverage](#)
- [Dataset Description](#)
  - [Methods & Sampling](#)
  - [Data Processing Description](#)
  - [BCO-DMO Processing Description](#)
- [Data Files](#)
- [Related Publications](#)
- [Related Datasets](#)
- [Parameters](#)
- [Instruments](#)
- [Deployments](#)
- [Project Information](#)
- [Funding](#)

## Coverage

**Location:** Subtropical North Atlantic, beginning at the Bermuda Atlantic Time-series Station (BATS) of the Sargasso Sea (31.666888 N64.166293 W) and ending in coastal Northeast US shelf waters (39.31658 N 71.123208 W)

**Spatial Extent:** N:38.527393 E:-64.165587 S:31.586387 W:-70.840703

**Temporal Extent:** 2019-06-16 - 2019-06-27

## Dataset Description

Related data table and dataset descriptions:

Total spectral counts refer to the total number of spectra with peptide to spectrum matches (PSMs) that matches to each entry within the FASTA sequence database. This approach allows each peptide to map to multiple closely related sequences. In contrast, with exclusive spectral counts each peptide is only allowed to

map to one sequence within the FASTA database, and when a peptide is found in multiple database sequences the one with the most peptides mapping (parsimony) to it is selected. There are pros and cons to each approach, where total spectral counts will double count peptides when two similar proteins are compared, and exclusive spectral counts will underrepresent less abundant proteins with shared peptides, favoring the most homolog with the most shared peptides. Considering protein groups with shared peptides or focusing on peptide-level analyses are alternative approaches that could be constructed from these results.

See "Related Datasets" section for:

\* "AE1913 Protein Spectral Counts" which includes the exclusive and total spectral counts associated with proteins.

\* "AE1913 Protein Identification FASTA"

CTD and other data from the same cruise are listed on deployment page AE1913: <https://www.bco-dmo.org/deployment/916412>

These data will become part of the Ocean Protein Portal (<https://proteinportal.whoi.edu/>; Saito et al., 2020).

## Methods & Sampling

Methods are reported in Cohen et al. 2023 (biorxiv preprint doi: [10.1101/2023.11.20.567900](https://doi.org/10.1101/2023.11.20.567900)) and are summarized below.

\* This section describes how this and related datasets were generated (see "Related Datasets" section).

One half of the 142 mm filters (0.2-51  $\mu\text{m}$ ) collected by Clio were processed for metaproteomics. Proteins were extracted in an 1% SDS-based detergent in 50 mM HEPES at pH 8.5, reduced with dithiothreitol, alkylated with iodoacetamide, and purified using a polyacrylamide electrophoresis tube gel method. Protein quantification was performed using a BSA assay. Trypsin was added to the protein-bead mixture in a 1:20 trypsin:protein ratio. Peptides were purified using C18 tips and diluted to a concentration of 0.1  $\mu\text{g } \mu\text{L}^{-1}$ .

Approximately 2-5  $\mu\text{g}$  of purified peptides were injected onto a Dionex UltiMate 3000 RSLCnano LC system with an additional RSLCnano pump, run in online 2D active modulation mode interfaced with a Thermo Fusion mass spectrometer. The mass spectrometer acquired MS1 scans from 380 to 1,580 m/z at 240K resolution in the Orbitrap. MS2 were collected in data dependent mode in the ion trap with a cycle time of 2 seconds between scans and acquisition of charge states 2 to 10. MS2 scans had 1.6 m/z isolation window, 50 ms maximum injection time and 5 s dynamic exclusion time.

## Data Processing Description

The metatranscriptomic ORFs were used as the protein database, and peptide-spectrum matches were performed using Sequest algorithm within IseNode Proteome Discoverer 2.2.0.388 with a parent ion tolerance of 10 ppm and fragment tolerance of 0.6 Da, and 0 max missed cleavage. Identification criteria consisted of a peptide threshold of 98% (false discovery rate [FDR] = 0.1%) and protein threshold of 99% (1 peptide minimum, FDR = 1.5%) in Scaffold 5.1.2 (Proteome Software) resulting in 77,438 proteins and 3,155,061 exclusive spectral counts.

"NA" in the spectral\_count\_sum column represents annotations not available, or no counts detected.

## BCO-DMO Processing Description

\* Submitted file "opp\_peptides.csv" was modified to name the first un-named column with name "row\_id". Data were not modified further. Due to table size the data were not imported into the data system and instead was attached to this dataset as a csv a Data File directly named .

\*\* In the BCO-DMO data system missing data identifiers are displayed according to the format of data you access. For example, in csv files it will be blank (null) values. In Matlab .mat files it will be NaN values. When viewing data online at BCO-DMO, the missing value will be shown as blank (null) values.

\* Column names adjusted to conform to BCO-DMO naming conventions designed to support broad re-use by a variety of research tools and scripting languages. [Only numbers, letters, and underscores. Can not start

with a number]

\* Date converted to ISO 8601 format

\* ISO DateTime with timezone (UTC) column added in ISO 8601 format.

\* Species names corrected. World Register of Marine Species taxa match tool used to find misspellings and unaccepted names. Name changes reviewed and accepted by the data contributor Species name spelling corrected to the accepted spelling (as of 2018-05-09)

[ [table of contents](#) | [back to top](#) ]

---

## Data Files

File
<b>Peptide information and spectral counts</b> filename: 934718_v1_ae1913-peptide-spectral-counts.csv (Comma Separated Values (.csv), 11.49 GB) MD5:3d20b0d47a2867c52ceaff32e5cb9f77 Primary data table for dataset 934718 version 1.

[ [table of contents](#) | [back to top](#) ]

---

## Related Publications

Saito, M. A., Saunders, J. K., Chagnon, M., Gaylord, D. A., Shepherd, A., Held, N. A., Dupont, C., Symmonds, N., York, A., Charron, M., & Kinkade, D. B. (2020). Development of an Ocean Protein Portal for Interactive Discovery and Education. *Journal of Proteome Research*, 20(1), 326–336.

<https://doi.org/10.1021/acs.jproteome.0c00382>

*Related Research*

[ [table of contents](#) | [back to top](#) ]

---

## Related Datasets

### IsRelatedTo

Cohen, N., Krinos, A., Alexander, H., & Saito, M. (2022). Protistan metabolism across the western North Atlantic Ocean revealed through autonomous underwater profiling (Version 2) [Data set]. Zenodo.  
<https://doi.org/10.5281/ZENODO.8287779> <https://doi.org/10.5281/zenodo.8287779>

Saito, M. A., Cohen, N. (2024) **Protein identification FASTA file (scaffold-derived metaproteomic proteins) from samples taken during R/V Atlantic Explorer cruise AE1913 from the Sargasso Sea to Northeast US shelf waters in June of 2019.** Biological and Chemical Oceanography Data Management Office (BCO-DMO). (Version 1) Version Date 2024-08-01 doi:10.26008/1912/bco-dmo.934727.1 [[view at BCO-DMO](#)]

*Relationship Description: These datasets are from the same collection and study and will be included in the Ocean Protein Portal (<https://proteinportal.whoi.edu>).*

Saito, M. A., Cohen, N. (2024) **Scaffold-derived metaproteomic exclusive and total spectral counts associated with proteins from samples taken during R/V Atlantic Explorer cruise AE1913 from the Sargasso Sea to Northeast US shelf waters in June of 2019.** Biological and Chemical Oceanography Data Management Office (BCO-DMO). (Version 1) Version Date 2024-08-01 doi:10.26008/1912/bco-dmo.934706.1 [[view at BCO-DMO](#)]

*Relationship Description: These datasets are from the same collection and study and will be included in the Ocean Protein Portal (<https://proteinportal.whoi.edu>).*

[ [table of contents](#) | [back to top](#) ]

---

## Parameters

Parameter	Description	Units
row_id	sequential row identifier	unitless
protein_id	Identification parameter for protein sequence, typically shared with metagenomic database. An identifier that uniquely identifies this protein within this dataset and the FASTA file.	unitless
peptide_sequence	Peptide genomic sequence (unique sequence in dataset)	unitless
sample_id	Identifies the sample associated with this annotation	unitless
spectral_count_sum	Sum of spectral counts (+2, +3, +4 ions)	unitless
cruise_id	Cruise identifier	unitless
station_id	Station identifier where sample was taken	unitless
depth_m	The depth in meters at which the sample as taken	meters
minimum_filter_size_microns	Minimum size of the collection filter	microns (um)
maximum_filter_size_microns	Maximum size of the collection filter	microns (um)
date_y_m_d	The date of sample collection	unitless
latitude_dd	The latitude at the station in decimal degrees (-90 to 90)	decimal degrees
longitude_dd	The longitude at the station in decimal degrees (-180 to 180)	decimal degrees

[ [table of contents](#) | [back to top](#) ]

---

## Instruments

<b>Dataset-specific Instrument Name</b>	
<b>Generic Instrument Name</b>	AUV Clio
<b>Generic Instrument Description</b>	Clio is an autonomous underwater vehicle (AUV) created to accomplish the dual goals of global ocean mapping and biochemistry sampling. The ability to sample dissolved and particulate seawater biochemistry across ocean basins while capturing fine-scale biogeochemical processes sets it apart from other AUVs. Clio is designed to efficiently and precisely move vertically through the ocean, drift laterally to observe water masses, and integrate with research vessel operations to map large horizontal scales up to a depth of 6,000 meters. More information is available at <a href="https://www2.whoi.edu/site/deepsubmergencelab/cliol/">https://www2.whoi.edu/site/deepsubmergencelab/cliol/</a>

<b>Dataset-specific Instrument Name</b>	Thermo Fusion mass spectrometer
<b>Generic Instrument Name</b>	Mass Spectrometer
<b>Generic Instrument Description</b>	General term for instruments used to measure the mass-to-charge ratio of ions; generally used to find the composition of a sample by generating a mass spectrum representing the masses of sample components.

<b>Dataset-specific Instrument Name</b>	Dionex UltiMate 3000 RSLCnano LC system
<b>Generic Instrument Name</b>	Ultra high-performance liquid chromatography
<b>Generic Instrument Description</b>	Ultra high-performance liquid chromatography: Column chromatography where the mobile phase is a liquid, the stationary phase consists of very small (< 2 microm) particles and the inlet pressure is relatively high.

[ [table of contents](#) | [back to top](#) ]

## Deployments

### AE1913

<b>Website</b>	<a href="https://www.bco-dmo.org/deployment/916412">https://www.bco-dmo.org/deployment/916412</a>
<b>Platform</b>	R/V Atlantic Explorer
<b>Start Date</b>	2019-06-16
<b>End Date</b>	2019-06-28
<b>Description</b>	coordinated deployments: McLane pumps, AUV Clio, CTD, trace metal rosette

[ [table of contents](#) | [back to top](#) ]

## Project Information

## **Collaborative Research: Direct Characterization of Adaptive Nutrient Stress Responses in the Sargasso Sea using Protein Biomarkers and a Biogeochemical AUV (Nutrient Stress Responses and AUV Clio)**

**Coverage:** Bermuda Atlantic Time Series

### *NSF Award Abstract:*

Microscopic communities in the ocean can be surprisingly diverse. This diversity makes it difficult to study the individual organisms and reactions that control specific reactions controlling nutrient cycles. Past studies confirm that iron and nitrogen are vital elements for biological growth. There is increasing evidence, however, that other chemicals such as silica, zinc, cobalt, and vitamin B12 may be just as important. This project will provide an unprecedented view of community distributions using new molecular methods to isolate and link active proteins to specific chemical cycles during the very first research deployment of a brand-new autonomous underwater vehicle (AUV). The AUV will collect samples in programmed patterns by pumping water directly into its filtering mechanism and then return the samples to the ship for analysis. The Bermuda Atlantic Time-series Study (BATS) station, which provides abundant supporting data, is the site for this innovative investigation into the microbial ecology and chemistry of the open oceans. Additionally, data will be widely distributed to other scientists through the Ocean Protein Portal website being developed by the Woods Hole Oceanographic Institute (WHOI) and the Biological and Chemical Oceanography Data Management Office. Data will also contribute a new teaching module in the Marine Bioinorganic Chemistry course at WHOI.

This first scientific deployment of the newly engineered and constructed biogeochemical AUV, Clio, will generate a novel dataset to examine marine microbial biogeochemical cycles in the Northwestern Atlantic oligotrophic ocean in unprecedented detail and at high vertical resolution. First the project proposes to understand if the microbial community reflects the varying chemical composition and cyanobacterial species through nutrient response adaptations. Additionally, the research will determine if iron stress in the low light *Prochlorococcus* ecotype found in the deep chlorophyll maximum is a persistent feature influenced by seasonal dust fluxes. The highly resolved vertical data from the in situ pumping capabilities of Clio are fundamental to a rigorous examination of these biogeochemical questions. This highly transformative dataset will greatly advance understanding of the nutrient and trace element cycling of this region and will be the first field validation of the potentially revolutionary capability these new approaches represent for the study of marine microbial biogeochemistry.

[ [table of contents](#) | [back to top](#) ]

---

## **Funding**

<b>Funding Source</b>	<b>Award</b>
<a href="#">NSF Division of Ocean Sciences (NSF OCE)</a>	<a href="#">OCE-1658030</a>

[ [table of contents](#) | [back to top](#) ]