

Decorator worm *Diopatra cuprea* single nucleotide polymorphism (SNP) genotypes from collections in the Gulf of Mexico and eastern United States estuaries from 2009 to 2022

Website: <https://www.bco-dmo.org/dataset/998297>

Data Type: Other Field Results

Version: 1

Version Date: 2026-05-08

Project

» [The genetic legacy of an Asian oyster introduction and its disease-causing parasite](#) (Oyster historical genetics)

Contributors	Affiliation	Role
Sotka, Erik	College of Charleston (CofC)	Principal Investigator
York, Amber D.	Woods Hole Oceanographic Institution (WHOI BCO-DMO)	BCO-DMO Data Manager

Abstract

The decorator worm *Diopatra cuprea* Bosc, 1802 (Annelid; Polychaete; Onuphidae) is an ecosystem engineer within high-salinity estuaries of the southern and eastern United States. A previous study revealed five morphologically cryptic mitochondrial lineages across its broad geographic distribution. We genotyped single-nucleotide polymorphisms (SNPs) with RADseq. This dataset includes metadata, methods, links to published code and processed data used to generate figures for the results publication Ziegler et al. (2025,doi:10.1007/s00227-025-04613-8) titled "Multiple cryptic lineages and restricted gene flow in the decorator worm *Diopatra cuprea*." This dataset also includes genetic accession identifiers for sequence data contributed to the National Center for Biotechnology Information (NCBI)'s Sequence Read Archive (SRA), available under BioProject PRJNA1103840.

Table of Contents

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

Coverage

Location: Eastern United States and Gulf of Mexico estuaries, soft sediment intertidal zone

Spatial Extent: N:42.04647 E:-69.99714 S:27.4575 W:-84.51061

Temporal Extent: 2009-01-01 - 2022-12-31

Dataset Description

See "Related Datasets" for sanger-sequenced a mitochondrial locus (COI) metadata and accession information at NCBI's Genbank database.

Methods & Sampling

Note: This metadata section includes description of a related dataset (sanger-sequenced a mitochondrial locus (COI) that was generated as part of this study. See "Related Datasets" section for metadata and NCBI Genbank accession identifiers in that dataset.

Diopatra cuprea were sampled from 19 locations on the United States east coast, from Duxbury MA (the genus' northern limit) to St. Teresa Beach FL. Tubes were excavated with a shovel, the worm removed, and antennae clipped off and preserved in 95% ethanol for later DNA extraction. or 2022 samples, twenty-five randomly selected individuals from each population were extracted for DNA, with the exception of Broad River Estuary populations from which 20 and 22 samples were collected. Approximately 25 mg of tissue wet weight was rinsed of ethanol with deionized water and extracted with the Qiagen DNeasy Blood and Tissue Kit (Qiagen, Valencia, CA, USA) following the manufacturer's extraction protocol. Extractions were screened with 1.5% agarose gels and DNA was quantified using a NanoDrop 200 spectrophotometer (Thermo Scientific, Waltham, MA, USA) to obtain concentration values and purity.

A portion of the COI gene was then PCR amplified using the protocol in Berke, et al. (2010, doi:10.1111/j.1466-8238.2009.00509.x). These were cleaned with an EXO-SAP-IT protocol and sent for Sanger sequencing with these same primers at a private company.

Double digest restriction-associated DNA sequencing (or ddRADSeq) library was prepared on 312 samples following the protocol in Parchman et al. (2012,doi: 10.1111/j.1365-294X.2012.05513.x). Briefly, we digested gDNA with two restriction enzymes, EcoRI and MseI, and ligated adaptors containing unique 8 to 10 bp barcodes to the digested DNA of each individual. The products were then PCR amplified in two independent reactions with standard Illumina primers. All amplicons were pooled and shipped to the University of Texas Genomic Sequencing and Analysis Facility or the Tufts University Core Facility, which used Pippin Prep[®] to isolate the 300–450 bp fraction. This fraction was then single-read sequenced (100 basepairs) with Illumina HiSeq 4000 machine. FASTQ sequences were demultiplexed using custom Unix code. 510 million reads contained barcode sequence (range = 39 to 6.8 million (M) reads per sample; mean = 1.6 M reads), and 234 individuals had at least 200 K reads that were analyzed.

FASTQ files were uploaded to the NCBI's Sequence Read Archive (SRA) that were generated as described above. See the "Related Datasets" section for SRA metadata and accessions.

Organism identifier (Life Science Identifier (LSID)):

Diopatra cuprea, Bosc, 1802, LSID(urn:lsid:marinespecies.org:taxname:157339)

BCO-DMO Processing Description

- Loaded TSV file SRAattributes+meta.txt as table 998297_v1_diopatra-sra-metadata, with empty string, "nd", and "NA" treated as missing values
- Renamed column: accession to BioSample
- Applied metadata (descriptions, standard name IDs, units) to 14 columns in 998297_v1_diopatra-sra-metadata including sample_ID, geographic.location, sites.Lat, sites.Lon, mtHaplotype, and related fields
- Renamed 13 columns in 998297_v1_diopatra-sra-metadata: sample_ID to sample_id, Pop. Defined to pop_defined, geographic.location to geographic_location, State.Name to state_name, EstuaryLocal to estuary_local, mtHaplotype to mt_haplotype, mtHap.Group to mt_hap_group, mtHap.Subgroup to mt_hap_subgroup, Year.Collecte d to year_collected, Collector to collector, sites.Lat to site_lat, sites.Lon to site_lon, sites.Collector to sites_collector
- Negated site_lon values (multiplied by -1) using a computed field mathematical operation (longitudes W are negative)
- Rounded site_lat and site_lon to 5 decimal places with trailing zeros preserved
- Loaded CSV file SraRunTable_PRJNA1103840.csv (output from NCBI Run Selector Tool) using filename as table name, with empty string and "nd" treated as missing values
- Joined SRA metadata from the NCBI Run selector into provided metadata table using BioSample as the join key. Tables both were checked first that they were unique by BioSample. Added SRA metadata not present in the attributes file including Assay Type, AvgSpotLen, Bases, BioProject, BioSampleModel, Bytes, Experiment, Instrument, Library Name, LibraryLayout, LibrarySelection, LibrarySource, Organism, Platform, SRA Run, SRA Study, isolate, isolation_source, tissue, and version
- Reordered 34 columns, placing year_collected and sample_id first, followed by sample metadata, SRA identifiers, and sequencing metadata
- Renamed columns: Assay Type to Assay_Type, Library Name to Library_Name
- Applied updated metadata (descriptions, standard name IDs, units) to all 34 columns

- Set final data types for columns
- Output written to 998297_v1_diopatra-sra-metadata.csv

Problem Description

NA

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

Related Publications

Berke, S. K., Mahon, A. R., Lima, F. P., Halanych, K. M., Wethey, D. S., & Woodin, S. A. (2010). Range shifts and species diversity in marine ecosystem engineers: patterns and predictions for European sedimentary habitats. *Global Ecology and Biogeography*, 19(2), 223–232. Portico. <https://doi.org/10.1111/j.1466-8238.2009.00509.x>
Methods

Parchman, T. L., Gompert, Z., Mudge, J., Schilkey, F. D., Benkman, C. W., & Buerkle, C. A. (2012). Genome-wide association genetics of an adaptive trait in lodgepole pine. *Molecular Ecology*, 21(12), 2991–3005. Portico. <https://doi.org/10.1111/j.1365-294x.2012.05513.x> <https://doi.org/10.1111/j.1365-294X.2012.05513.x>
Methods

Ziegler, A. J., Bell, T. M., Berke, S. K., Strand, A. E., & Sotka, E. E. (2025). Multiple cryptic lineages and restricted gene flow in the decorator worm *Diopatra Cuprea*. *Marine Biology*, 172(3). <https://doi.org/10.1007/s00227-025-04613-8>
Results

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

Related Datasets

IsRelatedTo

College of Charleston (2024). Multiple cryptic lineages and restricted gene flow in the decorator worm *Diopatra cuprea*. 2024/04. NCBI:BioProject: PRJNA1103840. In: BioProject [Internet]. Bethesda, MD: National Library of Medicine (US), National Center for Biotechnology Information. Available from: <http://www.ncbi.nlm.nih.gov/bioproject/PRJNA1103840>.

Sotka, E. (2026) **Decorator worm *Diopatra cuprea* mitochondrial locus (COI) sequences from collections in the Gulf of Mexico and eastern United States estuaries from 2009 to 2022.** Biological and Chemical Oceanography Data Management Office (BCO-DMO). (Version 1) Version Date 2026-05-19 <http://lod.bco-dmo.org/id/dataset/999053> [[view at BCO-DMO](#)]
Relationship Description: Data generated from the same study.

Software

Erik Sotka. (2026). *esotka/DiopatraSNPs: v1.0* (Version v1.0) [Computer software]. Zenodo. <https://doi.org/10.5281/ZENODO.18672742> <https://doi.org/10.5281/zenodo.18672742>

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

Parameters

Parameter	Description	Units
year_collected	Year sample was collected (yyyy)	unitless

sample_id	individual sample ID	unitless
pop_defined	Popular name of site	unitless
geographic_location	Sample Site	unitless
state_name	State name (USA)	unitless
estuary_local	Coastal vs Inland	unitless
mt_haplotype	mtDNA haplotype from Sotka et al. 2023 doi.org/10.3390/biology12040521	unitless
mt_hap_group	mtDNA haplotype from Sotka et al. 2023 doi.org/10.3390/biology12040521	unitless
mt_hap_subgroup	mtDNA haplotype from Sotka et al. 2023 doi.org/10.3390/biology12040521	unitless
collector	Person	unitless
SRA_Run	National Center for Biotechnology Information (NCBI) Run accession in the Sequence Read Archive (SRA).	unitless
SRA_Study	National Center for Biotechnology Information (NCBI) study accession in the Sequence Read Archive (SRA).	unitless
version	SRA accession version number	unitless
BioSample	National Center for Biotechnology Information (NCBI) BioSample accession	unitless
site_lat	Site latitude	decimal degrees
site_lon	Site longitude	decimal degrees
sites_collector	co-author that collected samples	unitless
Assay_Type	Assay Type	unitless
AvgSpotLen	Average Spot Length (SRA metadata from NCBI)	count
Bases	Bases	count

BioProject	National Center for Biotechnology Information (NCBI) BioProject identifier.	unitless
BioSampleModel	BioSample model (e.g. Invertebrate)	unitless
Bytes	Size in bytes	bytes
Experiment	National Center for Biotechnology Information (NCBI) Experiment accession in the Sequence Read Archive (SRA).	unitless
Instrument	Instrument name	unitless
isolate	Isolate	unitless
isolation_source	isolation source	unitless
Library_Name	Library Name	unitless
LibraryLayout	Library Layout	unitless
LibrarySelection	Library Selection	unitless
LibrarySource	Library Source	unitless
Organism	organism sampled (Diopatra cuprea)	unitless
Platform	platform	unitless
tissue	tissue sample location	unitless

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

Instruments

Dataset-specific Instrument Name	Illumina sequencing machine
Generic Instrument Name	Automated DNA Sequencer
Dataset-specific Description	Illumina HiSeq 4000 or Illumina NovaSeq 6000
Generic Instrument Description	A DNA sequencer is an instrument that determines the order of deoxynucleotides in deoxyribonucleic acid sequences.

Project Information

The genetic legacy of an Asian oyster introduction and its disease-causing parasite (Oyster historical genetics)

Coverage: Global

NSF abstract:

During the 20th century, the Pacific oyster *Crassostrea gigas* was deliberately introduced from its native range of coastal Asia to the estuaries of six continents. While the introduced Pacific oysters are widely aquacultured and thus can generate local economic wealth, they sometimes outcompete native oysters, and can carry microbial, animal and plant hitchhikers that negatively impact local economies and the ecological functioning of local estuaries. This study comprehensively assesses the pathways and sources of Pacific oyster introductions using a worldwide, population genetic survey. Simultaneously, the study also assesses the pathways and source of one hitchhiking protist (*Haplosporidium nelsoni*) that causes the disease MSX (multinucleated sphere X) in the Virginia oyster (*Crassostrea virginica*) along the eastern seaboard of the United States. One goal of this research is to generate management strategies that combat the negative impacts of the Pacific oyster and its associated invaders, and minimize future invasions. A second goal is to minimize some uncertainty about the population biology of the devastating *Haplosporidium* parasite, and thus, increase confidence of policy makers who are managing shellfish health, restoration and commerce. By quantifying the pathways and sources of *C. gigas*, this project may inform strategies to combat negative impacts of *C. gigas* and its associated invaders, as well as minimize future invasions. Moreover, quantifying dispersal within and among populations of *H. nelsoni* along the US East Coast will provide perspective on the effectiveness of regional biosecurity measures in preventing the ongoing dispersal of this destructive pathogen via aquaculture. In addition, the project lends itself well to programs that foster critical thinking and research experience among both undergraduate and K-12 students. The project provides opportunities for 6-9 undergraduates to perform research, includes a 2-day workshop on bioinformatics for the wider undergraduate community, and facilitates ongoing opportunities for K-12 students to participate in citizen-science research.

There is a wealth of information on the source, pathways and vectors of *C. gigas* based largely on historical documents but no study has comprehensively tested whether these historical accounts are correct using a worldwide, population genetic survey. Using >14K single-nucleotide polymorphisms (SNPs) from 41 populations across five continents a high level of spatial genetic differentiation was found within the native range and differences in source populations among non-native regions. Preliminary genetic data indicated that the parasitic protist, *Haplosporidium nelsoni* arrived with *C. gigas* imports to the US Atlantic coastline and then infected the native *C. virginica*, however the native source populations, the pathways and vector from which *H. nelsoni* arrived remain unknown. This project couples high-throughput sequencing technologies and Approximate Bayesian Computing (ABC)-based models to answer the following: What are the population genomic patterns among *C. gigas* from native and non-native regions? What are the population genomic patterns of *Haplosporidium nelsoni* among Asian and North American *Crassostrea gigas* and eastern North American *C. virginica*? What were the source populations and invasion pathways of *C. gigas* and *H. nelsoni*? Identifying source locations, pathways and vectors of introduction of *C. gigas* will provide researchers with a null-model of invasion history for dozens of other non-native species that were transported with *C. gigas*. Currently, there are no verified 'vector maps' for historical shipments of *C. gigas* that are similar to those generated from modern-day or historical shipping records.

This award reflects NSF's statutory mission and has been deemed worthy of support through evaluation using the Foundation's intellectual merit and broader impacts review criteria.

Funding

Funding Source	Award
NSF Division of Ocean Sciences (NSF OCE)	OCE-1924599

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