

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/359135545>

# A workflow for sequencing, annotating and curating aquatic symbiomes

Poster · March 2022

DOI: 10.5281/zenodo.5385621)

CITATIONS

0

READS

23

3 authors, including:



[Camilla Eldridge](#)  
Universidad Diego Portales

2 PUBLICATIONS 0 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Phylogeography of European and Asian populations of *Radix auricularia* [View project](#)



Evolutionary origins and the divergence of *Schistosoma turkestanicum* across Eurasia [View project](#)

# A workflow for sequencing, annotating and curating aquatic symbiomes

Camilla Eldridge<sup>1</sup>, Michael Paulini<sup>1</sup>, Caroline Howard<sup>1</sup>, Graeme Oathley<sup>1</sup>, Radka Platte<sup>1</sup>, Catherine McCarthy<sup>1</sup>, Nancy Holroyd<sup>1</sup>, Nicola Chapman<sup>1</sup>, Neha Wali<sup>1</sup>, Haoyu Niu<sup>1</sup>, Peter Harrison<sup>1</sup>, Alexey Sokolov<sup>1</sup>, Talal Ibrahim<sup>1</sup>, Claudia C. Weber<sup>1</sup>, Emmelien Vancaester<sup>1</sup>, Marcela Uliano-Silva<sup>1</sup>, Ksenia Krashenninnikova<sup>1</sup>, Eerik Aunin<sup>1</sup>, Shane McCarthy<sup>1</sup>, Will Chow<sup>1</sup>, Ying Simms<sup>1</sup>, James Torrance<sup>1</sup>, Alan Tracey<sup>1</sup>, Sarah Pelan<sup>1</sup>, Joana Collins<sup>1</sup>, Damon-Lee Pointon<sup>1</sup>, Jonathan Wood<sup>1</sup>, Kerstin Howe<sup>1</sup>, Michael Sweet<sup>2</sup>, Anne Thompson<sup>3</sup>, Felipe Porto<sup>4</sup>, Senjie Lin<sup>4</sup>, Patrick Keeling<sup>5</sup>, Helena Villela<sup>6</sup>, Raquel Peixoto<sup>6</sup>, Victoria McKenna<sup>1</sup>, Mara Lawnczak<sup>1</sup>, Mark Blaxter<sup>1</sup>, and members of the ASG consortium (see <https://doi.org/10.5281/zenodo.5385621>)

## The Aquatic Symbiosis Genomics Project

The Aquatic Symbiosis Genomics Project, a collaboration between fourteen international research hubs, is sequencing more than 1000 species from 500 aquatic symbiotic systems to provide insights into complex ecological and evolutionary relationships. The goal is to generate annotated, gold-standard reference genomes for host and symbiont. An initial workflow, from sample preparation to *de novo* assembly, curation and annotation, has been developed to overcome the challenges associated with large-scale sequencing of diverse aquatic symbiotic systems.

### Sample acquisition and management

Samples are collected by partners, snap frozen, and a detailed metadata manifest completed. We sample ethically and legally, including meeting all applicable regulatory compliance steps. These steps are completed before samples are shipped by specialist cold chain couriers in dry ice.

### DNA extraction

Weighed subsamples are homogenised in lysis buffer/trizol with BioMasher & PowerMasher. A good quality extraction features High Molecular Weight DNA as the main component. The level of RNA degradation is recorded.

### Sequencing

Long read data are generated from a low input library protocol followed by sequencing on the PacBio Sequel IIe. For HiC, DNA-protein complexes are cross-linked prior to sample fragmentation and DNA extraction using the Arima-HiC protocol. Libraries are sequenced on Illumina NovaSeq (PE 150 reads).

### Genome assembly and QC

Genomes are assembled from PacBio HiFi reads using HiCanu or Hifiasm and then purged to separate haplotypes. MitoHifi assembles the mitochondrion and eventual plastids. Salsascaffolds genomes using HiC reads (Arima or Qiagen).

## ASG Workflow

### Data publication and accessibility

Annotated genomes will be hosted by Ensembl where final annotations including protein domain information will be accessible and a Genome Note describing the assembly will be published in Wellcome Open Research.

### Genome Annotation

Curated genomes are annotated using the ASG annotation pipeline that generates a repeat library for each species and uses RNA-seq to train gene model predictors in Braker2. Initial annotations are viewable on the UCSC browser.

### Genome curation and postprocessing

Assembled genomes are curated using HiC maps and gEVAL where HiC contact information, coverage and assembly gaps are used in manual curation to scaffold long read assemblies to chromosomal level. Curated genomes are post-processed and submitted to the ENA.

### Host-symbiont separation

The assembled scaffolds are separated based on taxa using a combination of the NCBI BLAST based decontamination pipeline, the CobiontID workflow and BlobToolKit.

Figure 5: BlobToolKit plot (GC% x axis, coverage y axis) showing an assembly of a fish genome (upper right) also contains a myxozoan parasite (lower left)

## Conclusions and future work

- Aquatic symbiome samples have high levels of environmental contamination.
- Ongoing R & D needed for DNA extraction.
- Ongoing development of incorporating CobiontID into the workflow to enable separation of symbionts at read level.



Figure 1: STS and COPO sample management and tracking.



Figure 2: DNA extraction workflow.

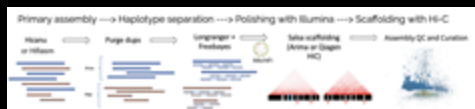


Figure 3: The assembly pipeline.

Fig 4: Read plot of Tetranucleotides coloured by coding density separating *Membranipora membranacea*, a sea rat, from its cobionts.

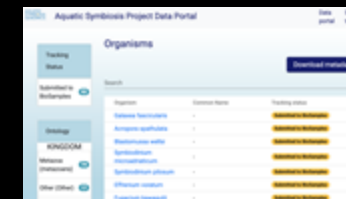
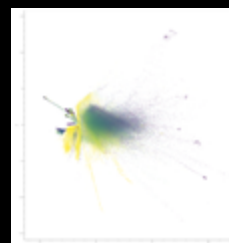


Figure 8: A dedicated data portal has been setup providing BioSample display, status tracking and taxonomy navigation.



Figure 7: Viewing annotations

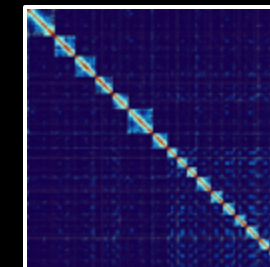


Figure 6: HiC contact map of a colonial sea squirt, *Aplidium turbinatum*