

Camilla Eldridge^{1*}, Michael Paulini^{1*}, Caroline Howard¹, Graeme Oathley¹, Radka Platte¹, Catherine McCarthy¹, Nancy Holroyd¹, Nicola Chapman¹, Neha Wali¹, Haoyu Niu¹, Peter Harrison¹, Alexey Sokolov¹, Talal Ibrahim¹, Claudia C. Weber¹, Emmelien Vancaester¹, Marcela Uliano-Silva¹, Ksenia Krasheninnikova¹, Eerik Aunin¹, Shane McCarthy¹, Will Chow, Ying Simms¹, James Torrance¹, Alan Tracey¹, Sarah Pelan¹, Joana Collins¹, Damon-Lee Pointon¹, Jonathan Wood¹, Kerstin Howe¹, Michael Sweet², Anne Thompson³, Felipe Porto⁴, Senjie Lin⁴, Patrick Keeling⁵, Helena Villega⁶, Raquel Peixoto⁶, Victoria McKenna¹, Mara Lawniczak¹, Mark Blaxter¹, 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 (PE150 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. Salsa scaffolds genomes using HiC reads (Arima or Qiagen).



Figure 1. STS and COPO sample management and tracking.



Figure 2. DNA extraction workflow.

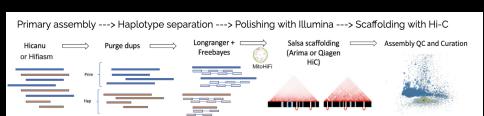
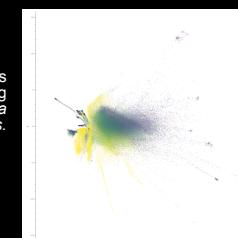


Figure 3. The assembly pipeline.

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

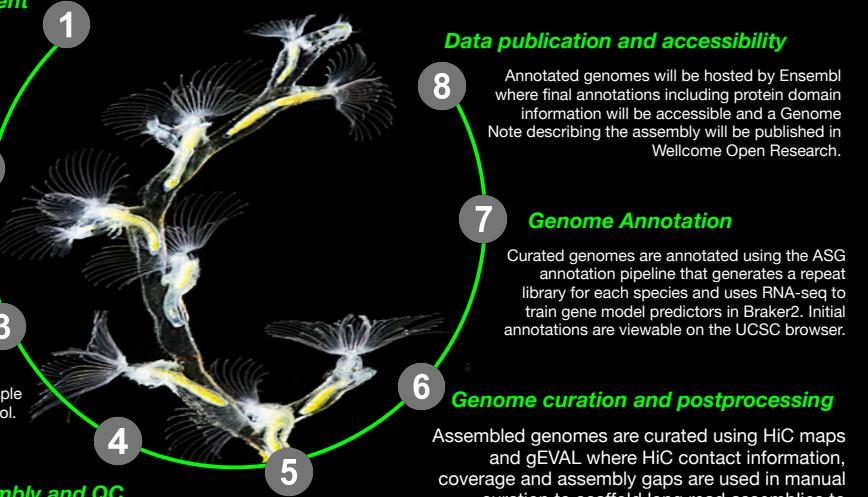


Project status

- 92 species: arrived at Sanger
- 64 species: completed DNA extractions
- 39 species: in sequencing

Sequencing and primary assembly QC data are available on TOLQC <https://tolqc.cog.sanger.ac.uk/>

ASG Workflow



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.

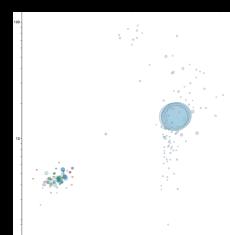


Figure 5: BlobToolKit plot (GC% x axis, coverage y axis) showing an assembly of a fish genome (upper right) also contains a myxozoan parasite (lower left)

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

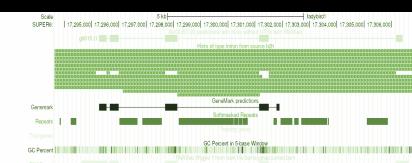


Figure 7: Viewing annotations

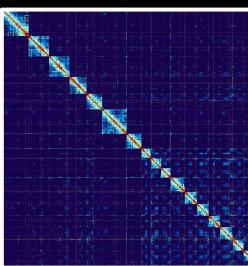


Figure 6. HiC contact map of a colonial sea squirt, *Aplidium turbinatum*.



GORDON AND BETTY MOORE FOUNDATION