



DATA NOTE

The genome sequence of the Birch Mocha moth, *Cyclophora albipunctata* (Hufnagel, 1767)

[version 1; peer review: awaiting peer review]

Tom Prescott¹, Marc Botham²,
Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory
team,

Wellcome Sanger Institute Scientific Operations: Sequencing Operations,
Wellcome Sanger Institute Tree of Life Core Informatics team,
Tree of Life Core Informatics collective, Darwin Tree of Life Consortium

¹Butterfly Conservation Scotland, Stirling, Scotland, UK

²UK Centre for Ecology & Hydrology, Wallingford, England, UK

V1 First published: 04 Nov 2024, 9:641
<https://doi.org/10.12688/wellcomeopenres.23257.1>

Latest published: 04 Nov 2024, 9:641
<https://doi.org/10.12688/wellcomeopenres.23257.1>

Open Peer Review

Approval Status AWAITING PEER REVIEW

Any reports and responses or comments on the article can be found at the end of the article.

Abstract

We present a genome assembly from a female specimen of *Cyclophora albipunctata* (the Birch Mocha; Arthropoda; Insecta; Lepidoptera; Geometridae). The genome sequence has a total length of 319.40 megabases. Most of the assembly is scaffolded into 32 chromosomal pseudomolecules, including the W and Z sex chromosomes. The mitochondrial genome has also been assembled and is 16.92 kilobases in length. Gene annotation of this assembly on Ensembl identified 16,542 protein-coding genes.

Keywords

Cyclophora albipunctata, Birch Mocha moth, genome sequence, chromosomal, Lepidoptera



This article is included in the [Tree of Life gateway](#).

Corresponding author: Darwin Tree of Life Consortium (mark.blaxter@sanger.ac.uk)

Author roles: **Prescott T:** Investigation, Resources; **Botham M:** Investigation, Resources, Writing – Review & Editing;

Competing interests: No competing interests were disclosed.

Grant information: This work was supported by Wellcome through core funding to the Wellcome Sanger Institute [206194, <https://doi.org/10.35802/206194>] and the Darwin Tree of Life Discretionary Award [218328, <https://doi.org/10.35802/218328>]. *The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

Copyright: © 2024 Prescott T *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

How to cite this article: Prescott T, Botham M, Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory team *et al.*

The genome sequence of the Birch Mocha moth, *Cyclophora albipunctata* (Hufnagel, 1767) [version 1; peer review: awaiting peer review] Wellcome Open Research 2024, 9:641 <https://doi.org/10.12688/wellcomeopenres.23257.1>

First published: 04 Nov 2024, 9:641 <https://doi.org/10.12688/wellcomeopenres.23257.1>

Species taxonomy

Eukaryota; Opisthokonta; Metazoa; Eumetazoa; Bilateria; Protostomia; Ecdysozoa; Panarthropoda; Arthropoda; Mandibulata; Pancrustacea; Hexapoda; Insecta; Dicondylia; Pterygota; Neoptera; Endopterygota; Amphiesmenoptera; Lepidoptera; Glossata; Neolepidoptera; Heteroneura; Ditrysia; Obtectomera; Geometroidea; Geometridae; Sterrhinae; *Cyclophora*; *Cyclophora albipunctata* (Hufnagel, 1767) (NCBI:txid505405).

Background

The Birch Mocha (*Cyclophora albipunctata*) (Figure 1) is a geometrid moth in the subfamily Sterrhinae (Mochas and Waves) found in the Palaearctic region. In the UK it is mostly associated with woodland and heathland habitats where the larvae feed on silver (*Betula pendula*) and downy (*B. pubescens*) birches (Waring *et al.*, 2017). Widespread throughout the UK it is most densely distributed in southern and south-east England and East Anglia. It is widespread in the northern half of Scotland where it has a single generation flying between May and July, whereas further south it tends to have two generations flying from early May to late June and then again from July to August. It is thinly scattered elsewhere in the UK, but with historic declines and an overall decline in abundance across monitored sites was observed from 1970–2016 (Randle *et al.*, 2019). Over the same time period there has been a significant increase in distribution and there has been an advance in the timing of the second generation in southern populations (Randle *et al.*, 2019). Winter is spent in the pupal stage.

The genome of the birch mocha, *Cyclophora albipunctata*, was sequenced as part of the Darwin Tree of Life Project, a collaborative effort to sequence all named eukaryotic species in the Atlantic Archipelago of Britain and Ireland. Here we present a chromosomally complete genome sequence for *Cyclophora albipunctata*, based on a female specimen from Glen Strathfarrar, Scotland, UK.



Figure 1. Photograph of *Cyclophora albipunctata* by Patrick Clement (not the specimen used for genome sequencing).

Genome sequence report

The genome of *Cyclophora albipunctata* was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating a total of 23.40 Gb (gigabases) from 2.11 million reads, providing an estimated 75-fold coverage. Primary assembly contigs were scaffolded with chromosome conformation Hi-C data, which produced 158.74 Gb from 1,051.23 million reads, yielding an approximate coverage of 497-fold. Specimen and sequencing details are summarised in Table 1.

Assembly errors were corrected by manual curation, including 8 missing joins or mis-joins. This reduced the scaffold number by 2.38%. The final assembly has a total length of 319.40 Mb in 40 sequence scaffolds, with 27 gaps, and a scaffold N50 of 11.0 Mb (Table 2).

The snail plot in Figure 2 provides a summary of the assembly statistics, indicating the distribution of scaffold lengths and other assembly metrics. Figure 3 shows the distribution of scaffolds by GC proportion and coverage. Figure 4 presents a cumulative assembly plot, with separate curves representing different scaffold subsets assigned to various phyla, illustrating the completeness of the assembly.

Most of the assembly sequence (99.95%) was assigned to 32 chromosomal-level scaffolds, representing 30 autosomes and the W and Z sex chromosomes. These chromosome-level scaffolds, confirmed by the Hi-C data, are named in order of size (Figure 5; Table 3). During manual curation it was noted that Chromosomes Z and W were assigned based on read coverage statistics and synteny to the assembly of *Cyclophora punctaria* (GCA_951394245.1) (Broad *et al.*, 2024).

While not fully phased, the assembly deposited is of one haplotype. Contigs corresponding to the second haplotype have also been deposited. The mitochondrial genome was also assembled and can be found as a contig within the multifasta file of the genome submission, and as a separate fasta file with accession OY720143.1.

The final assembly has a Quality Value (QV) of 68.0 and *k*-mer completeness of 100.0%. BUSCO (v5.4.3) analysis using the lepidoptera_odb10 reference set ($n = 5,286$) indicated a completeness score of 98.0% (single = 97.6%, duplicated = 0.4%).

Metadata for specimens, BOLD barcode results, spectra estimates, sequencing runs, contaminants and pre-curation assembly statistics are given at <https://links.tol.sanger.ac.uk/species/505405>.

Genome annotation report

The *Cyclophora albipunctata* genome assembly (GCA_963082685.1) was annotated at the European Bioinformatics Institute (EBI) on Ensembl Rapid Release. The resulting annotation includes 16,722 transcribed mRNAs from 16,542 protein-coding genes (Table 2; https://rapid.ensembl.org/Cyclophora_albipunctata_GCA_963082685.1/Info/Index). The

Table 1. Specimen and sequencing data for *Cyclophora albipunctata*.

| Project information | | | |
|-----------------------------|--|---------------------|----------------------|
| Study title | <i>Cyclophora albipunctata</i> (birch mocha) | | |
| Umbrella BioProject | PRJEB64086 | | |
| Species | <i>Cyclophora albipunctata</i> | | |
| BioSample | SAMEA112198466 | | |
| NCBI taxonomy ID | 505405 | | |
| Specimen information | | | |
| Technology | ToLID | BioSample accession | Organism part |
| PacBio long read sequencing | ilCycAlbi1 | SAMEA112198496 | Thorax |
| Hi-C sequencing | ilCycAlbi1 | SAMEA112198498 | Other somatic tissue |
| RNA sequencing | ilCycAlbi1 | SAMEA112198497 | Abdomen |
| Sequencing information | | | |
| Platform | Run accession | Read count | Base count (Gb) |
| Hi-C Illumina NovaSeq 6000 | ERR11679405 | 1.05e+09 | 158.74 |
| PacBio Sequel IIe | ERR11673240 | 2.11e+06 | 23.4 |
| RNA Illumina NovaSeq X | ERR12862081 | 6.11e+07 | 9.23 |

average transcript length is 6,348.32, and there are 5.86 exons per transcript.

Methods

Sample acquisition

An adult female specimen of *Cyclophora albipunctata* (specimen ID SAN00002586, ToLID ilCycAlbi1) was collected from Glen Strathfarrar, Highlands, Scotland, United Kingdom (latitude 57.41, longitude -4.73) on 2022-06-27, using a moth trap. The specimen was collected by Tom Prescott (Butterfly Conservation) and identified by Marc Botham (Centre for Ecology & Hydrology) and preserved by flash freezing.

Nucleic acid extraction

The workflow for high molecular weight (HMW) DNA extraction at the Wellcome Sanger Institute (WSI) Tree of Life Core Laboratory includes a sequence of core procedures: sample preparation and homogenisation, DNA extraction, fragmentation and purification. Detailed protocols are available on protocols.io (Denton *et al.*, 2023b). The ilCycAlbi1 sample was prepared for DNA extraction by weighing and dissecting it on dry ice (Jay *et al.*, 2023), and tissue from the thorax was homogenised using a PowerMasher II tissue disruptor (Denton *et al.*, 2023a).

HMW DNA was extracted in the WSI Scientific Operations core using the Automated MagAttract v2 protocol (Oatley *et al.*, 2023). The DNA was sheared into an average fragment size of 12–20 kb in a Megaruptor 3 system (Bates *et al.*, 2023). Sheared DNA was purified by solid-phase reversible

immobilisation, using AMPure PB beads to eliminate shorter fragments and concentrate the DNA (Strickland *et al.*, 2023). The concentration of the sheared and purified DNA was assessed using a Nanodrop spectrophotometer and Qubit Fluorometer using the Qubit dsDNA High Sensitivity Assay kit. Fragment size distribution was evaluated by running the sample on the FemtoPulse system.

RNA was extracted from abdomen tissue of ilCycAlbi1 in the Tree of Life Laboratory at the WSI using the RNA Extraction: Automated MagMax™ mirVana protocol (do Amaral *et al.*, 2023). The RNA concentration was assessed using a Nanodrop spectrophotometer and a Qubit Fluorometer using the Qubit RNA Broad-Range Assay kit. Analysis of the integrity of the RNA was done using the Agilent RNA 6000 Pico Kit and Eukaryotic Total RNA assay.

Hi-C preparation

Tissue from the ilCycAlbi1 sample was processed at the WSI Scientific Operations core, using the Arima-HiC v2 kit. Frozen tissue (stored at -80 °C) was fixed, and the DNA crosslinked using a TC buffer with 22% formaldehyde. After crosslinking, the tissue was homogenised using the Diagnocine Power Masher-II and BioMasher-II tubes and pestles. Following the kit manufacturer's instructions, crosslinked DNA was digested using a restriction enzyme master mix. The 5'-overhangs were then filled in and labelled with biotinylated nucleotides and proximally ligated. An overnight incubation was carried out for enzymes to digest remaining proteins and for crosslinks to reverse. A clean up was performed with SPRIselect beads prior to library preparation.

Table 2. Genome assembly data for *Cyclophora albipunctata*, ilCycAlbi1.1.

| Genome assembly | | |
|--|--|----------------------------|
| Assembly name | ilCycAlbi1.1 | |
| Assembly accession | GCA_963082685.1 | |
| Accession of alternate haplotype | GCA_963082705.1 | |
| Span (Mb) | 319.40 | |
| Number of contigs | 68 | |
| Contig N50 length (Mb) | 8.3 | |
| Number of scaffolds | 40 | |
| Scaffold N50 length (Mb) | 11.0 | |
| Longest scaffold (Mb) | 15.32 | |
| Assembly metrics* | | Benchmark |
| Consensus quality (QV) | 68.0 | ≥ 50 |
| k-mer completeness | 100.0% | ≥ 95% |
| BUSCO** | C:98.0%[S:97.6%,D:0.4%], F:0.6%,M:1.4%,n:5,286 | C ≥ 95% |
| Percentage of assembly mapped to chromosomes | 99.95% | ≥ 95% |
| Sex chromosomes | WZ | localised homologous pairs |
| Organelles | Mitochondrial genome: 16.92 kb | complete single alleles |
| Genome annotation of assembly GCA_963082685.1 at Ensembl | | |
| Number of protein-coding genes | 16,542 | |
| Number of gene transcripts | 16,722 | |

*Assembly metric benchmarks are adapted from column VGP-2020 of “Table 1: Proposed standards and metrics for defining genome assembly quality” from [Rhie et al. \(2021\)](#).

**BUSCO scores based on the lepidoptera_odb10 BUSCO set using version 5.3.2. C = complete [S = single copy, D = duplicated], F = fragmented, M = missing, n = number of orthologues in comparison. A full set of BUSCO scores is available at <https://blobtoolkit.genomehubs.org/view/CAUJAW01/dataset/CAUJAW01/busco>.

Library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core. Pacific Biosciences HiFi circular consensus DNA sequencing libraries were prepared using the PacBio Express Template Preparation Kit v2.0 (Pacific Biosciences, California, USA) as per the manufacturer’s instructions. The kit includes the reagents required for removal of single-strand overhangs, DNA damage repair, end repair/A-tailing, adapter ligation, and nuclease treatment. Library preparation also included a library purification step using AMPure PB beads (Pacific Biosciences, California, USA) and size selection step to remove templates shorter than 3 kb using AMPure PB modified SPRI. DNA concentration was quantified using the Qubit Fluorometer v2.0 and Qubit HS Assay Kit and the final library fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument and gDNA 165kb gDNA and 55kb BAC analysis kit. Samples

were sequenced using the Sequel IIe system (Pacific Biosciences, California, USA). The concentration of the library loaded onto the Sequel IIe was between 40–135 pM. The SMRT link software, a PacBio web-based end-to-end workflow manager, was used to set-up and monitor the run, as well as perform primary and secondary analysis of the data upon completion.

For Hi-C library preparation, DNA was fragmented to a size of 400 to 600 bp using a Covaris E220 sonicator. The DNA was then enriched, barcoded, and amplified using the NEBNext Ultra II DNA Library Prep Kit following manufacturers’ instructions. The Hi-C sequencing was performed using paired-end sequencing with a read length of 150 bp on an Illumina NovaSeq 6000 instrument.

Poly(A) RNA-Seq libraries were constructed using the NEB Ultra II RNA Library Prep kit, following the manufacturer’s

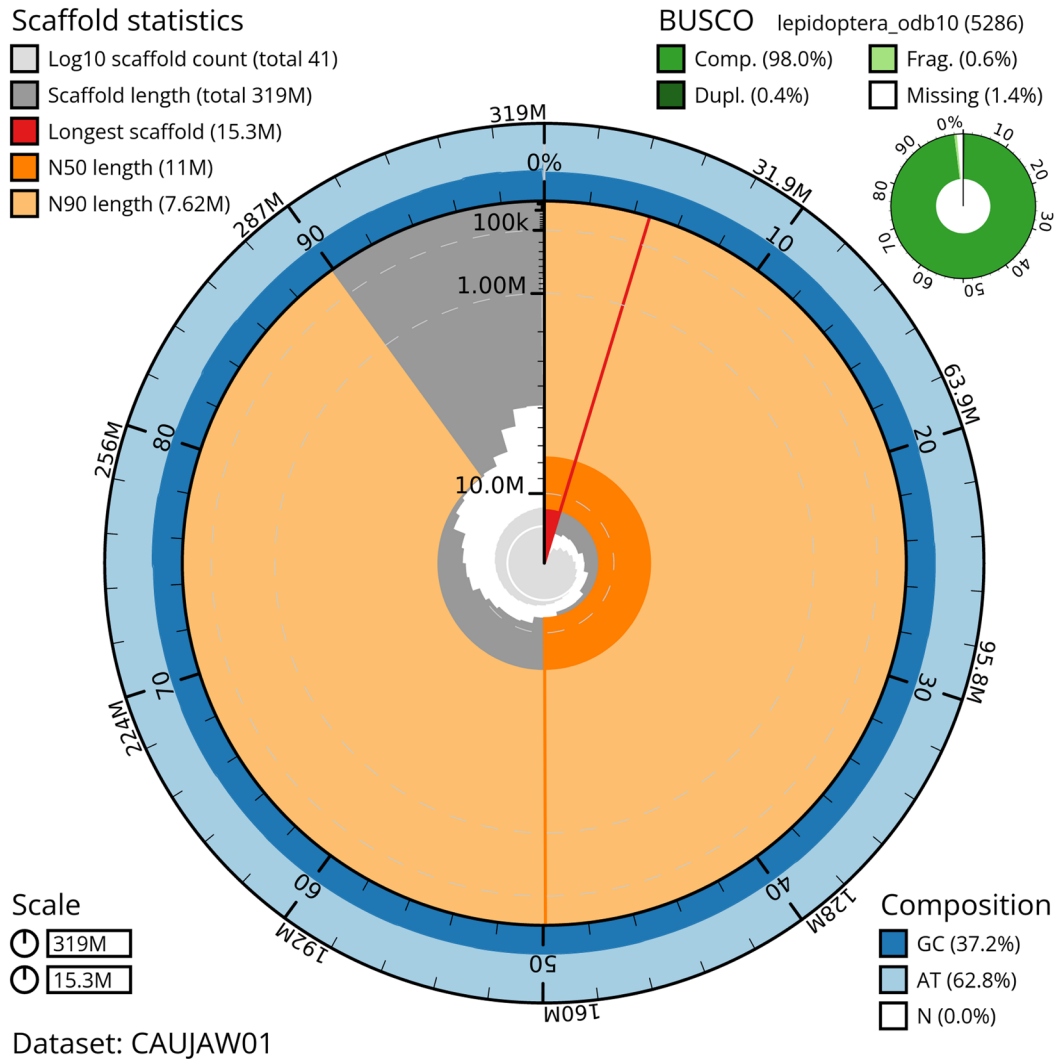


Figure 2. Genome assembly of *Cyclophora albipunctata*, ilCycAlbi1.1: metrics. The BlobToolKit snail plot shows N50 metrics and BUSCO gene completeness. The main plot is divided into 1,000 size-ordered bins around the circumference with each bin representing 0.1% of the 319,416,675 bp assembly. The distribution of scaffold lengths is shown in dark grey with the plot radius scaled to the longest scaffold present in the assembly (15,324,126 bp, shown in red). Orange and pale-orange arcs show the N50 and N90 scaffold lengths (10,985,054 and 7,618,004 bp), respectively. The pale grey spiral shows the cumulative scaffold count on a log scale with white scale lines showing successive orders of magnitude. The blue and pale-blue area around the outside of the plot shows the distribution of GC, AT and N percentages in the same bins as the inner plot. A summary of complete, fragmented, duplicated and missing BUSCO genes in the lepidoptera_odb10 set is shown in the top right. An interactive version of this figure is available at <https://blobtoolkit.genomehubs.org/view/CAUJAW01/dataset/CAUJAW01/snail>.

instructions. RNA sequencing was performed on the Illumina NovaSeq X instrument.

Genome assembly, curation and evaluation

Assembly

The HiFi reads were first assembled using Hifiasm (Cheng *et al.*, 2021) with the --primary option. Haplotypic duplications were identified and removed using purge_dups (Guan *et al.*, 2020). The Hi-C reads were mapped to the primary contigs using bwa-mem2 (Vasimuddin *et al.*, 2019). The contigs were

further scaffolded using the provided Hi-C data (Rao *et al.*, 2014) in YaHS (Zhou *et al.*, 2023) using the --break option. The scaffolded assemblies were evaluated using Gfastats (Formenti *et al.*, 2022), BUSCO (Manni *et al.*, 2021) and MERQURY.FK (Rhie *et al.*, 2020).

The mitochondrial genome was assembled using MitoHiFi (Uliano-Silva *et al.*, 2023), which runs MitoFinder (Allio *et al.*, 2020) and uses these annotations to select the final mitochondrial contig and to ensure the general quality of the sequence.

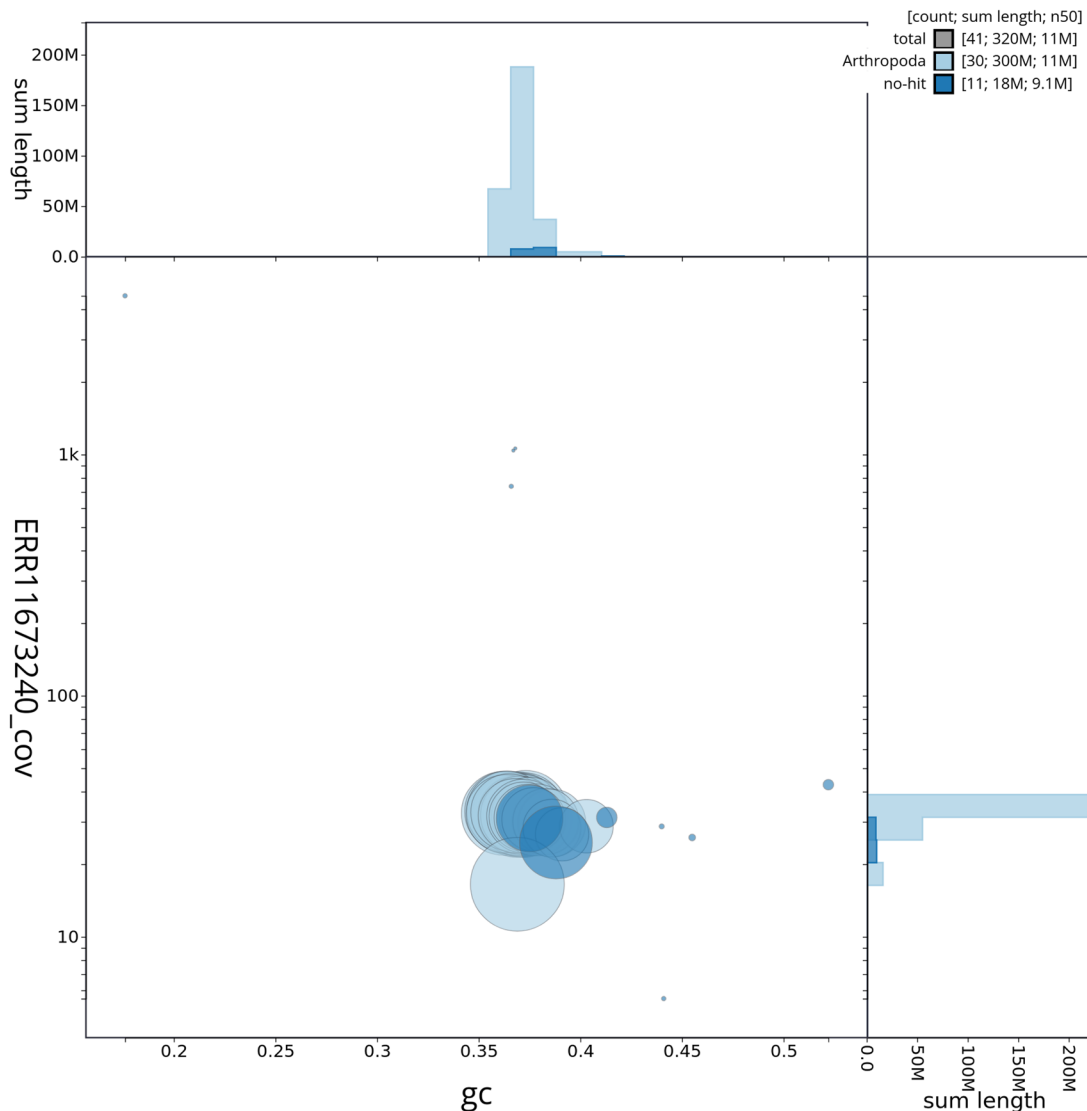


Figure 3. Genome assembly of *Cyclophora albipunctata*, iCycAlbi1.1. BlobToolKit GC-coverage plot showing sequence coverage (vertical axis) and GC content (horizontal axis). The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. The histograms along the axes display the total length of sequences distributed across different levels of coverage and GC content. An interactive version of this figure is available at <https://blobtoolkit.genomehubs.org/view/CAUJAW01/dataset/CAUJAW01/blob>.

Assembly curation

The assembly was decontaminated using the Assembly Screen for Cobionts and Contaminants (ASCC) pipeline (article in preparation). Manual curation was primarily conducted using PretextView (Harry, 2022), with additional insights provided by JBrowse2 (Diesh *et al.*, 2023) and HiGlass (Kerpedjiev *et al.*, 2018). Scaffolds were visually inspected and corrected as described by Howe *et al.* (2021). Any identified contamination, missed joins, and mis-joins were corrected, and duplicate sequences were tagged and removed. The curation process is documented at <https://gitlab.com/wtsi-grit/rapid-curation> (article in preparation).

Evaluation of the final assembly

A Hi-C map for the final assembly was produced using bwa-mem2 (Vasimuddin *et al.*, 2019) in the Cooler file format (Abdennur & Mirny, 2020). To assess the assembly metrics, the *k*-mer completeness and QV consensus quality values were calculated in Merqury (Rhie *et al.*, 2020). This work was done using the “sanger-tol/readmapping” (Surana *et al.*, 2023a) and “sanger-tol/genomenote” (Surana *et al.*, 2023b) pipelines. The genome readmapping pipelines were developed using the nf-core tooling (Ewels *et al.*, 2020), use MultiQC (Ewels *et al.*, 2016), and make extensive use of the Conda package manager, the Bioconda initiative (Grüning *et al.*, 2018), the

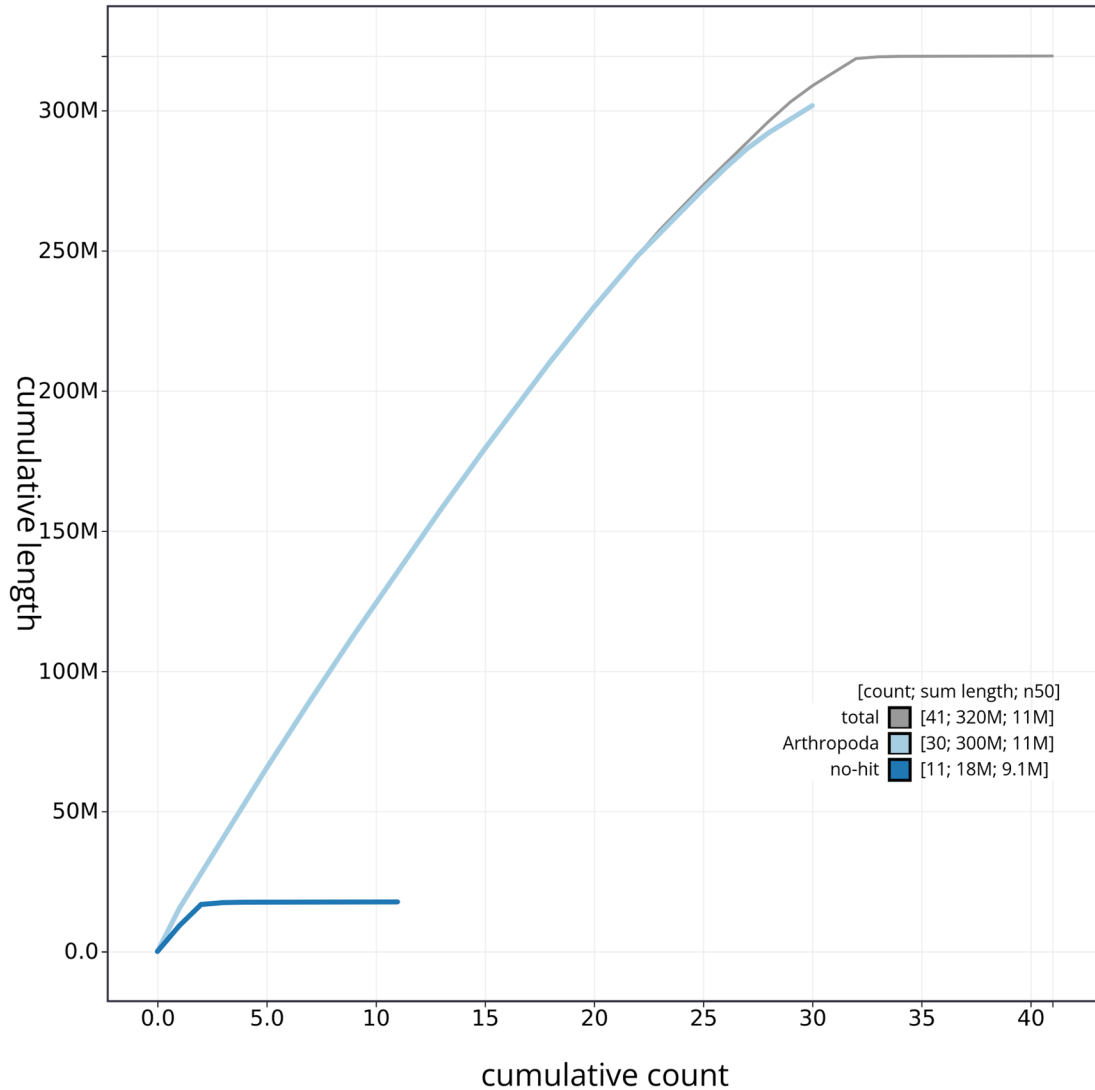


Figure 4. Genome assembly of *Cyclophora albipunctata* iCycAlbi1.1: BlobToolKit cumulative sequence plot. The grey line shows cumulative length for all sequences. Coloured lines show cumulative lengths of sequences assigned to each phylum using the buscogenes taxrule. An interactive version of this figure is available at <https://blobtoolkit.genomehubs.org/view/CAUJAW01/dataset/CAUJAW01/cumulative>.

Biocontainers infrastructure (da Veiga Leprevost *et al.*, 2017), and the Docker (Merkel, 2014) and Singularity (Kurtzer *et al.*, 2017) containerisation solutions. The genome was also analysed within the BlobToolKit environment (Challis *et al.*, 2020) and BUSCO scores (Manni *et al.*, 2021) were calculated.

Table 4 contains a list of relevant software tool versions and sources.

Genome annotation

The BRAKER2 pipeline (Brůna *et al.*, 2021) was used in the default protein mode to generate annotation for the *Cyclophora albipunctata* assembly (GCA_963082685.1) in Ensembl Rapid Release at the EBI.

Wellcome Sanger Institute – Legal and Governance

The materials that have contributed to this genome note have been supplied by a Darwin Tree of Life Partner. The submission of materials by a Darwin Tree of Life Partner is subject to the ‘**Darwin Tree of Life Project Sampling Code of Practice**’, which can be found in full on the Darwin Tree of Life website [here](#). By agreeing with and signing up to the Sampling Code of Practice, the Darwin Tree of Life Partner agrees they will meet the legal and ethical requirements and standards set out within this document in respect of all samples acquired for, and supplied to, the Darwin Tree of Life Project.

Further, the Wellcome Sanger Institute employs a process whereby due diligence is carried out proportionate to the nature

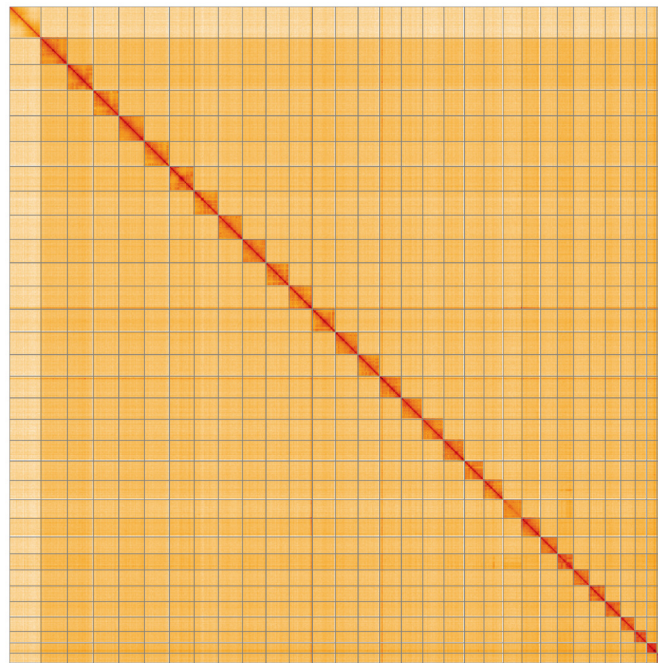


Figure 5. Genome assembly of *Cyclophora albipunctata* ilCycAlbi1.1: Hi-C contact map of the ilCycAlbi1.1 assembly, visualised using HiGlass. Chromosomes are shown in order of size from left to right and top to bottom. An interactive version of this figure may be viewed at <https://genome-note-higlass.tol.sanger.ac.uk/I/?d=fVGgzr-JQRem2LRYMUhsig>.

Table 3. Chromosomal pseudomolecules in the genome assembly of *Cyclophora albipunctata*, ilCycAlbi1.

| INSDC accession | Name | Length (Mb) | GC% |
|-----------------|------|-------------|------|
| OY720112.1 | 1 | 12.72 | 37.0 |
| OY720113.1 | 2 | 12.52 | 37.5 |
| OY720114.1 | 3 | 12.43 | 37.0 |
| OY720115.1 | 4 | 12.42 | 37.0 |
| OY720116.1 | 5 | 12.08 | 36.0 |
| OY720117.1 | 6 | 12.03 | 37.0 |
| OY720118.1 | 7 | 11.68 | 36.5 |
| OY720119.1 | 8 | 11.68 | 36.5 |
| OY720120.1 | 9 | 11.33 | 37.0 |
| OY720121.1 | 10 | 11.31 | 36.5 |
| OY720122.1 | 11 | 11.31 | 36.5 |
| OY720123.1 | 12 | 11.0 | 37.0 |
| OY720124.1 | 13 | 10.99 | 37.0 |
| OY720125.1 | 14 | 10.52 | 36.5 |
| OY720126.1 | 15 | 10.42 | 37.0 |

| INSDC accession | Name | Length (Mb) | GC% |
|-----------------|------|-------------|------|
| OY720127.1 | 16 | 10.41 | 37.0 |
| OY720128.1 | 17 | 10.22 | 36.5 |
| OY720129.1 | 18 | 9.95 | 37.5 |
| OY720131.1 | 19 | 9.47 | 37.0 |
| OY720132.1 | 20 | 9.33 | 37.0 |
| OY720133.1 | 21 | 8.96 | 38.0 |
| OY720134.1 | 22 | 8.15 | 37.0 |
| OY720135.1 | 23 | 7.87 | 38.5 |
| OY720136.1 | 24 | 7.7 | 37.5 |
| OY720137.1 | 25 | 7.62 | 37.5 |
| OY720138.1 | 26 | 7.49 | 38.5 |
| OY720139.1 | 27 | 6.95 | 37.5 |
| OY720140.1 | 28 | 5.72 | 38.5 |
| OY720141.1 | 29 | 4.88 | 40.5 |
| OY720142.1 | 30 | 4.84 | 39.0 |
| OY720130.1 | W | 9.13 | 39.0 |
| OY720111.1 | Z | 15.32 | 37.0 |
| OY720143.1 | MT | 0.02 | 17.5 |

Table 4. Software tools: versions and sources.

| Software tool | Version | Source |
|------------------------|-------------|---|
| BlobToolKit | 4.2.1 | https://github.com/blobtoolkit/blobtoolkit |
| BUSCO | 5.3.2 | https://gitlab.com/ezlab/busco |
| bwa-mem2 | 2.2.1 | https://github.com/bwa-mem2/bwa-mem2 |
| Cooler | 0.8.11 | https://github.com/open2c/cooler |
| Gfastats | 1.3.6 | https://github.com/vgl-hub/gfastats |
| Hifiasm | 0.19.8-r587 | https://github.com/chhylp123/hifiasm |
| HiGlass | 1.11.6 | https://github.com/higlass/higlass |
| Merqury | MerquryFK | https://github.com/thegenemyers/MERQURY.FK |
| MitoHiFi | 2 | https://github.com/marcelauliano/MitoHiFi |
| PretextView | 0.2 | https://github.com/wtsi-hpag/PretextView |
| purge_dups | 1.2.3 | https://github.com/dfguan/purge_dups |
| sanger-tol/ascc | - | https://github.com/sanger-tol/ascc |
| sanger-tol/genomenote | v1.0 | https://github.com/sanger-tol/genomenote |
| sanger-tol/readmapping | 1.1.0 | https://github.com/sanger-tol/readmapping/tree/1.1.0 |
| Singularity | 3.9.0 | https://github.com/sylabs/singularity |
| YaHS | 1.2a.2 | https://github.com/c-zhou/yahs |

of the materials themselves, and the circumstances under which they have been/are to be collected and provided for use. The purpose of this is to address and mitigate any potential legal and/or ethical implications of receipt and use of the materials as part of the research project, and to ensure that in doing so we align with best practice wherever possible. The overarching areas of consideration are:

- Ethical review of provenance and sourcing of the material
- Legality of collection, transfer and use (national and international)

Each transfer of samples is further undertaken according to a Research Collaboration Agreement or Material Transfer Agreement entered into by the Darwin Tree of Life Partner, Genome Research Limited (operating as the Wellcome Sanger Institute), and in some circumstances other Darwin Tree of Life collaborators.

Data availability

European Nucleotide Archive: *Cyclophora albipunctata* (birch mocha). Accession number PRJEB64086; <https://identifiers.org/ena.embl/PRJEB64086> (Wellcome Sanger Institute, 2023). The genome sequence is released openly for reuse.

The *Cyclophora albipunctata* genome sequencing initiative is part of the Darwin Tree of Life (DToL) project. All raw sequence data and the assembly have been deposited in INSDC databases. Raw data and assembly accession identifiers are reported in [Table 1](#) and [Table 2](#).

Author information

Members of the Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory team are listed here: <https://doi.org/10.5281/zenodo.12162482>.

Members of Wellcome Sanger Institute Scientific Operations: Sequencing Operations are listed here: <https://doi.org/10.5281/zenodo.12165051>.

Members of the Wellcome Sanger Institute Tree of Life Core Informatics team are listed here: <https://doi.org/10.5281/zenodo.12160324>.

Members of the Tree of Life Core Informatics collective are listed here: <https://doi.org/10.5281/zenodo.12205391>.

Members of the Darwin Tree of Life Consortium are listed here: <https://doi.org/10.5281/zenodo.4783558>.

References

- Abdennur N, Mirny LA: **Cooler: scalable storage for Hi-C data and other genomically labeled arrays.** *Bioinformatics.* 2020; **36**(1): 311–316.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Allio R, Schomaker-Bastos A, Romiguier J, et al.: **MitoFinder: efficient automated large-scale extraction of mitogenomic data in target enrichment phylogenomics.** *Mol Ecol Resour.* 2020; **20**(4): 892–905.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Bates A, Clayton-Lucey I, Howard C: **Sanger Tree of Life HMW DNA fragmentation: diagenode Megaruptor[®]3 for LI PacBio.** *protocols.io.* 2023.
[Publisher Full Text](#)
- Broad GR, Januszczak I, Natural History Museum Genome Acquisition Lab, et al.: **The genome sequence of the Maiden's Blush moth, *Cyclophora punctaria* (Linnaeus, 1758) [version 1; peer review: 1 approved].** *Wellcome Open Res.* 2024; **9**: 406.
[Publisher Full Text](#)
- Brůna T, Hoff KJ, Lomsadze A, et al.: **BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database.** *NAR Genom Bioinform.* 2021; **3**(1): lqaa108.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Challis R, Richards E, Rajan J, et al.: **BlobToolKit – interactive quality assessment of genome assemblies.** *G3 (Bethesda).* 2020; **10**(4): 1361–1374.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Cheng H, Concepcion GT, Feng X, et al.: **Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm.** *Nat Methods.* 2021; **18**(2): 170–175.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- da Veiga Leprevost F, Grüning BA, Alves Aflitos S, et al.: **BioContainers: an open-source and community-driven framework for software standardization.** *Bioinformatics.* 2017; **33**(16): 2580–2582.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Denton A, Oatley G, Cornwell C, et al.: **Sanger Tree of Life sample homogenisation: PowerMash.** *protocols.io.* 2023a.
[Publisher Full Text](#)
- Denton A, Yatsenko H, Jay J, et al.: **Sanger Tree of Life wet laboratory protocol collection V.1.** *protocols.io.* 2023b.
[Publisher Full Text](#)
- Diesch C, Stevens GJ, Xie P, et al.: **JBrowse 2: a modular genome browser with views of synteny and structural variation.** *Genome Biol.* 2023; **24**(1): 74.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- do Amaral RJV, Denton A, Yatsenko H, et al.: **Sanger Tree of Life RNA extraction: automated MagMax[™] mirVana.** *protocols.io.* 2023.
[Publisher Full Text](#)
- Ewels P, Magnusson M, Lundin S, et al.: **MultiQC: summarize analysis results for multiple tools and samples in a single report.** *Bioinformatics.* 2016; **32**(19): 3047–3048.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ewels PA, Peltzer A, Fillinger S, et al.: **The nf-core framework for community-curated bioinformatics pipelines.** *Nat Biotechnol.* 2020; **38**(3): 276–278.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Formenti G, Abueg L, Brajuka A, et al.: **Gfstats: conversion, evaluation and manipulation of genome sequences using assembly graphs.** *Bioinformatics.* 2022; **38**(17): 4214–4216.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Grüning B, Dale R, Sjödin A, et al.: **Bioconda: sustainable and comprehensive software distribution for the life sciences.** *Nat Methods.* 2018; **15**(7): 475–476.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Guan D, McCarthy SA, Wood J, et al.: **Identifying and removing haplotypic duplication in primary genome assemblies.** *Bioinformatics.* 2020; **36**(9): 2896–2898.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Harry E: **PretextView (Paired REad TEXTure Viewer): a desktop application for viewing pretext contact maps.** 2022.
[Reference Source](#)
- Howe K, Chow W, Collins J, et al.: **Significantly improving the quality of genome assemblies through curation.** *GigaScience.* 2021; **10**(1): g1aa153.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Jay J, Yatsenko H, Narváez-Gómez JP, et al.: **Sanger Tree of Life sample preparation: triage and dissection.** *protocols.io.* 2023.
[Publisher Full Text](#)
- Kerpedjiev P, Abdennur N, Lekschas F, et al.: **HiGlass: web-based visual exploration and analysis of genome interaction maps.** *Genome Biol.* 2018; **19**(1): 125.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Kurtzer GM, Sochat V, Bauer MW: **Singularity: scientific containers for mobility of compute.** *PLoS One.* 2017; **12**(5): e0177459.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Manni M, Berkeley MR, Seppely M, et al.: **BUSCO update: novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes.** *Mol Biol Evol.* 2021; **38**(10): 4647–4654.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Merkel D: **Docker: lightweight Linux containers for consistent development and deployment.** *Linux J.* 2014; **2014**(239): 2. [Accessed 2 April 2024].
[Reference Source](#)
- Oatley G, Denton A, Howard C: **Sanger Tree of Life HMW DNA extraction: automated MagAttract v.2.** *protocols.io.* 2023.
[Publisher Full Text](#)
- Randle Z, Evans-Hill LJ, Parsons MS, et al.: **Atlas of Britain & Ireland's Larger Moths.** Newbury: NatureBureau, 2019.
[Reference Source](#)
- Rao SSP, Huntley MH, Durand NC, et al.: **A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping.** *Cell.* 2014; **159**(7): 1665–1680.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Rhie A, McCarthy SA, Fedrigo O, et al.: **Towards complete and error-free genome assemblies of all vertebrate species.** *Nature.* 2021; **592**(7856): 737–746.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Rhie A, Walenz BP, Koren S, et al.: **Mercury: reference-free quality, completeness, and phasing assessment for genome assemblies.** *Genome Biol.* 2020; **21**(1): 245.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Strickland M, Cornwell C, Howard C: **Sanger Tree of Life fragmented DNA clean up: manual SPRI.** *protocols.io.* 2023.
[Publisher Full Text](#)
- Surana P, Muffato M, Qi G: **sanger-tol/readmapping: sanger-tol/readmapping v1.1.0 - Hebridean Black (1.1.0).** *Zenodo.* 2023a.
[Publisher Full Text](#)
- Surana P, Muffato M, Sadasivan Baby C: **sanger-tol/genomenote (v1.0.dev).** *Zenodo.* 2023b.
[Publisher Full Text](#)
- Uliano-Silva M, Ferreira JGRN, Krashennikova K, et al.: **MitoHiFi: a python pipeline for mitochondrial genome assembly from PacBio high fidelity reads.** *BMC Bioinformatics.* 2023; **24**(1): 288.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Vasimuddin M, Misra S, Li H, et al.: **Efficient architecture-aware acceleration of BWA-MEM for multicore systems.** In: *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS).* IEEE, 2019; 314–324.
[Publisher Full Text](#)
- Waring P, Townsend M, Lewington R: **Field Guide to the Moths of Great Britain and Ireland: Third Edition.** Bloomsbury Wildlife Guides, 2017.
[Reference Source](#)
- Wellcome Sanger Institute: **The genome sequence of the Birch Mocha moth, *Cyclophora albipunctata* (Hufnagel, 1767).** European Nucleotide Archive. [dataset], accession number PRJEB64086, 2023.
- Zhou C, McCarthy SA, Durbin R: **YaHS: yet another Hi-C scaffolding tool.** *Bioinformatics.* 2023; **39**(1): btac808.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)