

DATA NOTE

The genome sequence of the Ox-tongue Conch, Neocochylis

molliculana (Zeller, 1847) (Lepidoptera: Tortricidae)

[version 1; peer review: awaiting peer review]

Douglas Boyes¹⁺, Gavin R. Broad ¹⁰⁻², Clare Boyes³, University of Oxford and Wytham Woods Acquisition Lab, Natural History Museum Genome Acquisition Lab, Darwin Tree of Life Barcoding Collective, 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

V1

First published: 08 Sep 2025, **10**:490

https://doi.org/10.12688/wellcomeopenres.24868.1

Latest published: 08 Sep 2025, **10**:490

https://doi.org/10.12688/wellcomeopenres.24868.1

Abstract

We present a genome assembly from an individual female *Neocochylis molliculana* (Ox-tongue Conch; Arthropoda; Insecta; Lepidoptera; Tortricidae). The genome sequence has a total length of 420.54 megabases. Most of the assembly (99.69%) is scaffolded into 29 chromosomal pseudomolecules, including the W and Z sex chromosomes. The mitochondrial genome has also been assembled, with a length of 15.43 kilobases. This assembly was generated as part of the Darwin Tree of Life project, which produces reference genomes for eukaryotic species found in Britain and Ireland.

Keywords

Neocochylis molliculana; Ox-tongue Conch; genome sequence; chromosomal; Lepidoptera

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.

¹UK Centre for Ecology & Hydrology, Wallingford, England, UK

²Natural History Museum, London, England, UK

³Independent researcher, Welshpool, Wales, UK

⁺ Deceased author



This article is included in the Tree of Life gateway.

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

Author roles: Boyes D: Investigation, Resources; **Broad GR**: Investigation, Resources; **Boyes C**: Writing – Original Draft Preparation; **Competing interests:** No competing interests were disclosed.

Grant information: This work was supported by Wellcome through core funding to the Wellcome Sanger Institute (220540) 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: © 2025 Boyes D *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: Boyes D, Broad GR, Boyes C *et al.* The genome sequence of the Ox-tongue Conch, *Neocochylis molliculana* (Zeller, 1847) (Lepidoptera: Tortricidae) [version 1; peer review: awaiting peer review] Wellcome Open Research 2025, 10:490 https://doi.org/10.12688/wellcomeopenres.24868.1

First published: 08 Sep 2025, 10:490 https://doi.org/10.12688/wellcomeopenres.24868.1

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; Apoditrysia; Tortricoidea; Tortricidae; Tortricinae; Cochylini; Neocochylis; Neocochylis molliculana (Zeller, 1847) (NCBI:txid2769477)

Background

Neocochylis molliculana is a micro-moth in the family Tortricidae. It is a recent colonist of the UK, first recorded in 1993 and now spreading through England and occurs in a variety of habitats (Sterling *et al.*, 2023). It is believed to have colonised from southern Europe although there are only scattered records of this species in Europe, with the highest number of records from Portugal (GBIF Secretariat, 2025).

This small moth (forewing length 6–8 mm) is one of the 'bird-dropping' Tortix moths whose appearance offers some protection from predators. It is buff-brown with darker brown patterning which can be variable (Clifton & Wheeler, 2011). The foodplant is the Bristly Oxtongue and the moths flies from May to October with two generations (Sterling *et al.*, 2023).

We present a chromosome-level genome sequence for *Neocochylis molliculana*, the Ox-tongue Conch. This is the first genome for the genus *Neocochylis* and one of 100 genomes available for the family Tortricidae as of August 2025 (data obtained via NCBI datasets, O'Leary *et al.*, 2024). The assembly was produced using the Tree of Life pipeline from a specimen collected in Wytham Woods, Oxfordshire, United Kingdom (Figure 1). This assembly was generated as part of the Darwin Tree of Life Project, which aims to generate high-quality reference genomes for all named eukaryotic species in Britain and Ireland to support research, conservation, and the sustainable use of biodiversity (Blaxter *et al.*, 2022).



Figure 1. Photograph of the *Neocochylis molliculana* (ilNeoMoll1) specimen used for genome sequencing.

Methods

Sample acquisition and DNA barcoding

The specimen used for genome sequencing was an adult female *Neocochylis molliculana* (specimen ID Ox000808, ToLID ilNeoMoll1; Figure 1), collected from Wytham Woods, Oxfordshire, United Kingdom (latitude 51.772, longitude –1.338) on 2020-08-01. The specimen was collected and identified by Douglas Boyes. A second specimen was used for Hi-C and RNA sequencing (specimen ID NHMUK014584805, ToLID ilNeoMoll2). It was collected from Tudeley Woods, England, United Kingdom (latitude 51.16, longitude 0.31) on 2022-06-10. The specimen was collected by Gavin Broad and identified by David Lees. For the Darwin Tree of Life sampling and metadata approach, refer to Lawniczak *et al.* (2022).

The initial identification was verified by an additional DNA barcoding process according to the framework developed by Twyford *et al.* (2024). A small sample was dissected from the specimen and stored in ethanol, while the remaining parts were shipped on dry ice to the Wellcome Sanger Institute (WSI) (see the protocol). The tissue was lysed, the COI marker region was amplified by PCR, and amplicons were sequenced and compared to the BOLD database, confirming the species identification (Crowley *et al.*, 2023). Following whole genome sequence generation, the relevant DNA barcode region was also used alongside the initial barcoding data for sample tracking at the WSI (Twyford *et al.*, 2024). The standard operating procedures for Darwin Tree of Life barcoding are available on protocols.io.

Nucleic acid extraction

Protocols for high molecular weight (HMW) DNA extraction developed at the Wellcome Sanger Institute (WSI) Tree of Life Core Laboratory are available on protocols.io (Howard et al., 2025). The ilNeoMoll1 sample was weighed and triaged to determine the appropriate extraction protocol. Tissue from the whole organism was homogenised by powermashing using a PowerMasher II tissue disruptor.

RNA was extracted from tissue of ilNeoMoll2 in the Tree of Life Laboratory at the WSI using the RNA Extraction: Automated MagMaxTM mirVana protocol. 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.

PacBio HiFi library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core. Libraries were prepared using the SMRTbell Prep Kit 3.0 (Pacific Biosciences, California, USA), following the manufacturer's instructions. The kit includes reagents for end repair/A-tailing, adapter ligation, post-ligation SMRTbell bead clean-up, and nuclease treatment. Size selection and clean-up were performed using diluted AMPure PB beads (Pacific Biosciences). DNA concentration was quantified using a Qubit Fluorometer v4.0 (ThermoFisher Scientific) and the Qubit 1X dsDNA HS assay kit. Final library fragment

size was assessed with the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) using the gDNA 55 kb BAC analysis kit.

The sample was sequenced using the Sequel IIe system (Pacific Biosciences, California, USA). The concentration of the library loaded onto the Sequel IIe was in the range 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, and to perform primary and secondary analysis of the data upon completion.

Hi-C

Sample preparation and crosslinking

The Hi-C sample was prepared from 20–50 mg of frozen whole organism tissue of the ilNeoMoll2 sample using the Arima-HiC v2 kit (Arima Genomics). Following the manufacturer's instructions, tissue was fixed and DNA crosslinked using TC buffer to a final formaldehyde concentration of 2%. The tissue was homogenised using the Diagnocine Power Masher-II. Crosslinked DNA was digested with a restriction enzyme master mix, biotinylated, and ligated. Clean-up was performed with SPRISelect beads before library preparation. DNA concentration was measured with the Qubit Fluorometer (Thermo Fisher Scientific) and Qubit HS Assay Kit. The biotinylation percentage was estimated using the Arima-HiC v2 QC beads.

Hi-C library preparation and sequencing

Biotinylated DNA constructs were fragmented using a Covaris E220 sonicator and size selected to 400-600 bp using SPRISelect beads. DNA was enriched with Arima-HiC v2 kit Enrichment beads. End repair, A-tailing, and adapter ligation were carried out with the NEBNext Ultra II DNA Library Prep Kit (New England Biolabs), following a modified protocol where library preparation occurs while DNA remains bound to the Enrichment beads. Library amplification was performed using KAPA HiFi HotStart mix and a custom Unique Dual Index (UDI) barcode set (Integrated DNA Technologies). Depending on sample concentration and biotinylation percentage determined at the crosslinking stage, libraries were amplified with 10-16 PCR cycles. Post-PCR clean-up was performed with SPRISelect beads. Libraries were quantified using the AccuClear Ultra High Sensitivity dsDNA Standards Assay Kit (Biotium) and a FLUOstar Omega plate reader (BMG Labtech).

Prior to sequencing, libraries were normalised to 10 ng/µL. Normalised libraries were quantified again and equimolar and/or weighted 2.8 nM pools. Pool concentrations were checked using the Agilent 4200 TapeStation (Agilent) with High Sensitivity D500 reagents before sequencing. Sequencing was performed using paired-end 150 bp reads on the Illumina NovaSeq X.

RNA library preparation and sequencing

Libraries were prepared using the NEBNext® Ultra™ II Directional RNA Library Prep Kit for Illumina (New England Biolabs), following the manufacturer's instructions. Poly(A) mRNA in the total RNA solution was isolated using oligo(dT)

beads, converted to cDNA, and uniquely indexed; 14 PCR cycles were performed. Libraries were size-selected to produce fragments between 100–300 bp. Libraries were quantified, normalised, pooled to a final concentration of 2.8 nM, and diluted to 150 pM for loading. Sequencing was carried out on the Illumina NovaSeq X to generate 150-bp paired-end reads.

Genome assembly

Prior to assembly of the PacBio HiFi reads, a database of k-mer counts (k = 31) was generated from the filtered reads using FastK. GenomeScope2 (Ranallo-Benavidez *et al.*, 2020) was used to analyse the k-mer frequency distributions, providing estimates of genome size, heterozygosity, and repeat content.

The HiFi reads were assembled using Hifiasm (Cheng et al., 2021) with the --primary option. The Hi-C reads (Rao et al., 2014) were mapped to the primary contigs using bwa-mem2 (Vasimuddin et al., 2019), and the contigs were scaffolded in YaHS (Zhou et al., 2023) with the --break option for handling potential misassemblies. 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.

Assembly curation

The assembly was decontaminated using the Assembly Screen for Cobionts and Contaminants (ASCC) pipeline. TreeVal was used to generate the flat files and maps for use in curation. Manual curation was conducted primarily in PretextView and HiGlass (Kerpedjiev *et al.*, 2018). Scaffolds were visually inspected and corrected as described by Howe *et al.* (2021). The manual corrections included 10 joins. The curation process is documented at https://gitlab.com/wtsi-grit/rapid-curation. Pretext-Snapshot was used to generate a Hi-C contact map of the final assembly.

Assembly quality assessment

The Merqury.FK tool (Rhie *et al.*, 2020) was run in a Singularity container (Kurtzer *et al.*, 2017) to evaluate k-mer completeness and assembly quality for the primary and alternate haplotypes using the k-mer databases (k = 31) computed prior to genome assembly. The analysis outputs included assembly QV scores and completeness statistics.

The genome was analysed using the BlobToolKit pipeline, a Nextflow implementation of the earlier Snakemake version (Challis *et al.*, 2020). The pipeline aligns PacBio reads using minimap2 (Li, 2018) and SAMtools (Danecek *et al.*, 2021) to generate coverage tracks. It runs BUSCO (Manni *et al.*, 2021) using lineages identified from the NCBI Taxonomy (Schoch *et al.*, 2020). For the three domain-level lineages, BUSCO genes are aligned to the UniProt Reference Proteomes database

(Bateman et al., 2023) using DIAMOND blastp (Buchfink et al., 2021). The genome is divided into chunks based on the density of BUSCO genes from the closest taxonomic lineage, and each chunk is aligned to the UniProt Reference Proteomes database with DIAMOND blastx. Sequences without hits are chunked using seqtk and aligned to the NT database with blastn (Altschul et al., 1990). The BlobToolKit suite consolidates all outputs into a blobdir for visualisation. The BlobToolKit pipeline was developed using nf-core tooling (Ewels et al., 2020) and MultiQC (Ewels et al., 2016), with containerisation through Docker (Merkel, 2014) and Singularity (Kurtzer et al., 2017).

Genome sequence report

Sequence data

PacBio sequencing of the *Neocochylis molliculana* specimen generated 16.06 Gb (gigabases) from 1.34 million reads, which were used to assemble the genome. GenomeScope2.0 analysis estimated the haploid genome size at 394.66 Mb, with a heterozygosity of 0.62% and repeat content of 33.10% (Figure 2). These estimates guided expectations for the assembly. Based on the estimated genome size, the sequencing data provided approximately 39× coverage. Hi-C sequencing produced 98.95 Gb from 655.32 million reads, which were used to scaffold the assembly. RNA sequencing data were also generated and are available in public sequence repositories. Table 1 summarises the specimen and sequencing details.

Assembly statistics

The primary haplotype was assembled, and contigs corresponding to an alternate haplotype were also deposited in INSDC databases. The final assembly has a total length of 420.54 Mb in 71 scaffolds, with 22 gaps, and a scaffold N50 of 14.43 Mb (Table 2).

Most of the assembly sequence (99.69%) was assigned to 29 chromosomal-level scaffolds, representing 27 autosomes and the W and Z sex chromosomes. These chromosome-level scaffolds, confirmed by Hi-C data, are named according to size (Figure 3; Table 3). Chromosome Z was identified by BUSCO gene painting with ancestral Merian elements (Wright *et al.*, 2024). Chromosome W was identified by read coverage.

The mitochondrial genome was also assembled. This sequence is included as a contig in the multifasta file of the genome submission and as a standalone record.

The combined primary and alternate assemblies achieve an estimated QV of 62.3. The *k*-mer completeness is 89.45% for the primary assembly, 64.35% for the alternate haplotype, and 98.78% for the combined assemblies (Figure 4).

BUSCO v.5.5.0 analysis using the lepidoptera_odb10 reference set (n = 5286) identified 98.3% of the expected gene set

GenomeScope Profile

len:394,660,592bp uniq:67% aa:99.4% ab:0.617% kcov:19.3 err:0.164% dup:0.454 k:31 p:2

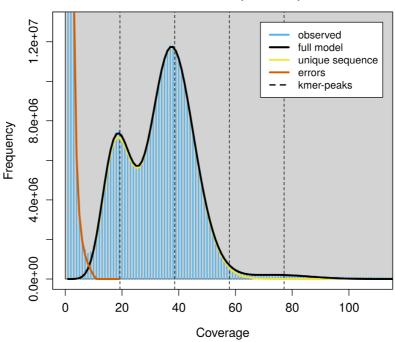


Figure 2. Frequency distribution of *k***-mers generated using GenomeScope2.** The plot shows observed and modelled *k*-mer spectra, providing estimates of genome size, heterozygosity, and repeat content based on unassembled sequencing reads.

Table 1. Specimen and sequencing data for BioProject PRJEB71298.

Platform	PacBio HiFi	Hi-C	RNA-seq
ToLID	ilNeoMoll1	ilNeoMoll2	ilNeoMoll2
Specimen ID	Ox000808	NHMUK014584805	NHMUK014584805
BioSample (source individual)	SAMEA7746615	SAMEA115574939	SAMEA115574939
BioSample (tissue)	SAMEA7746702	SAMEA115575058	SAMEA115575058
Tissue	whole organism	whole organism	whole organism
Instrument	Sequel IIe	Illumina NovaSeq X	Illumina NovaSeq X
Run accessions	ERR12370425	ERR13623229	ERR13999091
Read count total	1.34 million	655.32 million	81.61 million
Base count total	16.06 Gb	98.95 Gb	12.32 Gb

Table 2. Genome assembly statistics.

Assembly name	ilNeoMoll1.1
Assembly accession	GCA_964656585.1
Alternate haplotype accession	GCA_964656565.1
Assembly level	chromosome
Span (Mb)	420.54
Number of chromosomes	29
Number of contigs	93
Contig N50	11.51 Mb
Number of scaffolds	71
Scaffold N50	14.43 Mb
Sex chromosomes	W and Z
Organelles	Mitochondrion: 15.43 kb

(single = 97.9%, duplicated = 0.4%). The snail plot in Figure 5 summarises the scaffold length distribution and other assembly statistics for the primary assembly. The blob plot in Figure 6 shows the distribution of scaffolds by GC proportion and coverage.

Table 4 lists the assembly metric benchmarks adapted from Rhie *et al.* (2021) the Earth BioGenome Project Report on Assembly Standards September 2024. The EBP metric, calculated for the primary assembly, is **7.C.Q64**, meeting the recommended reference standard.

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. 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 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

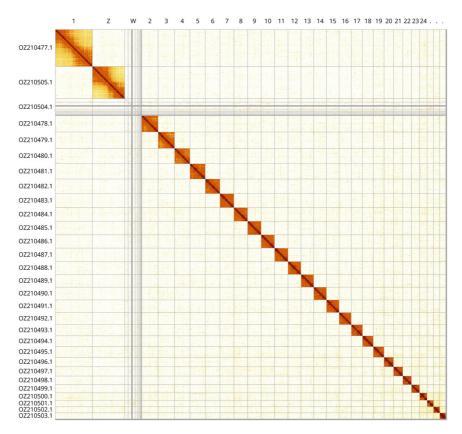


Figure 3. Hi-C contact map of the *Neocochylis molliculana* **genome assembly.** Assembled chromosomes are shown in order of size and labelled along the axes. The plot was generated using PretextSnapshot.

Table 3. Chromosomal pseudomolecules in the primary genome assembly of *Neocochylis molliculana* ilNeoMoll1.

INSDC accession	Molecule	Length (Mb)	GC%
OZ210477.1	1	39.92	37.50
OZ210478.1	2	17.68	38
OZ210479.1	3	17.53	38
OZ210480.1	4	16.60	38.50
OZ210481.1	5	16.41	37.50
OZ210482.1	6	15.78	37.50
OZ210483.1	7	15.01	38
OZ210484.1	8	14.67	37.50
OZ210485.1	9	14.50	38
OZ210486.1	10	14.43	37.50
OZ210487.1	11	14.39	38
OZ210488.1	12	13.91	37.50
OZ210489.1	13	13.71	38

INSDC accession	Molecule	Length (Mb)	GC%
OZ210490.1	14	13.68	37.50
OZ210491.1	15	13.60	37.50
OZ210492.1	16	13.08	38.50
OZ210493.1	17	12.05	38
OZ210494.1	18	11.56	38
OZ210495.1	19	11.52	38.50
OZ210496.1	20	10.14	37.50
OZ210497.1	21	9.95	37.50
OZ210498.1	22	9.33	38
OZ210499.1	23	8.75	37.50
OZ210500.1	24	8.05	38
OZ210501.1	25	6.84	39
OZ210502.1	26	6.68	38
OZ210503.1	27	6.54	37.50
OZ210504.1	W	18.31	38.50
OZ210505.1	Z	34.60	38

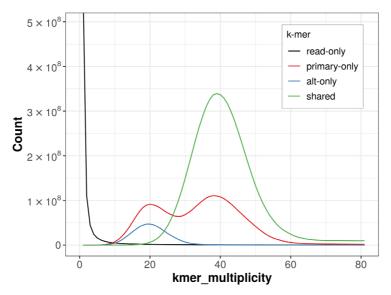


Figure 4. Evaluation of *k***-mer completeness using MerquryFK.** This plot illustrates the recovery of *k*-mers from the original read data in the final assemblies. The horizontal axis represents *k*-mer multiplicity, and the vertical axis shows the number of *k*-mers. The black curve represents *k*-mers that appear in the reads but are not assembled. The green curve corresponds to *k*-mers shared by both haplotypes, and the red and blue curves show *k*-mers found only in one of the haplotypes.

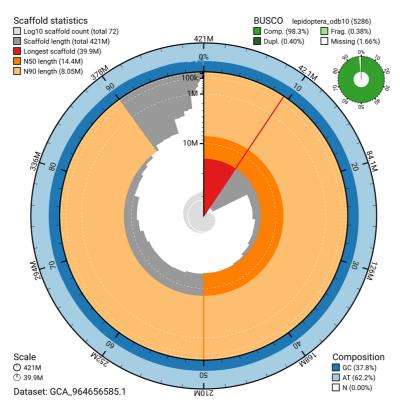


Figure 5. Assembly metrics for ilNeoMoll1.1. The BlobToolKit snail plot provides an overview of assembly metrics and BUSCO gene completeness. The circumference represents the length of the whole genome sequence, and the main plot is divided into 1 000 bins around the circumference. The outermost blue tracks display the distribution of GC, AT, and N percentages across the bins. Scaffolds are arranged clockwise from longest to shortest and are depicted in dark grey. The longest scaffold is indicated by the red arc, and the deeper orange and pale orange arcs represent the N50 and N90 lengths. A light grey spiral at the centre shows the cumulative scaffold count on a logarithmic scale. A summary of complete, fragmented, duplicated, and missing BUSCO genes in the lepidoptera_odb10 set is presented at the top right. An interactive version of this figure can be accessed on the BlobToolKit viewer.

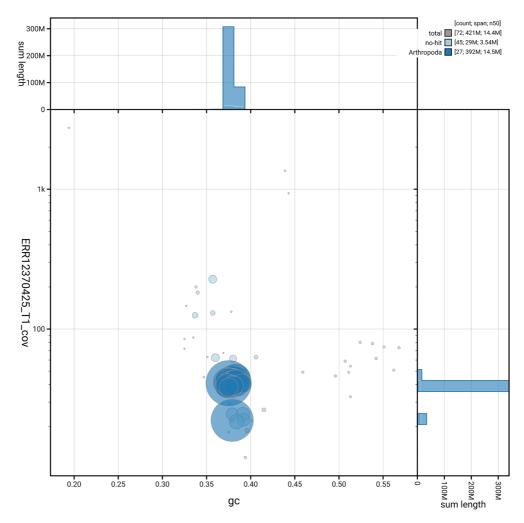


Figure 6. BlobToolKit GC-coverage plot for ilNeoMoll1.1. Blob plot showing sequence coverage (vertical axis) and GC content (horizontal axis). The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. The histograms along the axes display the total length of sequences distributed across different levels of coverage and GC content. An interactive version of this figure is available on the BlobToolKit viewer.

Table 4. Earth Biogenome Project summary metrics for the Neocochylis molliculana assembly.

Measure	Value	Benchmark
EBP summary (primary)	7.C.Q64	6.C.Q40
Contig N50 length	11.51 Mb	≥ 1 Mb
Scaffold N50 length	14.43 Mb	= chromosome N50
Consensus quality (QV)	Primary: 64.1; alternate: 61.4; combined: 62.3	≥ 40
k-mer completeness	Primary: 89.45%; alternate: 64.35%; combined: 98.78%	≥ 95%
BUSCO	C:98.3% [S:97.9%; D:0.4%]; F:0.4%; M:1.3%; n:5 286	S > 90%; D < 5%
Percentage of assembly assigned to chromosomes	99.69%	≥ 90%

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: Neocochylis molliculana (ox-tongue conch). Accession number PRJEB71298. The genome sequence is released openly for reuse. The *Neocochylis molliculana* genome sequencing initiative is part of the Darwin Tree of Life Project (PRJEB40665), the Sanger Institute Tree of Life Programme (PRJEB43745) and Project Psyche (PRJEB71705). All raw sequence data and the assembly have been deposited in INSDC databases. The genome will be annotated using available RNA-Seq data and presented through the Ensembl pipeline at the European Bioinformatics Institute.

Raw data and assembly accession identifiers are reported in Table 1 and Table 2.

Production code used in genome assembly at the WSI Tree of Life is available at https://github.com/sanger-tol. Table 5 lists software versions used in this study.

Author information

Contributors are listed at the following links:

- Members of the University of Oxford and Wytham Woods Genome Acquisition Lab
- Members of the Natural History Museum Genome Acquisition Lab
- Members of the Darwin Tree of Life Barcoding collective
- Members of the Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory team
- Members of Wellcome Sanger Institute Scientific Operations – Sequencing Operations
- Members of the Wellcome Sanger Institute Tree of Life Core Informatics team
- Members of the Tree of Life Core Informatics collective
- Members of the Darwin Tree of Life Consortium

Table 5. Software versions and sources.

Software	Version	Source
BEDTools	2.30.0	https://github.com/arq5x/bedtools2
BLAST	2.14.0	ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/
BlobToolKit	4.3.9	https://github.com/blobtoolkit/blobtoolkit
BUSCO	5.5.0	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
DIAMOND	2.1.8	https://github.com/bbuchfink/diamond
fasta_windows	0.2.4	https://github.com/tolkit/fasta_windows
FastK	1.1	https://github.com/thegenemyers/FASTK
GenomeScope2.0	2.0.1	https://github.com/tbenavi1/genomescope2.0
Gfastats	1.3.6	https://github.com/vgl-hub/gfastats
GoaT CLI	0.2.5	https://github.com/genomehubs/goat-cli
Hifiasm	0.19.8-r603	https://github.com/chhylp123/hifiasm
HiGlass	1.13.4	https://github.com/higlass/higlass
MerquryFK	1.1.2	https://github.com/thegenemyers/MERQURY.FK
Minimap2	2.24-r1122	https://github.com/lh3/minimap2

Software	Version	Source
MitoHiFi	3	https://github.com/marcelauliano/MitoHiFi
MultiQC	1.14; 1.17 and 1.18	https://github.com/MultiQC/MultiQC
Nextflow	23.10.0	https://github.com/nextflow-io/nextflow
PretextSnapshot	N/A	https://github.com/sanger-tol/PretextSnapshot
PretextView	0.2.5	https://github.com/sanger-tol/PretextView
samtools	1.19.2	https://github.com/samtools/samtools
sanger-tol/ascc	0.1.0	https://github.com/sanger-tol/ascc
sanger-tol/blobtoolkit	0.6.0	https://github.com/sanger-tol/blobtoolkit
sanger-tol/curationpretext	1.4.2	https://github.com/sanger-tol/curationpretext
Seqtk	1.3	https://github.com/lh3/seqtk
Singularity	3.9.0	https://github.com/sylabs/singularity
TreeVal	1.4.0	https://github.com/sanger-tol/treeval
YaHS	1.2a.2	https://github.com/c-zhou/yahs

References

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

Altschul SF, Gish W, Miller W, et al.: Basic Local Alignment Search Tool. J Mol Biol. 1990; **215**(3): 403–410. **PubMed Abstract | Publisher Full Text**

Bateman A, Martin MJ, Orchard S, et al.: UniProt: the universal protein knowledgebase in 2023. Nucleic Acids Res. 2023; 51(D1): D523-D531. PubMed Abstract | Publisher Full Text | Free Full Text

Blaxter M, Mieszkowska N, Di Palma F, et al.: Sequence locally, think globally: the Darwin Tree of Life Project. Proc Natl Acad Sci U S A. 2022; 119(4): e2115642118.

PubMed Abstract | Publisher Full Text | Free Full Text

Buchfink B, Reuter K, Drost HG: Sensitive protein alignments at Tree-of-Life scale using DIAMOND. Nat Methods. 2021; 18(4): 366-368.

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):

PubMed Abstract | Publisher Full Text | Free Full Text

Clifton J, Wheeler J: Bird-dropping tortrix moths of the British Isles. Dorchester: Clifton & Wheeler, 2011.

Crowley L, Allen H, Barnes I, et al.: A sampling strategy for genome sequencing the British terrestrial Arthropod fauna [version 1; peer review: 2 approved]. Wellcome Open Res. 2023; 8: 123.

PubMed Abstract | Publisher Full Text | Free Full Text

Danecek P, Bonfield JK, Liddle J, et al.: Twelve years of SAMtools and BCFtools. GigaScience. 2021; 10(2): giab008

PubMed Abstract | Publisher Full Text | Free 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 communitycurated bioinformatics pipelines. *Nat Biotechnol*. 2020; **38**(3): 276–278. PubMed Abstract | Publisher Full Text

Formenti G, Abueg L, Brajuka A, et al.: Gfastats: conversion, evaluation and

manipulation of genome sequences using assembly graphs. Bioinformatics. 2022; 38(17): 4214–4216.

PubMed Abstract | Publisher Full Text | Free Full Text

GBIF Secretariat: Neocochylis molliculana (Zeller, 1847). in the GBIF Backbone Taxonomy. 2025.

Howard C, Denton A, Jackson B, et al.: On the path to reference genomes for all biodiversity: lessons learned and laboratory protocols created in the Sanger Tree of Life core laboratory over the first 2000 species. *bioRxiv.* 2025. **Publisher Full Text**

Howe K, Chow W, Collins J, et al.: Significantly improving the quality of genome assemblies through curation. *GigaScience*. 2021; **10**(1): giaa153. PubMed Abstract | Publisher Full Text | Free 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

Lawniczak MKN, Davey RP, Rajan J, et al.: Specimen and sample metadata standards for biodiversity genomics: a proposal from the darwin Tree of Life Project [version 1; peer review: 2 approved with reservations]. Wellcome Open Res. 2022; 7: 187. **Publisher Full Text**

Li H: Minimap2: pairwise alignment for nucleotide sequences.

Bioinformatics. 2018; 34(18): 3094–3100.

PubMed Abstract | Publisher Full Text | Free Full Text

Manni M, Berkeley MR, Seppey 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.

Reference Source

O'Leary NA, Cox E, Holmes JB, et al.: Exploring and retrieving sequence and metadata for species across the Tree of Life with NCBI datasets. Sci Data. 2024; 11(1): 732.

PubMed Abstract | Publisher Full Text | Free Full Text

Ranallo-Benavidez TR, Jaron KS, Schatz MC: GenomeScope 2.0 and

Smudgeplot for reference-free profiling of polyploid genomes. *Nat Commun.* 2020; **11**(1): 1432. **PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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.: Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. Genome Biol. 2020; 21(1): 245

PubMed Abstract | Publisher Full Text | Free Full Text

Schoch CL, Ciufo S, Domrachev M, et al.: NCBI Taxonomy: a comprehensive update on curation, resources and tools. Database (Oxford). 2020; 2020: baaa062.

PubMed Abstract | Publisher Full Text | Free Full Text

Sterling P, Parsons M, Lewington R: Field guide to the micro moths of great

Britain and Ireland. Dorset: British Wildlife Publishing, 2023.

Twyford AD, Beasley J, Barnes I, et al.: A DNA barcoding framework for taxonomic verification in the Darwin Tree of Life Project [version 1; peer review: 2 approved]. Wellcome Open Res. 2024; 9: 339.
PubMed Abstract | Publisher Full Text | Free Full Text

Uliano-Silva M, Ferreira JGRN, Krasheninnikova 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**

Wright CJ, Stevens L, Mackintosh A, et al.: Comparative genomics reveals the dynamics of chromosome evolution in Lepidoptera. Nat Ecol Evol. 2024; 8(4): 777-90

PubMed Abstract | Publisher Full Text | Free Full Text

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