



DATA NOTE

# The genome sequence of the Small Grey moth, *Eudonia*

## *mercurella* Linnaeus, 1758

[version 1; peer review: awaiting peer review]

Douglas Boyes<sup>1+</sup>, Maxwell V. L. Barclay<sup>2</sup>, Lily V. M. Lewis<sup>3</sup>,  
University of Oxford and Wytham Woods Genome 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

<sup>1</sup>UK Centre for Ecology & Hydrology, Wallingford, England, UK<sup>2</sup>Natural History Museum, London, England, UK<sup>3</sup>The University of Sheffield, Sheffield, England, UK

+ Deceased author

---

**V1** First published: 19 Feb 2025, 10:76  
<https://doi.org/10.12688/wellcomeopenres.23699.1>  
Latest published: 19 Feb 2025, 10:76  
<https://doi.org/10.12688/wellcomeopenres.23699.1>

---

### 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 *Eudonia mercurella* (Small Grey; Arthropoda; Insecta; Lepidoptera; Crambidae). The genome sequence has a total length of 591.50 megabases. Most of the assembly (96.86%) is scaffolded into 33 chromosomal pseudomolecules, including the Z sex chromosome. The mitochondrial genome has also been assembled and is 15.31 kilobases in length. Gene annotation of this assembly on Ensembl identified 13,075 protein-coding genes.

### Keywords

*Eudonia mercurella*, small grey, 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](mailto:mark.blaxter@sanger.ac.uk))

**Author roles:** **Boyes D:** Investigation, Resources; **Barclay MVL:** Investigation, Resources; **Lewis LVM:** 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 [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:** © 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, Barclay MVL, Lewis LVM *et al.* **The genome sequence of the Small Grey moth, *Eudonia mercurella* Linnaeus, 1758 [version 1; peer review: awaiting peer review]** Wellcome Open Research 2025, 10:76 <https://doi.org/10.12688/wellcomeopenres.23699.1>

**First published:** 19 Feb 2025, 10:76 <https://doi.org/10.12688/wellcomeopenres.23699.1>

## 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; Pyraloidea; Crambidae; Scopariinae; *Eudonia*; *Eudonia mercurella* Linnaeus, 1758 (NCBI:txid1100992).

## Background

The Small Grey or Garden Grey, *Eudonia mercurella*, is a moth in the family Crambidae, subfamily Scopariinae. The species is variable in colour and pattern but generally shaded speckled dark brown or grey with a pale postmedian line (Sterling *et al.*, 2023).

Most records of *Eudonia mercurella* in the UK are from the south, but the distribution extends to northern Scotland (Sterling *et al.*, 2023). It has been recorded widely across much of Europe (GBIF Secretariat, 2023).

In the UK, adults are often found in urban gardens and parks, grasslands and moorlands from June to mid-October. Larvae feed on mosses growing on tree trunks, rocks and walls. From September to April, larvae can be seen in silk tubes on these mosses (Sterling *et al.*, 2023).

The genome of *Eudonia mercurella* 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. The sequence data collected will contribute to a growing data set for understanding lepidopteran biology.

## Genome sequence report

The genome of *Eudonia mercurella* (Figure 1) was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating a total of 21.30 Gb (gigabases) from 2.15 million reads, providing an estimated 37-fold coverage. Chromosome



**Figure 1.** Photograph of the *Eudonia mercurella* (ilEudMerc1) specimen used for genome sequencing.

conformation Hi-C sequencing produced 115.40 Gb from 764.22 million reads. Specimen and sequencing details are summarised in Table 1.

Assembly errors were corrected by manual curation, including 166 missing joins or mis-joins and 8 haplotypic duplications. This reduced the scaffold number by 9.48% and increased the scaffold N50 by 1.77%. The final assembly has a total length of 591.50 Mb in 553 sequence scaffolds, with 1,710 gaps, and a scaffold N50 of 19.1 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 (96.86%) was assigned to 33 chromosomal-level scaffolds, representing 30 autosomes and the Z sex chromosome. These chromosome-level scaffolds, confirmed by the Hi-C data, are named in order of size (Figure 5; Table 3). During manual curation the Z chromosome was assigned by the observed half coverage of PacBio reads. This appears to be a ZO female moth.

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.

The final assembly has a Quality Value (QV) of 58.1. The *k*-mer completeness values were 78.83% for the primary assembly, 72.68% for the alternate haplotype, and 98.48% for the combined assemblies. BUSCO (5.3.2) analysis using the lepidoptera\_odb10 reference set ( $n = 5,286$ ) indicated a completeness score of 94.9% (single = 93.9%, duplicated = 1.0%). The assembly achieves the EBP reference standard of 5.C.Q58. Other quality metrics are given in Table 2.

## Genome annotation report

The *Eudonia mercurella* genome assembly (GCA\_963082485.1) was annotated at the European Bioinformatics Institute (EBI) on Ensembl Rapid Release. The resulting annotation includes 24,087 transcribed mRNAs from 13,075 protein-coding and 2,316 non-coding genes (Table 2; [https://rapid.ensembl.org/Eudonia\\_mercurella\\_GCA\\_963082485.1/Info/Index](https://rapid.ensembl.org/Eudonia_mercurella_GCA_963082485.1/Info/Index)). The average transcript length is 17,745.90. There are 1.57 coding transcripts per gene and 7.00 exons per transcript.

## Methods

### Sample acquisition and DNA barcoding

An adult female *Eudonia mercurella* (specimen ID Ox000941, ToLID ilEudMerc1) was collected from Wytham Woods, Berkshire, United Kingdom (latitude 51.77, longitude -1.34) on 2020-08-31, using a light trap. The specimen was collected

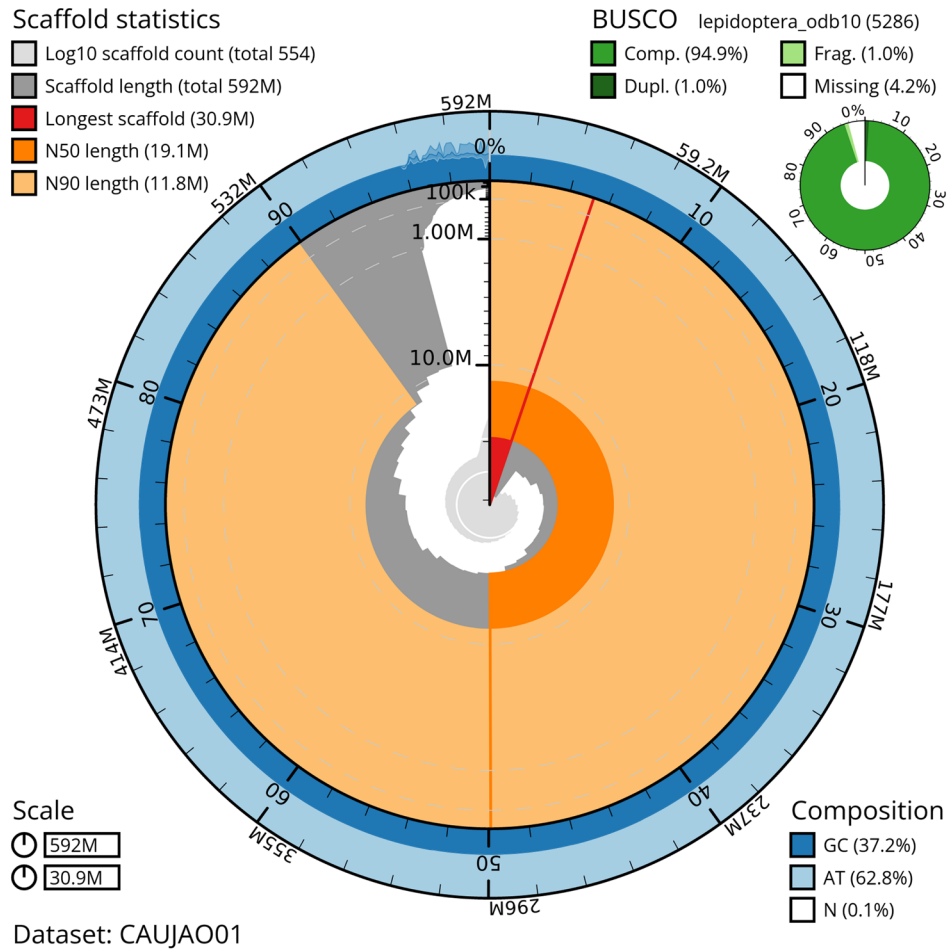
**Table 1. Specimen and sequencing data for *Eudonia mercurella*.**

Project information			
Study title	Eudonia mercurella (small grey)		
Umbrella BioProject	PRJEB62417		
Species	<i>Eudonia mercurella</i>		
BioSample	SAMEA7701587		
NCBI taxonomy ID	1100992		
Specimen information			
Technology	ToLID	BioSample accession	Organism part
PacBio long read sequencing	ilEudMerc1	SAMEA7701787	Whole organism
Hi-C sequencing	ilEudMerc2	SAMEA112222388	Whole organism
RNA sequencing	ilEudMerc2	SAMEA112222388	Whole organism
Sequencing information			
Platform	Run accession	Read count	Base count (Gb)
Hi-C Illumina NovaSeq 6000	ERR11468760	7.64e+08	115.4
PacBio Sequel Iie	ERR11458826	2.15e+06	21.3
RNA Illumina NovaSeq 6000	ERR11837493	6.92e+07	10.45

**Table 2. Genome assembly data for *Eudonia mercurella*, ilEudMerc1.1.**

Genome assembly		
Assembly name	ilEudMerc1.1	
Assembly accession	GCA_963082485.1	
Accession of alternate haplotype	GCA_963082515.1	
Span (Mb)	591.50	
Number of contigs	2,264	
Number of scaffolds	553	
Longest scaffold (Mb)	31.11	
Assembly metrics*	Benchmark	
Contig N50 length (Mb)	0.6	≥ 1 Mb
Scaffold N50 length (Mb)	19.1	= chromosome N50
Consensus quality (QV)	58.1	≥ 40
k-mer completeness	Primary: 78.83%; alternate: 72.68%; combined: 98.48%	≥ 95%
BUSCO v5.4.3 lineage lepidoptera_odb10	C:94.9%[S:93.9%,D:1.0%], F:1.0%,M:4.1%,n:5,286	S > 90%, D < 5%
Percentage of assembly mapped to chromosomes	96.86%	≥ 90%
Sex chromosomes	ZO	localised homologous pairs
Organelles	Mitochondrial genome: 15.31 kb	complete single alleles
Genome annotation of assembly GCA_963082485.1 at Ensembl		
Number of protein-coding genes	13,075	
Number of non-coding genes	2,316	
Number of gene transcripts	24,087	

\* Assembly metric benchmarks are adapted from [Rhie et al. \(2021\)](#) and the Earth BioGenome Project Report on Assembly Standards [September 2024](#).



**Figure 2. Genome assembly of *Eudonia mercurella*, ilEudMerc1.1: metrics.** The BlobToolKit snail plot shows N50 metrics and BUSCO gene completeness. \$BTK\_SNAIL\_LEG An interactive version of this figure is available at <https://blobtoolkit.genomehubs.org/view/Eudonia%20mercurella/dataset/CAUJAO01/snail/>.

and identified by Douglas Boyes (University of Oxford) and preserved on dry ice.

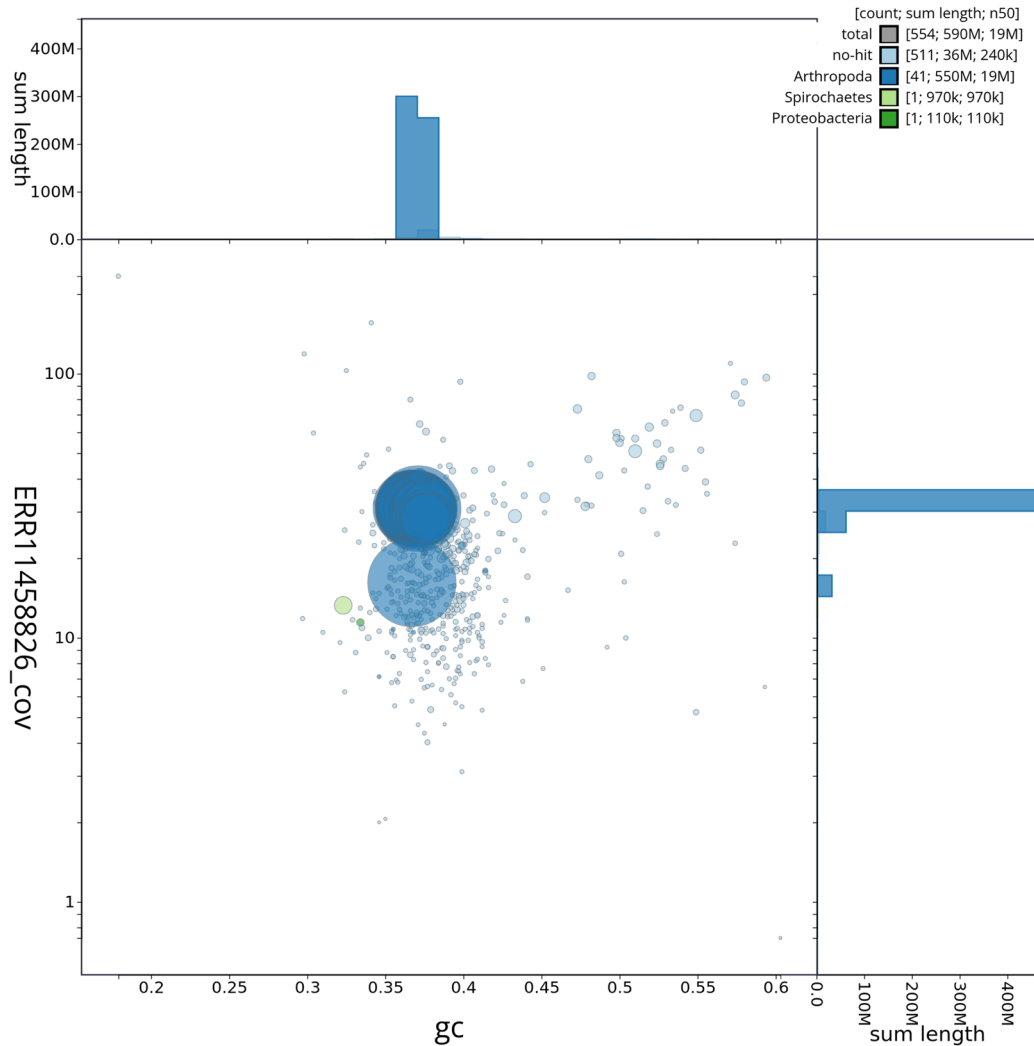
The specimens used for Hi-C sequencing (specimen ID NHMUK013696792, ToLID ilEudMerc2) and RNA sequencing (specimen ID NHMUK013696792, ToLID ilEudMerc2) were collected from the Natural History Museum Wildlife Garden, London, England, United Kingdom (latitude 51.5, longitude -0.18) on 2021-06-29, using an aerial net. The specimen was collected and identified by Maxwell Barclay (Natural History Museum) and preserved by dry freezing (-80 °C).

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 specimens and stored in ethanol, while the remaining parts were shipped on dry ice to the Wellcome Sanger Institute (WSI). 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 have been deposited on protocols.io (Beasley *et al.*, 2023).

#### 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 ilEudMerc1 sample was prepared for DNA extraction by weighing and dissecting it on dry ice (Jay *et al.*, 2023). Tissue from the whole organism was homogenised using a PowerMasher II tissue disruptor (Denton *et al.*, 2023a). HMW DNA was extracted using the



**Figure 3. Genome assembly of *Eudonia mercurella*, iLEudMerc1.1: 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/Eudonia%20mercurella/dataset/CAUJAO01/blob>.

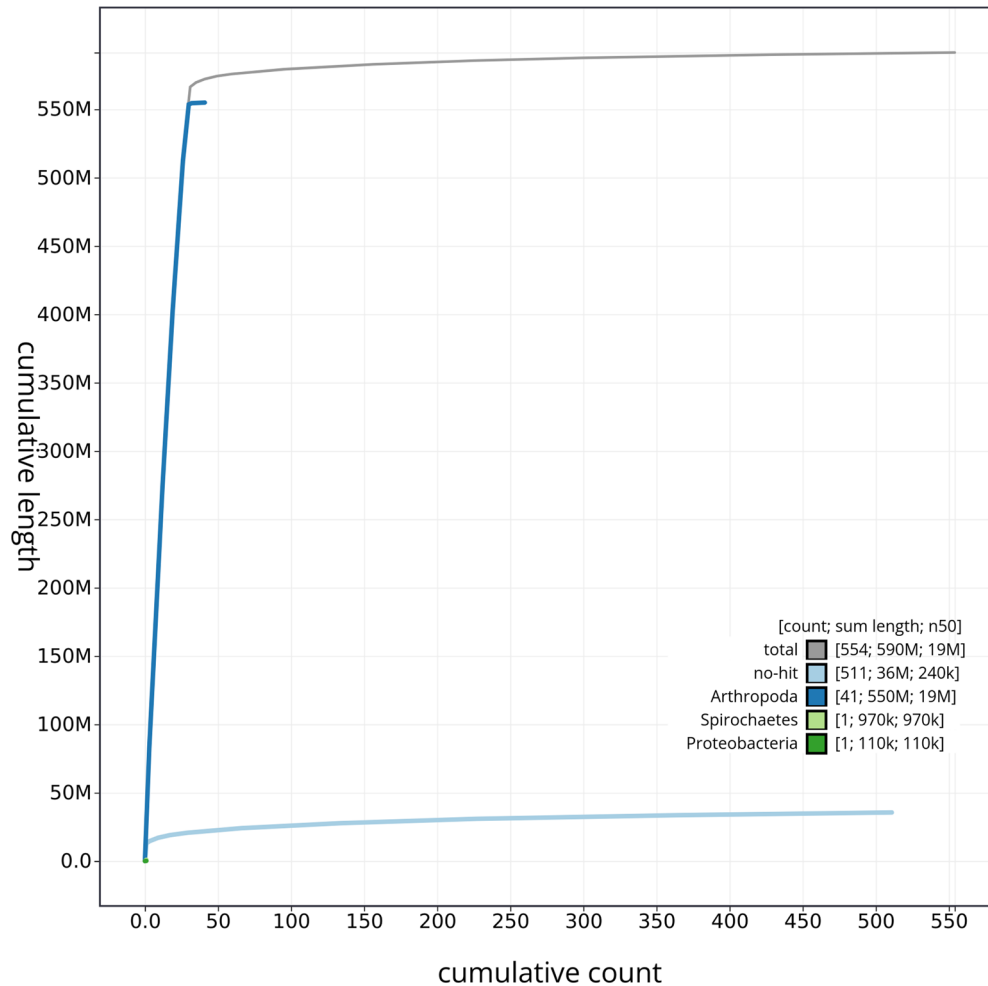
Automated MagAttract v1 protocol (Sheerin *et al.*, 2023). For ultra-low input (ULI) PacBio sequencing, DNA was fragmented using the Covaris g-TUBE method (Oatley *et al.*, 2023). Sheared DNA was purified by solid-phase reversible immobilisation, using AMPure PB beads to eliminate shorter fragments and concentrate the DNA (Strickland *et al.*, 2023). The concentration of the sheared and purified DNA was assessed using a Nanodrop spectrophotometer and a Qubit Fluorometer using the Qubit dsDNA High Sensitivity Assay kit. The fragment size distribution was evaluated by running the sample on the FemtoPulse system.

RNA was extracted from iLEudMerc2 in the Tree of Life Laboratory at the WSI using the RNA Extraction: Automated MagMax™ *mir*-Vana 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 sample preparation

Tissue from the iLEudMerc2 sample was processed at the WSI Scientific Operations core, using the Arima-HiC v2 kit. Tissue (stored at  $-80^{\circ}\text{C}$ ) was fixed, and the DNA crosslinked using a TC buffer with 22% formaldehyde. After crosslinking, the tissue was homogenised using the Diagenode 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



**Figure 4. Genome assembly of *Eudonia mercurella* iEudMerc1.1: BlobToolKit cumulative sequence plot.** The grey line shows cumulative length for all scaffolds. Coloured lines show cumulative lengths of scaffolds assigned to each phylum using the buscodegenes taxrule. An interactive version of this figure is available at <https://blobtoolkit.genomehubs.org/view/Eudonia%20mercurella/dataset/CAUJAO01/cumulative>.

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.

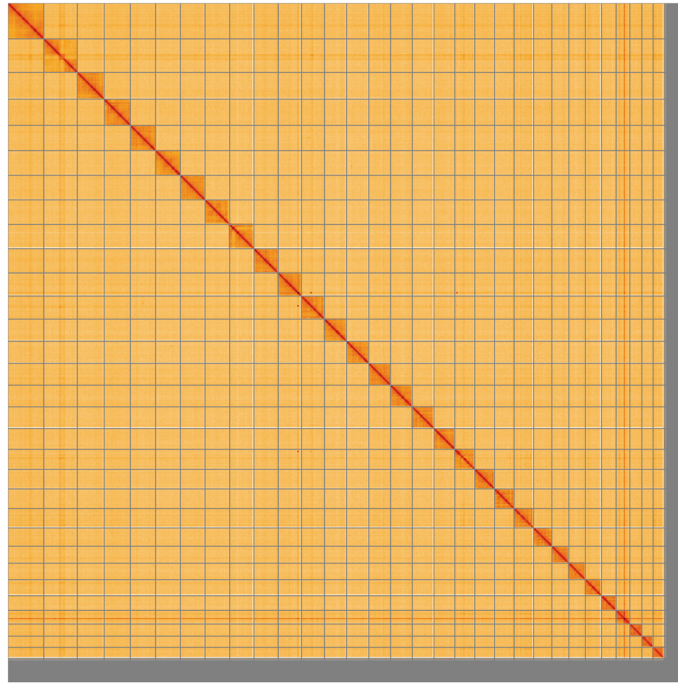
#### Library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core.

#### **PacBio HiFi**

The sample requires Covaris g-TUBE shearing to approximately 10 kb prior to library preparation. Ultra-low input libraries were prepared using PacBio SMRTbell® Express Template Prep Kit 2.0 and PacBio SMRTbell® gDNA Sample Amplification Kit. To begin, samples were normalised to 20 ng of DNA. Initial removal of single-strand overhangs, DNA damage repair, and

end repair/A-tailing were performed per manufacturer's instructions. From the SMRTbell® gDNA Sample Amplification Kit, amplification adapters were then ligated. A 0.85X pre-PCR clean-up was performed with Promega ProNex beads and the sample was then divided into two for a dual PCR. PCR reactions A and B each followed the PCR programs as described in the manufacturer's protocol. A 0.85X post-PCR clean-up was performed with ProNex beads for PCR reactions A and B and DNA concentration was quantified using the Qubit Fluorometer v4.0 (Thermo Fisher Scientific) and Qubit HS Assay Kit and fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) and gDNA 55kb BAC analysis kit. PCR reactions A and B were then pooled, ensuring the total mass was  $\geq 500$  ng in 47.4  $\mu$ l. The pooled sample then repeated the process for DNA damage repair, end repair/A-tailing and additional hairpin adapter ligation. A 1X clean-up was performed



**Figure 5. Genome assembly of *Eudonia mercurella* iEudMerc1.1: Hi-C contact map of the iEudMerc1.1 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=Uub4EoYSS-OI5PtiiEkd4A>.

**Table 3. Chromosomal pseudomolecules in the genome assembly of *Eudonia mercurella*, iEudMerc1.**

INSDC accession	Name	Length (Mb)	GC%
OY720023.1	1	28.96	37.0
OY720024.1	2	23.14	37.0
OY720025.1	3	22.74	37.0
OY720026.1	4	21.72	36.5
OY720027.1	5	21.38	37.0
OY720028.1	6	21.24	37.0
OY720029.1	7	21.28	37.0
OY720030.1	8	20.89	37.0
OY720031.1	9	20.96	36.5
OY720032.1	10	20.1	36.5
OY720033.1	11	19.85	37.0
OY720034.1	12	19.11	36.5
OY720035.1	13	18.03	37.0
OY720036.1	14	19.1	37.0
OY720037.1	15	18.85	37.5
OY720038.1	16	18.74	37.0

INSDC accession	Name	Length (Mb)	GC%
OY720039.1	17	18.71	37.0
OY720040.1	18	17.21	37.5
OY720041.1	19	17.03	37.5
OY720042.1	20	16.96	37.0
OY720043.1	21	16.78	37.5
OY720044.1	22	15.77	37.5
OY720045.1	23	14.63	37.5
OY720046.1	24	14.31	37.5
OY720047.1	25	13.77	37.5
OY720048.1	26	12.8	37.5
OY720049.1	27	11.77	38.0
OY720050.1	28	9.48	37.5
OY720051.1	29	10.47	37.5
OY720052.1	30	9.62	37.5
OY720053.1	31	0.01	37.0
OY720054.1	32	0.0	60.5
OY720022.1	Z	30.88	36.5
OY720055.1	MT	0.02	18.0



with ProNex beads and DNA concentration was quantified using the Qubit and fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies). Size selection was performed using Sage Sciences' PippinHT system with target fragment size determined by analysis from the Femto Pulse, usually a value between 4000 and 9000 bp. Size selected libraries were then cleaned-up using 1.0X ProNex beads and normalised to 2 nM before proceeding to sequencing.

Samples were 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, as well as perform primary and secondary analysis of the data upon completion.

### Hi-C

For Hi-C library preparation, DNA was fragmented using the Covaris E220 sonicator (Covaris) and size selected using SPRISelect beads to 400 to 600 bp. The DNA was then enriched using the Arima-HiC v2 kit Enrichment beads. Using the NEBNext Ultra II DNA Library Prep Kit (New England Biolabs) for end repair, a-tailing, and adapter ligation. This uses a custom protocol which resembles the standard NEBNext Ultra II DNA Library Prep protocol but where library preparation occurs while DNA is bound to the Enrichment beads. For library amplification, 10 to 16 PCR cycles were required, determined by the sample biotinylation percentage. The Hi-C sequencing was performed using paired-end sequencing with a read length of 150 bp on an Illumina NovaSeq 6000 instrument.

### RNA

Poly(A) RNA-Seq libraries were constructed using the NEB Ultra II RNA Library Prep kit, following the manufacturer's instructions. RNA sequencing was performed on the Illumina NovaSeq 6000 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 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 (article in preparation). Flat files and maps used in curation were generated in TreeVal (Pointon *et al.*, 2023). 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).

### Assembly quality assessment

The Merqury.FK tool (Rhie *et al.*, 2020), run in a Singularity container (Kurtzer *et al.*, 2017), was used to evaluate *k*-mer completeness and assembly quality for the primary and alternate haplotypes using the *k*-mer databases ( $k = 31$ ) that were computed prior to genome assembly. The analysis outputs included assembly QV scores and completeness statistics. The genome was also analysed within the BlobToolKit environment (Challis *et al.*, 2020) and BUSCO scores (Manni *et al.*, 2021) were calculated.

A Hi-C contact map was produced for the final version of the assembly. The Hi-C reads were aligned using bwa-mem2 (Vasimuddin *et al.*, 2019) and the alignment files were combined using SAMtools (Danecek *et al.*, 2021). The Hi-C alignments were converted into a contact map using BEDTools (Quinlan & Hall, 2010) and the Cooler tool suite (Abdennur & Mirny, 2020). The contact map was visualised in HiGlass (Kerpedjiev *et al.*, 2018).

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

### Genome annotation

The Ensembl Genebuild annotation system (Aken *et al.*, 2016) was used to generate annotation for the *Eudonia mercurella* assembly (GCA\_963082485.1) in Ensembl Rapid Release at the EBI. Annotation was created primarily through alignment of transcriptomic data to the genome, with gap filling via protein-to-genome alignments of a select set of proteins from UniProt (UniProt Consortium, 2019).

### 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](https://www.darwintreeoflife.org/). 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.

**Table 4. Software tools: versions and sources.**

Software tool	Version	Source
BEDTools	2.30.0	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>
BLAST	2.14.0	<a href="ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/">ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/</a>
BlobToolKit	4.3.7	<a href="https://github.com/blobtoolkit/blobtoolkit">https://github.com/blobtoolkit/blobtoolkit</a>
BUSCO	5.3.2	<a href="https://gitlab.com/ezlab/busco">https://gitlab.com/ezlab/busco</a>
bwa-mem2	2.2.1	<a href="https://github.com/bwa-mem2/bwa-mem2">https://github.com/bwa-mem2/bwa-mem2</a>
Cooler	0.8.11	<a href="https://github.com/open2c/cooler">https://github.com/open2c/cooler</a>
fasta_windows	0.2.4	<a href="https://github.com/tolkit/fasta_windows">https://github.com/tolkit/fasta_windows</a>
FastK	427104ea91c78c3b8b8b49f1a7d6bbeaa869ba1c	<a href="https://github.com/thegenemyers/FASTK">https://github.com/thegenemyers/FASTK</a>
Gfastats	1.3.6	<a href="https://github.com/vgl-hub/gfastats">https://github.com/vgl-hub/gfastats</a>
Hifiasm	0.16.1-r375	<a href="https://github.com/chhylp123/hifiasm">https://github.com/chhylp123/hifiasm</a>
HiGlass	44086069ee7d4d3f6f3f0012569789ec138f42b84aa44357826c0b6753eb28de	<a href="https://github.com/higlass/higlass">https://github.com/higlass/higlass</a>
Mercury.FK	d00d98157618f4e8d1a9190026b19b471055b22e	<a href="https://github.com/thegenemyers/MERQURY.FK">https://github.com/thegenemyers/MERQURY.FK</a>
MitoHiFi	3	<a href="https://github.com/marcelauliano/MitoHiFi">https://github.com/marcelauliano/MitoHiFi</a>
PretextView	0.2.5	<a href="https://github.com/sanger-tol/PretextView">https://github.com/sanger-tol/PretextView</a>
purge_dups	1.2.5	<a href="https://github.com/dfguan/purge_dups">https://github.com/dfguan/purge_dups</a>
samtools	1.16.1, 1.17, and 1.18	<a href="https://github.com/samtools/samtools">https://github.com/samtools/samtools</a>
sanger-tol/ascc	-	<a href="https://github.com/sanger-tol/ascc">https://github.com/sanger-tol/ascc</a>
Seqtk	1.3	<a href="https://github.com/lh3/seqtk">https://github.com/lh3/seqtk</a>
Singularity	3.9.0	<a href="https://github.com/sylabs/singularity">https://github.com/sylabs/singularity</a>
YaHS	1.2a.2	<a href="https://github.com/c-zhou/yahs">https://github.com/c-zhou/yahs</a>

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 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: *Eudonia mercurella* (small grey). Accession number PRJEB62417; <https://identifiers.org/ena.embl/PRJEB62417>. The genome sequence is released openly for reuse. The *Eudonia mercurella* 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 University of Oxford and Wytham Woods Genome Acquisition Lab are listed here: <https://doi.org/10.5281/zenodo.12157525>.

Members of the Natural History Museum Genome Acquisition Lab are listed here: <https://doi.org/10.5281/zenodo.12159242>.

Members of the Darwin Tree of Life Barcoding collective are listed here: <https://doi.org/10.5281/zenodo.12158331>.

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](#)
- Aken BL, Ayling S, Barrell D, *et al.*: **The ensembl gene annotation system.** *Database (Oxford).* 2016; **2016**: baw093.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Allio R, Schomaker-Bastos A, Romiguier J, *et al.*: **MitoFinder: efficient automated large-scale extraction of mitochondrial data in target enrichment phylogenomics.** *Mol Ecol Resour.* 2020; **20**(4): 892–905.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Beasley J, Uhl R, Forrest LL, *et al.*: **DNA barcoding SOPs for the Darwin Tree of Life project.** *protocols.io.* 2023; [Accessed 25 June 2024].  
[Publisher 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](#)
- 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): gjab008.  
[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](#)
- Diesh 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, Bates A, Denton A, *et al.*: **Sanger Tree of Life RNA extraction: automated MagMax™mirVana.** *protocols.io.* 2023.  
[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: ***Eudonia mercuriella* (Linnaeus, 1758).** Checklist dataset, *GBIF Backbone Taxonomy.* 2023; [Accessed 30 January 2025].  
[Reference Source](#)
- 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): gjaa153.  
[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](#)
- Oatley G, Sampaio F, Kitchin L, *et al.*: **Sanger Tree of Life HMW DNA fragmentation: Covaris g-TUBE for ULI PacBio.** 2023; [Accessed 13 June 2024].  
[Publisher Full Text](#)
- Pointon DL, Eagles W, Sims Y, *et al.*: **sanger-tol/treeval v1.0.0 – Ancient Atlantis.** 2023.  
[Publisher Full Text](#)
- Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic features.** *Bioinformatics.* 2010; **26**(6): 841–842.  
[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](#)
- Sheerin E, Sampaio F, Oatley G, *et al.*: **Sanger Tree of Life HMW DNA extraction: automated MagAttract v.1.** *protocols.io.* 2023.  
[Publisher Full Text](#)
- Sterling P, Parsons M, Lewington R: **Field guide to the micro-moths of Great Britain and Ireland, second edition.** London: Bloomsbury Publishing, 2023.  
[Reference Source](#)
- Strickland M, Cornwell C, Howard C: **Sanger Tree of Life fragmented DNA**

clean up: manual SPRI. *protocols.io*. 2023.

[Publisher Full Text](#)

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, 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](#)

UniProt Consortium: **UniProt: a worldwide hub of protein knowledge**. *Nucleic Acids Res.* 2019; **47**(D1): D506–D515.

[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](#)

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](#)