

DATA NOTE

The genome sequence of the Pale Oak Beauty, Hypomecis punctinalis (Scopoli, 1763)

[version 1; peer review: 2 approved]

Douglas Boyes¹⁺, Clare Boyes⁰²,

University of Oxford and Wytham Woods 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: 18 Sep 2024, 9:531

https://doi.org/10.12688/wellcomeopenres.23061.1

Latest published: 18 Sep 2024, 9:531

https://doi.org/10.12688/wellcomeopenres.23061.1

Abstract

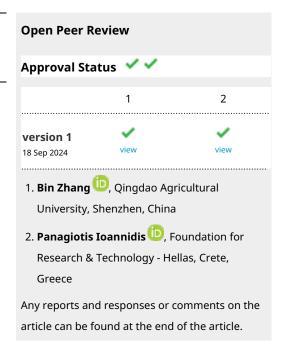
We present a genome assembly from an individual male Hypomecis punctinalis (the Pale Oak Beauty; Arthropoda; Insecta; Lepidoptera; Geometridae). The genome sequence has a total length of 741.20 megabases. Most of the assembly is scaffolded into 30 chromosomal pseudomolecules, including the Z sex chromosome. The mitochondrial genome has also been assembled and is 15.64 kilobases in length. Gene annotation of this assembly on Ensembl identified 13,897 protein-coding genes.

Keywords

Hypomecis punctinalis, Pale Oak Beauty moth, genome sequence, chromosomal, Lepidoptera



This article is included in the Tree of Life gateway.



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

²Independent researcher, Welshpool, Wales, UK

⁺ Deceased author

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

Author roles: Boyes D: 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 [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 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, Boyes C, University of Oxford and Wytham Woods Genome Acquisition Lab *et al.* The genome sequence of the Pale Oak Beauty, *Hypomecis punctinalis* (Scopoli, 1763) [version 1; peer review: 2 approved] Wellcome Open Research 2024, 9:531 https://doi.org/10.12688/wellcomeopenres.23061.1

First published: 18 Sep 2024, **9**:531 https://doi.org/10.12688/wellcomeopenres.23061.1

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; Ennominae; Hypomecis; Hypomecis; punctinalis (Scopoli, 1763) (NCBI:txid439567).

Background

Pale Oak Beauty (*Hypomecis punctinalis*) is a moth in the family Geometridae, which is widespread in south-eastern England (Randle *et al.*, 2019). It is found throughout Europe and has a disjointed population across Asia, as far east as Japan (GBIF Secretariat, 2024).

The adult moth has a forewing length of between 22–26 mm and is greyish brown with speckling and a dark, scalloped, central crossline. It has a dark central spot on the hindwing which helps to distinguish it from similar species. The adult is on the wing between May and mid-July and may have a partial second generation in southern England. It readily comes to light and during the day, adults can be found resting on tree-trunks. The main habitat is broadleaved woodland where its larval foodplants are common. These include pedunculate oak, downy and silver birch, hawthorn and sallows (Waring et al., 2017).

The genome of *Hypomecis punctinalis* 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 *Hypomecis punctinalis* based on one male specimen from Wytham Woods, Oxfordshire, UK.

Genome sequence report

The genome of an adult male *Hypomecis punctinalis* (Figure 1) was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating a total of 18.84 Gb (gigabases) from 1.66 million reads, providing approximately 25-fold coverage. Primary assembly contigs were scaffolded with chromosome conformation Hi-C data, which produced 143.19 Gb from 948.29 million reads, yielding an approximate coverage of 193-fold. Specimen and sequencing information is summarised in Table 1.

Manual assembly curation corrected 28 missing joins or mis-joins and 16 haplotypic duplications, reducing the assembly length by 1.15% and the scaffold number by 7.95%, and increasing the scaffold N50 by 1.81%. The final assembly has a total length of 741.20 Mb in 80 sequence scaffolds with a scaffold N50 of 25.9 Mb (Table 2). The snail plot in Figure 2 provides a summary of the assembly statistics, while the distribution of assembly scaffolds on GC proportion and coverage is shown in Figure 3. The cumulative assembly plot in Figure 4 shows curves for subsets of scaffolds assigned to different phyla. Most (99.43%) of the assembly sequence was assigned to 30



Figure 1. Photograph of the *Hypomecis punctinalis* (ilHypPunc1) specimen used for genome sequencing.

chromosomal-level scaffolds, representing 29 autosomes and the Z sex chromosome. Chromosome-scale scaffolds confirmed by the Hi-C data are named in order of size (Figure 5; Table 3). Chromosome Z was assigned by alignment to Lycia hirtaria (GCA_947563715.1) (Boyes et al., 2023). 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.

The estimated Quality Value (QV) of the final assembly is 63.9 with k-mer completeness of 100.0%, and the assembly has a BUSCO v5.3.2 completeness of 98.3% (single = 97.4%, duplicated = 1.0%), using the lepidoptera_odb10 reference set (n = 5,286).

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/439567.

Genome annotation report

The *Hypomecis punctinalis* genome assembly (GCA_949316475.1) was annotated at the European Bioinformatics Institute (EBI) on Ensembl Rapid Release. The resulting annotation includes 25,642 transcribed mRNAs from 13,897 protein-coding and 2,235 non-coding genes (Table 2; https://rapid.ensembl.org/Hypomecis_punctinalis_GCA_949316475.1/Info/Index). The average transcript length is 20,486.41. There are 1.59 coding transcripts per gene and 7.31 exons per transcript.

Methods

Sample acquisition

An adult male *Hypomecis punctinalis* (specimen ID Ox001903, ToLID ilHypPunc1) was collected from Wytham Woods,

Table 1. Specimen and sequencing data for Hypomecis punctinalis.

Project information			
Study title	Hypomecis punctinalis (pale oak beauty)		
Umbrella BioProject	PRJEB59306		
Species	Hypomecis punctinalis		
BioSample	SAMEA10979165		
NCBI taxonomy ID	439567		
Specimen information			
Technology	ToLID	BioSample accession	Organism part
PacBio long read sequencing	ilHypPunc1	SAMEA10979596	thorax
Hi-C sequencing	ilHypPunc1	SAMEA10979595	head
RNA sequencing	ilHypPunc1	SAMEA10979597	abdomen
Sequencing information			
Platform	Run accession	Read count	Base count (Gb)
Hi-C Illumina NovaSeq 6000	ERR10818325	9.48e+08	143.19
PacBio Sequel IIe	ERR10809409	1.66e+06	18.84
RNA Illumina NovaSeq 6000	ERR11641124	8.09e+07	12.22

Oxfordshire (biological vice-county Berkshire), UK (latitude 51.77, longitude –1.34) on 2021-06-16, using a light trap. The specimen was collected and identified by Douglas Boyes (University of Oxford) and preserved on dry ice.

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 of the specimen 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 y available on protocols.io (Denton et al., 2023). In sample preparation, the ilHypPunc1 sample was weighed and dissected on dry ice (Jay et al., 2023). For sample homogenisation, thorax tissue

was cryogenically disrupted using the Covaris cryoPREP® Automated Dry Pulverizer (Narváez-Gómez et al., 2023).

HMW DNA was extracted using the Automated MagAttract v1 protocol (Sheerin et al., 2023). DNA was sheared into an average fragment size of 12–20 kb in a Megaruptor 3 system (Todorovic 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 ilHypPunc1 in the Tree of Life Laboratory at the WSI using the RNA Extraction: Automated MagMaxTM *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.

Library preparation and sequencing

Pacific Biosciences HiFi circular consensus DNA sequencing libraries were constructed according to the manufacturers' instructions. Poly(A) RNA-Seq libraries were constructed using the NEB Ultra II RNA Library Prep kit. DNA and RNA

Table 2. Genome assembly data for *Hypomecis punctinalis*, ilHypPunc1.1.

Genome assembly			
Assembly name	ilHypPunc1.1		
Assembly accession	GCA_949316475.1		
Accession of alternate haplotype	GCA_949316335.1		
Span (Mb)	741.20		
Number of contigs	391		
Contig N50 length (Mb)	3.6		
Number of scaffolds	80		
Scaffold N50 length (Mb)	25.9		
Longest scaffold (Mb)	33.65		
Assembly metrics*		Benchmark	
Consensus quality (QV)	63.9	≥ 50	
k-mer completeness	100.0%	≥ 95%	
BUSCO**	C:98.3%[S:97.4%,D:1.0%], F:0.5%,M:1.2%,n:5,286	<i>C</i> ≥ 95%	
Percentage of assembly mapped to chromosomes	99.43%	≥ 95%	
Sex chromosomes	Z	localised homologous pairs	
Organelles	Mitochondrial genome: 15.64 kb	complete single alleles	
Genome annotation of assembly GCA_949316475.1 at Ensembl			
Number of protein-coding genes	13,897		
Number of non-coding genes	2,235		
Number of gene transcripts	25,642		

^{*} 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).

sequencing was performed by the Scientific Operations core at the WSI on Pacific Biosciences Sequel IIe (HiFi) and Illumina NovaSeq 6000 (RNA-Seq) instruments.

Hi-C data were generated from frozen head tissue of the ilHypPunc1 sample, using the Arima-HiC v2 kit. In brief, frozen tissue (-80 °C) was fixed, and the DNA crosslinked using a TC buffer containing formaldehyde. The crosslinked DNA was then digested using a restriction enzyme master mix. The 5'-overhangs were then filled in and labelled with a biotinylated nucleotide and proximally ligated. The biotinylated DNA construct was fragmented to a fragment 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.

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

^{**} BUSCO scores based on the lepidoptera_odb10 BUSCO set using version 5.3.2. $C = complete [S = single copy, D = duplicated], F = fragmented, M = missing, n = number of orthologues in comparison. A full set of BUSCO scores is available at https://blobtoolkit.genomehubs.org/view/ilHypPunc1_1/dataset/ilHypPunc1_1/busco.$

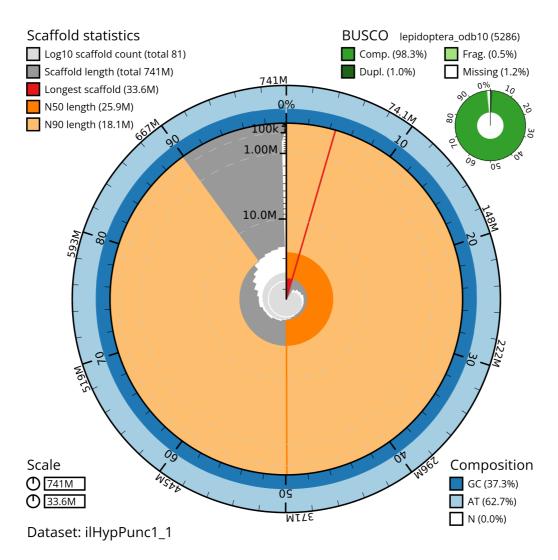


Figure 2. Genome assembly of *Hypomecis punctinalis*, **ilHypPunc1.1: metrics.** The BlobToolKit snail plot shows N50 metrics and BUSCO gene completeness. The main plot is divided into 1,000 size-ordered bins around the circumference with each bin representing 0.1% of the 741,177,338 bp assembly. The distribution of scaffold lengths is shown in dark grey with the plot radius scaled to the longest scaffold present in the assembly (33,645,120 bp, shown in red). Orange and pale-orange arcs show the N50 and N90 scaffold lengths (25,902,484 and 18,132,001 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/ilHypPunc1_1/snail.

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.

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

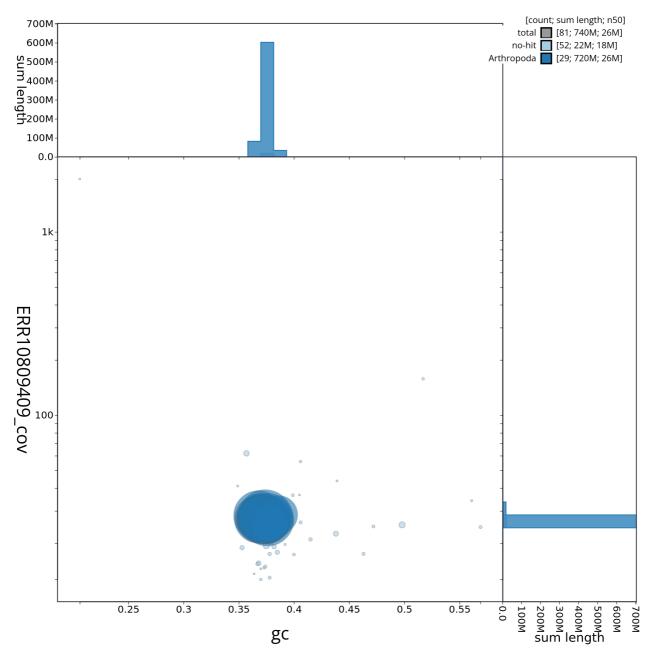


Figure 3. Genome assembly of *Hypomecis punctinalis*, ilHypPunc1.1: Blob plot of base coverage against GC proportion for sequences in the assembly. Sequences are coloured by phylum. Circles are sized in proportion to sequence length. Histograms show the distribution of sequence length sum along each axis. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/ilHypPunc1_1/dataset/ilHypPunc1_1/blob.

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 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

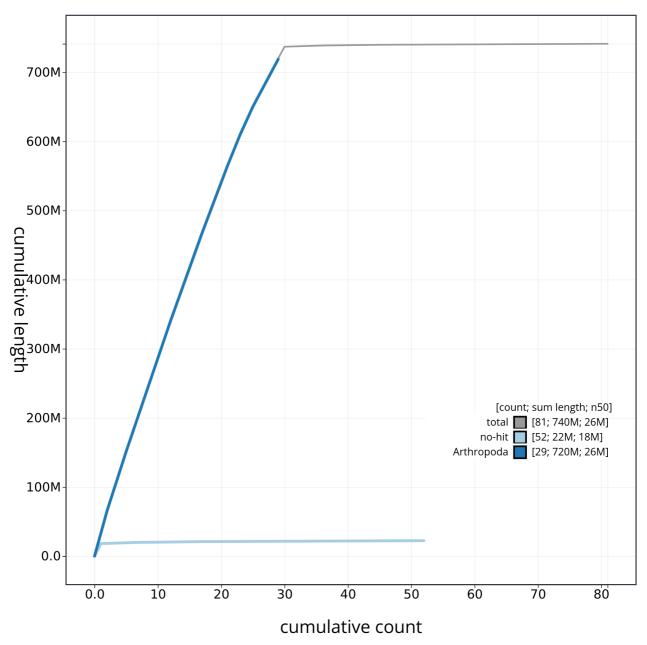


Figure 4. Genome assembly of *Hypomecis punctinalis* **ilHypPunc1.1: BlobToolKit cumulative sequence plot.** The grey line shows cumulative length for all sequences. Coloured lines show cumulative lengths of sequences assigned to each phylum using the buscogenes taxrule. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/ilHypPunc1_1/dataset/ilHypPunc1_1/cumulative.

(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 Ensembl Genebuild annotation system (Aken *et al.*, 2016) was used to generate annotation for the *Hypomecis punctinalis* assembly (GCA_949316475.1) 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. By agreeing with and signing up to the Sampling

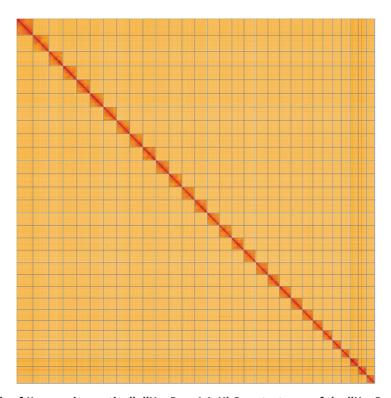


Figure 5. Genome assembly of *Hypomecis punctinalis* ilHypPunc1.1: Hi-C contact map of the ilHypPunc1.1 assembly, visualised using HiGlass. Chromosomes are shown in order of size from left to right and top to bottom. An interactive version of this figure may be viewed at https://genome-note-higlass.tol.sanger.ac.uk/l/?d=HKFWaArvTqSZ6sFBNt7h1A.

Table 3. Chromosomal pseudomolecules in the genome assembly of *Hypomecis punctinalis*, ilHypPunc1.

INSDC accession	Name	Length (Mb)	GC%
OX438795.1	1	32.57	37.5
OX438796.1	2	29.12	37.0
OX438797.1	3	28.27	37.5
OX438798.1	4	27.81	37.0
OX438799.1	5	27.23	37.5
OX438800.1	6	27.17	37.0
OX438801.1	7	27.03	37.5
OX438802.1	8	26.99	37.0
OX438803.1	9	26.93	37.5
OX438804.1	10	26.9	37.0
OX438805.1	11	26.6	37.0
OX438806.1	12	26.18	37.0
OX438807.1	13	25.9	37.0
OX438808.1	14	25.44	37.0

INSDC accession	Name	Length (Mb)	GC%
OX438809.1	15	25.36	37.5
OX438810.1	16	24.81	37.0
OX438811.1	17	24.76	37.0
OX438812.1	18	24.53	37.0
OX438813.1	19	23.82	37.5
OX438814.1	20	23.78	37.0
OX438815.1	21	23.22	37.5
OX438816.1	22	21.91	37.5
OX438817.1	23	20.59	37.5
OX438818.1	24	19.68	37.5
OX438819.1	25	18.13	37.5
OX438820.1	26	17.78	37.5
OX438821.1	27	17.28	38.0
OX438822.1	28	17.01	38.5
OX438823.1	29	16.5	38.0
OX438794.1	Z	33.65	37.5
OX438824.1	MT	0.02	20.5

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	0.16.1-r375	https://github.com/chhylp123/hifiasm
HiGlass	1.11.6	https://github.com/higlass/higlass
Merqury.FK	d00d98157618f4e8d1a9 190026b19b471055b22e	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/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	https://github.com/c-zhou/yahs

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 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: Hypomecis punctinalis (pale oak beauty). Accession number PRJEB59306; https://identifiers.org/ena.embl/PRJEB59306 (Wellcome Sanger Institute, 2023).

The genome sequence is released openly for reuse. The *Hypomecis punctinalis* 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 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 mitogenomic 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. Publisher Full Text

Boyes D, Holland PWH, University of Oxford and Wytham Woods Genome Acquisition Lab, et al.: The genome sequence of the Brindled Beauty, Lycia hirtaria (Clerck, 1759) [version 1; peer review: 3 approved]. Wellcome Open Res. 2023; 8: 303.

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

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

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, Yatsenko H, Jay J, et al.: Sanger Tree of Life wet laboratory protocol collection V.1. protocols.io. 2023.

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 $^{
m M}$ 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.: 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: *Hypomecis punctinalis* (Scopoli, 1763). *GBIF Backbone Taxonomy*. 2024; [Accessed 14 August 2024].

Reference Source

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): giaa153. 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, 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, [Accessed 2 April 2024]. **Reference Source**

Narváez-Gómez JP, Mbye H, Oatley G, et al.: Sanger Tree of Life sample homogenisation: covaris cryoPREP® automated dry pulverizer V.1. protocols. in 2023

Publisher Full Text

Randle Z, Evans-Hill LJ, Parsons MS, et al.: Atlas of Britain & Ireland's larger moths. Newbury: NatureBureau, 2019.

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):

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**

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

Todorovic M, Sampaio F, Howard C: Sanger Tree of Life HMW DNA fragmentation: diagenode Megaruptor®3 for PacBio HiFi. protocols.io. 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: awaiting peer review]. Wellcome Open Res. 2024; 9: 339. **Publisher 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

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 Md, 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 Pale Oak Beauty**, **Hypomecis punctinalis** (**Scopoli**, **1763**). European Nucleotide Archive. [dataset], accession number PRJEB59306, 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

Open Peer Review

Current Peer Review Status:





Version 1

Reviewer Report 06 October 2025

https://doi.org/10.21956/wellcomeopenres.25397.r132281

© 2025 Ioannidis P. This is an open access peer review report 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.



Panagiotis Ioannidis 🗓



Foundation for Research & Technology - Hellas, Crete, Greece

This manuscript describes the genome sequencing, assembly and annotation of the lepidopteran Hypomecis punctinalis.

All methodology is the standard one used in DToL projects, and the resulting assembly looks good with virtually the entire assembly being found in a few chromosome-level scaffolds.

My only comment has to do with reporting the genome annotation results. The community needs to understand that simply reporting the number of mRNAs, protein-coding genes etc is not enough! They should be also reporting the BUSCO score of this gene set, for example. Or some other metric of quality.

Also, what is a "disjointed population"? Maybe the authors could explain more? Was it transferred there by transportation from Europe? Is it considered invasive? Is it causing problems in Asia?

Is the rationale for creating the dataset(s) clearly described?

Yes

Are the protocols appropriate and is the work technically sound?

Are sufficient details of methods and materials provided to allow replication by others?

Are the datasets clearly presented in a useable and accessible format?

Yes

Competing Interests: No competing interests were disclosed.

Reviewer Expertise: insect genomics; bioinformatics

I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.

Reviewer Report 03 November 2024

https://doi.org/10.21956/wellcomeopenres.25397.r104954

© **2024 Zhang B.** This is an open access peer review report 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.



Bin Zhang 🗓

China-Australia Joint Institute of Agricultural and Environmental Health, Qingdao Agricultural University, Shenzhen, China

Pale Oak Beauty is an important moth, which is widespread in south-eastern England. The sample for genome assembly is suitable to conduct the sequencing using Pacific Biosciences single-molecule HiFi long reads. According to the results, the quality of assembly is good and amount of annotation is in noraml range. I believe this genome assembly would be a excellent data basic for further research on this moth.

Is the rationale for creating the dataset(s) clearly described?

Yes

Are the protocols appropriate and is the work technically sound?

Yes

Are sufficient details of methods and materials provided to allow replication by others? Yes

Are the datasets clearly presented in a useable and accessible format?

Yes

Competing Interests: No competing interests were disclosed.

Reviewer Expertise: insect ecology and bioinformatics

I confirm that I have read this submission and believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.