

**DATA NOTE** 

# The genome sequence of the Shy Cosmet Moth, Limnaecia phragmitella Stainton, 1851 (Lepidoptera: Cosmopterigidae)

[version 1; peer review: 2 approved]

Douglas Boyes<sup>1+</sup>, Clare Boyes<sup>2</sup>,

University of Oxford and Wytham Woods 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: 15 Sep 2025, **10**:509

https://doi.org/10.12688/wellcomeopenres.24881.1

Latest published: 15 Sep 2025, 10:509

https://doi.org/10.12688/wellcomeopenres.24881.1

#### **Abstract**

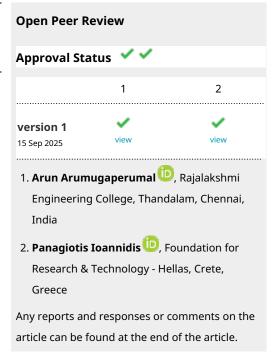
We present a genome assembly from an individual male *Limnaecia* phragmitella (Shy Cosmet Moth; Arthropoda; Insecta; Lepidoptera; Cosmopterigidae). The assembly contains two haplotypes with total lengths of 539.22 megabases and 553.78 megabases. Most of haplotype 1 (98.9%) is scaffolded into 30 chromosomal pseudomolecules, including the Z sex chromosome. Haplotype 2 was assembled to scaffold level. The mitochondrial genome has also been assembled, with a length of 15.26 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**

Limnaecia phragmitella; Shy Cosmet Moth; genome sequence; chromosomal; Lepidoptera



This article is included in the Tree of Life gateway.



<sup>&</sup>lt;sup>1</sup>UK Centre for Ecology & Hydrology, Wallingford, England, UK <sup>2</sup>Independent researcher, Welshpool, Wales, UK

<sup>&</sup>lt;sup>+</sup> 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 (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, Boyes C, University of Oxford and Wytham Woods Acquisition Lab *et al.* The genome sequence of the Shy Cosmet Moth, *Limnaecia phragmitella* Stainton, 1851 (Lepidoptera: Cosmopterigidae) [version 1; peer review: 2 approved] Wellcome Open Research 2025, 10:509 https://doi.org/10.12688/wellcomeopenres.24881.1

First published: 15 Sep 2025, 10:509 https://doi.org/10.12688/wellcomeopenres.24881.1

#### 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; Gelechioidea; Cosmopterigidae; Cosmopteriginae; Limnaecia; Limnaecia phragmitella Stainton, 1851 (NCBI:txid687102)

## **Background**

Limnaecia phragmitella is a micro-moth in the family Cosmopterigidae. It is a moth of wet areas in which its foodplants, Bulrush (*Typha latifolia*) and Lesser Bulrush (*Typha angustifolia*), grow. It is common in suitable habitats throughout the UK. It is occurs throughout Europe and North America, with scattered records from New Zealand, Russia and Japan (GBIF Secretariat, 2025).

The moth is small (forewing length 8–11 mm) with a pale brown forewing with indistinct markings (Sterling *et al.*, 2023). It is more easily found in its larval stage as a miner in the seed heads and stems of its foodplants. The early instar larvae mines in the leaf sheaf, later moving to the seedheads. The presence of the larvae can be recognised by the down hanging from the seed head. The larval silk prevents the down from being dispersed by the wind (Emmet & Langmaid, 2002).

We present a chromosome-level genome sequence for Limnaecia phragmitella, the Shy Cosmet Moth. This assembly is the first high-quality genome for the genus Limnaecia and one of two genomes available for the family Cosmopterigidae 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 *Limnaecia phragmitella* (ilLimPhra2) specimen used for genome sequencing.

#### Methods

#### Sample acquisition and DNA barcoding

The specimen used for genome sequencing was an adult male *Limnaecia phragmitella* (specimen ID Ox001870, ToLID ilLimPhra2; Figure 1), collected from Wytham Woods, Oxfordshire, United Kingdom (latitude 51.765, longitude –1.335) on 2021-07-20. A second specimen collected on the same occasionwas used for Hi-C sequencing (specimen ID Ox001869, ToLID ilLimPhra1). The specimens were collected and identified by Douglas Boyes. 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 ilLimPhra2 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. HMW DNA was extracted using the Automated MagAttract v2 protocol. We used centrifuge-mediated fragmentation to produce DNA fragments in the 8-10 kb range, following the Covaris g-TUBE protocol for ultra-low input (ULI). Sheared DNA was purified by automated SPRI (solid-phase reversible immobilisation). 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. For this sample, the final post-shearing DNA had a Qubit concentration of 1.08 ng/µL and a yield of 140.40 ng.

#### PacBio HiFi library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core. Prior to library preparation, the DNA was fragmented to ~10 kb. Ultra-low-input (ULI) libraries were prepared using the PacBio SMRTbell® Express Template Prep Kit 2.0 and gDNA Sample Amplification Kit. Samples were normalised to 20 ng DNA. Single-strand overhang removal, DNA damage repair, and end-repair/A-tailing

were performed according to the manufacturer's instructions, followed by adapter ligation. A 0.85× pre-PCR clean-up was carried out with Promega ProNex beads.

The DNA was evenly divided into two aliquots for dual PCR (reactions A and B), both following the manufacturer's protocol. A  $0.85\times$  post-PCR clean-up was performed with ProNex beads. DNA concentration was measured using a Qubit Fluorometer v4.0 (Thermo Fisher Scientific) with the Qubit HS Assay Kit, and fragment size was assessed on an Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) using the gDNA 55 kb BAC analysis kit. PCR reactions A and B were then pooled, ensuring a total mass of  $\geq 500$  ng in 47.4 µl.

The pooled sample underwent another round of DNA damage repair, end-repair/A-tailing, and hairpin adapter ligation. A 1× clean-up was performed with ProNex beads, followed by DNA quantification using the Qubit and fragment size analysis using the Agilent Femto Pulse. Size selection was performed on the Sage Sciences PippinHT system, with target fragment size determined by Femto Pulse analysis (typically 4–9 kb). Size-selected libraries were cleaned with 1.0× ProNex beads and normalised to 2 nM before sequencing.

The sample was sequenced on a Revio instrument (Pacific Biosciences). The prepared library was normalised to 2 nM, and  $15\,\mu\text{L}$  was used for making complexes. Primers were annealed and polymerases bound to generate circularised complexes, following the manufacturer's instructions. Complexes were purified using 1.2X SMRTbell beads, then diluted to the Revio loading concentration (200–300 pM) and spiked with a Revio sequencing internal control. The sample was sequenced on a Revio 25M SMRT cell. The SMRT Link software (Pacific Biosciences), a web-based workflow manager, was used to configure and monitor the run and to carry out primary and secondary data analysis.

#### Hi-C

#### Sample preparation and crosslinking

The Hi-C sample was prepared from 20–50 mg of frozen tissue of the ilLimPhra1 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 6000.

#### 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 in Hi-C phasing mode (Cheng et al., 2021; Cheng et al., 2022), producing two haplotypes. Hi-C reads (Rao et al., 2014) were mapped to the primary contigs using bwa-mem2 (Vasimuddin et al., 2019). Contigs were further scaffolded with Hi-C data 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 MEROURY.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). Manual corrections included 99 breaks and 177 joins. The curation process is documented at <a href="https://gitlab.com/wtsi-grit/rapid-curation">https://gitlab.com/wtsi-grit/rapid-curation</a>. PretextSnapshot 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 both 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 *Limnaecia phragmitella* specimen generated 46.22 Gb (gigabases) from 5.65 million reads, which were used to assemble the genome. GenomeScope2.0 analysis estimated the haploid genome size at 486.08 Mb, with a heterozygosity of 0.66% and repeat content of 24.49% (Figure 2). These estimates guided expectations for the assembly. Based on the estimated genome size, the sequencing data provided approximately 87× coverage. Hi-C sequencing produced 117.68 Gb from 779.36 million reads, which were used to scaffold the assembly. Table 1 summarises the specimen and sequencing details.

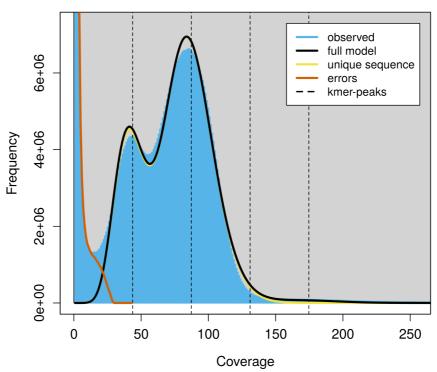
## Assembly statistics

The genome was assembled into two haplotypes using Hi-C phasing. Haplotype 1 was curated to chromosome level, while haplotype 2 was assembled to scaffold level. The final assembly has a total length of 539.22 Mb in 170 scaffolds, with 1 496 gaps, and a scaffold N50 of 18.83 Mb (Table 2).

Most of the assembly sequence (98.9%) was assigned to 30 chromosomal-level scaffolds, representing 29 autosomes

## GenomeScope Profile

len:486,082,089bp uniq:75.7% aa:99.3% ab:0.656% kcov:43.7 err:0.0869% dup:2.55 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 PRJEB85382.

Platform	PacBio HiFi	Hi-C
ToLID	ilLimPhra2	ilLimPhra1
Specimen ID	Ox001870	Ox001869
BioSample (source individual)	SAMEA10979130	SAMEA10979129
BioSample (tissue)	SAMEA10979542	SAMEA10979541
Tissue	whole organism	whole organism
Instrument	Revio	Illumina NovaSeq 6000
Run accessions	ERR14231588	ERR14242295
Read count total	5.65 million	779.36 million
Base count total	46.22 Gb	117.68 Gb

Table 2. Genome assembly statistics.

Assembly name	ilLimPhra2.hap1.1	ilLimPhra2.hap2.1
Assembly accession	GCA_965195335.1	GCA_965195365.1
Assembly level	chromosome	scaffold
Span (Mb)	539.22	553.78
Number of chromosomes	30	N/A
Number of contigs	1 666	11 249
Contig N50	0.56 Mb	0.09 Mb
Number of scaffolds	170	4 445
Scaffold N50	18.83 Mb	14.83 Mb
Longest scaffold length (Mb)	28.59	N/A
Sex chromosomes	Z	N/A
Organelles	Mitochondrion: 15.26 kb	N/A

and the Z sex chromosome. These chromosome-level scaffolds, confirmed by Hi-C data, are named according to size (Figure 3; Table 3). The Z chromosome was identified based on BUSCO gene painting with ancestral Merian elements (Wright *et al.*, 2024).

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.

For haplotype 1, the estimated QV is 55.5, and for haplotype 2, 57.1. When the two haplotypes are combined, the assembly achieves an estimated QV of 56.2. The k-mer completeness is 91.42% for haplotype 1, 82.76% for haplotype 2, and 99.36% for the combined haplotypes (Figure 4).

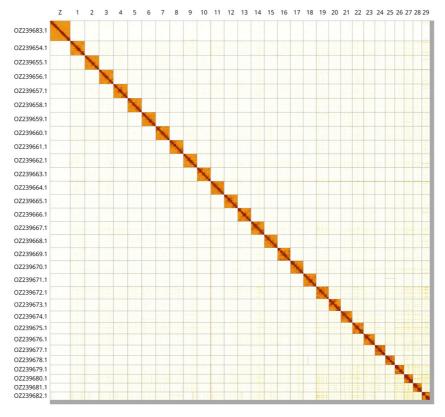
BUSCO analysis using the lepidoptera\_odb10 reference set (n = 5 286) identified 98.5% of the expected gene set

(single = 97.9%, duplicated = 0.6%) for haplotype 1. The snail plot in Figure 5 summarises the scaffold length distribution and other assembly statistics for haplotype 1. The blob plot in Figure 6 shows the distribution of scaffolds by GC proportion and coverage for haplotype 1.

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 haplotype 1, is **5.C.Q55**.

## 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

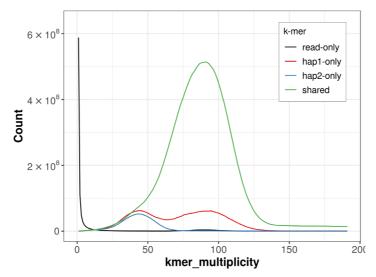


**Figure 3. Hi-C contact map of the** *Limnaecia phragmitella* **genome assembly.** Assembled chromosomes are shown in order of size and labelled along the axes. The plot was generated using PretextSnapshot.

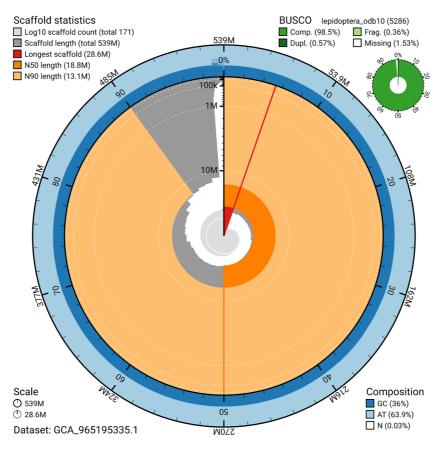
Table 3. Chromosomal pseudomolecules in the haplotype 1 genome assembly of *Limnaecia phragmitella* ilLimPhra2.

INSDC accession	Molecule	Length (Mb)	GC%
OZ239654.1	1	20.39	36
OZ239655.1	2	20.10	35.50
OZ239656.1	3	20.09	36
OZ239657.1	4	19.94	35.50
OZ239658.1	5	19.77	35.50
OZ239659.1	6	19.69	35.50
OZ239660.1	7	19.37	36
OZ239661.1	8	19.28	35.50
OZ239662.1	9	19.27	36.50
OZ239663.1	10	19.12	35.50
OZ239664.1	11	18.97	36.50
OZ239665.1	12	18.89	35.50
OZ239666.1	13	18.83	36
OZ239667.1	14	18.60	36

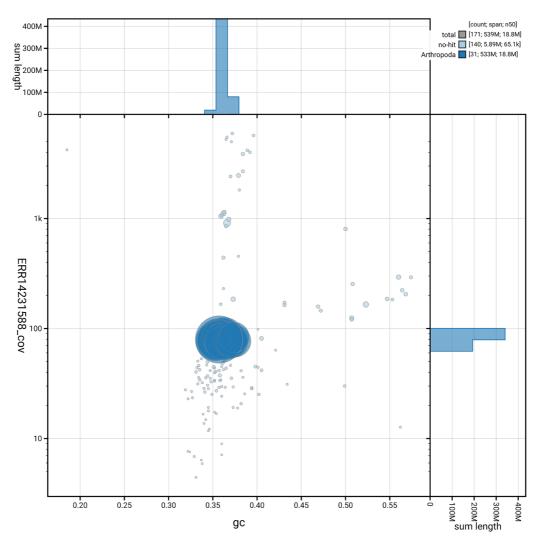
INSDC accession	Molecule	Length (Mb)	GC%
OZ239668.1	15	18.42	35.50
OZ239669.1	16	18.24	36
OZ239670.1	17	18.08	36
OZ239671.1	18	18.04	35.50
OZ239672.1	19	17.64	36
OZ239673.1	20	16.75	36
OZ239674.1	21	16.31	36
OZ239675.1	22	15.96	37
OZ239676.1	23	15.69	36
OZ239677.1	24	14.57	37
OZ239678.1	25	13.33	36
OZ239679.1	26	13.06	37
OZ239680.1	27	12.55	37.50
OZ239681.1	28	12.35	37
OZ239682.1	29	11.41	37.50
OZ239683.1	Z	28.59	35.50



**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 ilLimPhra2.hap1.1.** 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 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 ilLimPhra2.hap1.1.** 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.

 $\label{thm:continuous} \textbf{Table 4.} \ \textbf{Earth Biogenome Project summary metrics for the } \textit{Limnaecia phragmitella} \\ \textbf{assembly.}$ 

Measure	Value	Benchmark
EBP summary (haplotype 1)	5.C.Q55	6.C.Q40
Contig N50 length	0.56 Mb	≥ 1 Mb
Scaffold N50 length	18.83 Mb	= chromosome N50
Consensus quality (QV)	Haplotype 1: 55.5; haplotype 2: 57.1; combined: 56.2	≥ 40
k-mer completeness	Haplotype 1: 91.42%; Haplotype 2: 82.76%; combined: 99.36%	≥ 95%
BUSCO	C:98.5% [S:97.9%; D:0.6%]; F:0.4%; M:1.2%; n:5 286	S > 90%; D < 5%
Percentage of assembly assigned to chromosomes	98.90%	≥ 90%

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: Limnaecia phragmitella (bulrush cosmet). Accession number PRJEB85382. The genome sequence is released openly for reuse. The *Limnaecia phragmitella* 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 <a href="https://github.com/sanger-tol">https://github.com/sanger-tol</a>. Table 5 lists software versions used in this study.

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	4.4.5	https://github.com/blobtoolkit/blobtoolkit
BUSCO	5.7.1	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	2.28-r1209	https://github.com/lh3/minimap2
MitoHiFi	3	https://github.com/marcelauliano/MitoHiFi
MultiQC	1.14; 1.17 and 1.18	https://github.com/MultiQC/MultiQC
Nextflow	24.10.4	https://github.com/nextflow-io/nextflow
PretextSnapshot	N/A	https://github.com/sanger-tol/PretextSnapshot

Software	Version	Source
PretextView	0.2.5	https://github.com/sanger-tol/PretextView
samtools	1.21	https://github.com/samtools/samtools
sanger-tol/ascc	0.1.0	https://github.com/sanger-tol/ascc
sanger-tol/blobtoolkit	v0.7.1	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	1.2.2	https://github.com/c-zhou/yahs

#### Author information

Contributors are listed at the following links:

- Members of the University of Oxford and Wytham Woods 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

#### 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): 170.175.

PubMed Abstract | Publisher Full Text | Free Full Text

Cheng H, Jarvis ED, Fedrigo O, *et al.*: **Haplotype-resolved assembly of diploid genomes without parental data**. *Nat Biotechnol*. 2022; **40**(9): 1332–35. **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): giab008.

PubMed Abstract | Publisher Full Text | Free Full Text

Emmet AM, Langmaid JR: **The moths and butterflies of Great Britain and Ireland. Oecophoridae – scythrididae (Excluding Gelechiidae).** Colchester: Harley Books, 2002.

Reference Source

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: *Limnaecia phragmitella* Stainton, 1851 on the GBIF Backbone Taxonomy. 2025.

**Reference Source** 

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

Kerpedijev P. Abdennur N. Lekschas F. et al.: HiGlass: web-based visual exploration and analysis of genome interaction maps, Genome Biol. 2018:

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):

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

## **Open Peer Review**

## **Current Peer Review Status:**





## **Version 1**

Reviewer Report 24 October 2025

https://doi.org/10.21956/wellcomeopenres.27406.r135613

© 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 sequencing and assembly of Limnaecia phragmitella, a lepidopteran insect.

The methodology is identical to the one used by other DToL projects. All assembly metrics (N50, completeness, QV, BUSCO, blobplot) show that the assembly quality is very good.

Recently, I have come to realize that providing a genome assembly without an accompanying gene set is of limited to most scientists; very few are either interest in the genome sequence per se, or know how to perform gene prediction.

Related to the above, the authors mention that "The genome will be annotated using available RNA-Seg data...". What are these RNAseg data? Are they from the same organism? Because if they're not I would guess that the resulting gene set would be of relatively low quality.

## 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 24 October 2025

https://doi.org/10.21956/wellcomeopenres.27406.r135610

© 2025 Arumugaperumal A. 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.



## Arun Arumugaperumal 🗓



Department of Biotechnology, Rajalakshmi Engineering College, Thandalam, Chennai, Tamil NaduIndia, 602105, India

The data note describes the genome sequencing of the Shy Cosmet Moth, *Limnaecia phragmitella*. The authors have used long-read sequencing and Hi-C mapping to obtain a high-quality genome sequence. The assembly reported here is of size 539.22 Mb spread among 30 chromosomes. Since a male specimen has been used for sequencing, the authors have captured only the Z sex chromosome. A female specimen could have been used. The mitogenome has also been presented of size 15.26 kb. The annotations were yet to be released in Ensembl. The genome sequence achieved a EBP summary metric of 5.C.Q55. However, the contig N50 value is less compared to the benchmark. The BUSCO completeness percentage of 98.5 indicates that the genome is near-complete with respect to sequences in the database.

In the introduction part, check the phrase 'The early instar larvae mines in the leaf sheaf..'. Is it sheath or sheaf?

The article can be indexed.

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:** Bioinformatics; Genomics

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.