DATA NOTE

# The genome sequence of the Knot Grass moth, *Acronicta rumicis* Linnaeus, 1758 (Lepidoptera: Noctuidae)

[version 1; peer review: awaiting peer review]

Ian Sims[1], Douglas Boyes[2+], Natural History Museum Genome Acquisition Lab,
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

[1]Syngenta International Research Station, Jealott's Hill, Berkshire, England, UK
[2]UK Centre for Ecology & Hydrology, Wallingford, Oxfordshire, England, UK

[+] Deceased author

**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 an individual female *Acronicta rumicis* (Knot Grass moth; Arthropoda; Insecta; Lepidoptera; Noctuidae). The assembly contains two haplotypes with total lengths of 582.86 megabases and 528.05 megabases. Most of haplotype 1 (99.6%) is scaffolded into 32 chromosomal pseudomolecules, including the W and Z sex chromosomes. Haplotype 2 was assembled to scaffold level. The mitochondrial genome has also been assembled, with a length of 15.39 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

Acronicta rumicis; Knot Grass moth; genome sequence; chromosomal; Lepidoptera

This article is included in the Tree of Life gateway.

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

**Author roles: Sims I**: Investigation, Resources; **Boyes D**: Investigation, Resources;

**How to cite this article:** Sims I, Boyes D, Natural History Museum Genome Acquisition Lab *et al.* **The genome sequence of the Knot Grass moth,** *Acronicta rumicis* **Linnaeus, 1758 (Lepidoptera: Noctuidae) [version 1; peer review: awaiting peer review]** Wellcome Open Research 2025, **10**:484 https://doi.org/10.12688/wellcomeopenres.24843.1

**First published:** 08 Sep 2025, **10**:484 https://doi.org/10.12688/wellcomeopenres.24843.1

## 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; Noctuoidea; Noctuidae; Acronictinae; *Acronicta*; *Acronicta rumicis* Linnaeus, 1758 (NCBI:txid753146)

## Background

The Knot Grass moth *Acronicta rumicis* (Linnaeus, 1758) is a widespread noctuid moth occurring throughout the UK in a range of open habitats, including grasslands, heathlands, wetlands, and gardens (The Wildlife Trusts, 2025; Waring *et al.*, 2017). Adults are typically on the wing from May to July, with a second brood flying in August and September in the southern parts of England (Kimber, 2025; The Wildlife Trusts, 2025; Waring *et al.*, 2017). Adults are nocturnal, frequently visiting flowers and are attracted to light.

The forewings are usually mottled grey with a small curved white mark near the rear edge, a helpful feature for distinguishing this species even in melanic forms (Kimber, 2025; The Wildlife Trusts, 2025). The caterpillars are black with red markings, long brown hairs, and a distinctive resting posture, hunched just behind the head. Larvae feed on a wide variety of woody and herbaceous plants, including knotgrass, dock, plantains, bramble, hawthorn, sorrel, heather, and purple loosestrife. They are present from June to September, with later broods extending into October in the south. The species overwinters as a pupa in a cocoon amongst ground litter (The Wildlife Trusts, 2025).

Despite its widespread distribution, *A. rumicis* has undergone a major decline. It is now considered **Vulnerable** under IUCN criteria following an estimated 68% decline over 25 years, based on long-term population monitoring (Butterfly Conservation & Rothamsted Research, 2006; Fox *et al.*, 2006). It is included on the UK Biodiversity Action Plan priority list and has also been designated as a Northern Ireland Priority Species (Joint Nature Conservation Committee (JNCC), 2007; National Museums Northern Ireland, 2025).

We present the first genome sequence for *A. rumicis*. The assembly was produced using the Tree of Life pipeline from a specimen collected in Hartslock Nature Reserve, England, United Kingdom (Figure 1), as part of the Darwin Tree of Life project.

## Methods

### Sample acquisition and DNA barcoding

The specimen used for genome sequencing was an adult female *Acronicta rumicis* (specimen ID NHMUK013805967, ToLID ilAcrRumi2; Figure 1), collected from Hartslock Nature Reserve, England, United Kingdom (latitude 51.51, longitude –1.11) on 2021-07-29. The specimen was collected by Ian Sims and formally identified by Ian Sims and David Lees. A second specimen was used for Hi-C sequencing (specimen



**Figure 1. Photograph of the *Acronicta rumicis* (ilAcrRumi2) specimen used for genome sequencing.**

ID Ox000692, ToLID ilAcrRumi1). It was collected from Wytham Woods, Oxfordshire, United Kingdom (latitude 51.772, longitude –1.338) on 2020-07-20. This specimen was collected and identified by Douglas Boyes. All sample metadata were collected in line with the Darwin Tree of Life project standards described by 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 ilAcrRumi2 sample was weighed and triaged to determine the appropriate extraction protocol. Tissue from the thorax was homogenised by powermashing using a PowerMasher II tissue disruptor. HMW DNA was extracted using the MagAttract v3 protocol. DNA was sheared into an average fragment size of 12–20 kb following the Megaruptor®3 for LI PacBio protocol. 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 6.8 ng/µL and a yield of 2 720.00 ng.

## PacBio HiFi library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core. Libraries were prepared using the SMRTbell Prep Kit 3.0 (Pacific Biosciences, California, USA), following the manufacturer's instructions. The kit includes reagents for end repair/A-tailing, adapter ligation, post-ligation SMRTbell bead clean-up, and nuclease treatment. Size selection and clean-up were performed using diluted AMPure PB beads (Pacific Biosciences). DNA concentration was quantified using a Qubit Fluorometer v4.0 (ThermoFisher Scientific) and the Qubit 1X dsDNA HS assay kit. Final library fragment size was assessed with the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) using the gDNA 55 kb BAC analysis kit.

The sample was sequenced on a Revio instrument (Pacific Biosciences). The prepared library was normalised to 2 nM, and 15 µL 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 head and thorax tissue of the ilAcrRumi1 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 to 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 MERQURY.FK (Rhie et al., 2020).

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

## Assembly curation

The assembly was decontaminated using the Assembly Screen for Cobionts and Contaminants (ASCC) pipeline. TreeVal was used to generate the flat files and maps for use in curation.

Manual curation was conducted primarily in PretextView and HiGlass (Kerpedjiev et al., 2018). Scaffolds were visually inspected and corrected as described by Howe et al. (2021). Manual corrections included 10 breaks and 24 joins. The curation process is documented at https://gitlab.com/wtsi-grit/rapid-curation. 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 Blob-ToolKit 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 *Acronicta rumicis* specimen generated 49.66 Gb (gigabases) from 5.23 million reads, which were used to assemble the genome. GenomeScope2.0 analysis estimated the haploid genome size at 574.01 Mb, with a heterozygosity of 2.56% and repeat content of 37.25% (Figure 2). These estimates guided expectations for the assembly. Based on the estimated genome size, the sequencing data provided approximately 84× coverage. Hi-C sequencing produced 109.54 Gb from 725.44 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 582.86 Mb in 53 scaffolds, with 20 gaps, and a scaffold N50 of 19.31 Mb (Table 2).

Most of the assembly sequence (99.6%) was assigned to 32 chromosomal-level scaffolds, representing 31 autosomes and the W and Z sex chromosomes. These chromosome-level scaffolds, confirmed by Hi-C data, are named according to size (Figure 3; Table 3). The Z and W chromosomes were identified by copy number in the diploid assembly. During curation, we noted that the exact order and orientation of the contigs on chromosome 23 (10.8–11.63Mbp) are uncertain. 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.



**GenomeScope Profile**

len:574,012,976bp uniq:62.8%
aa:97.4% ab:2.56%
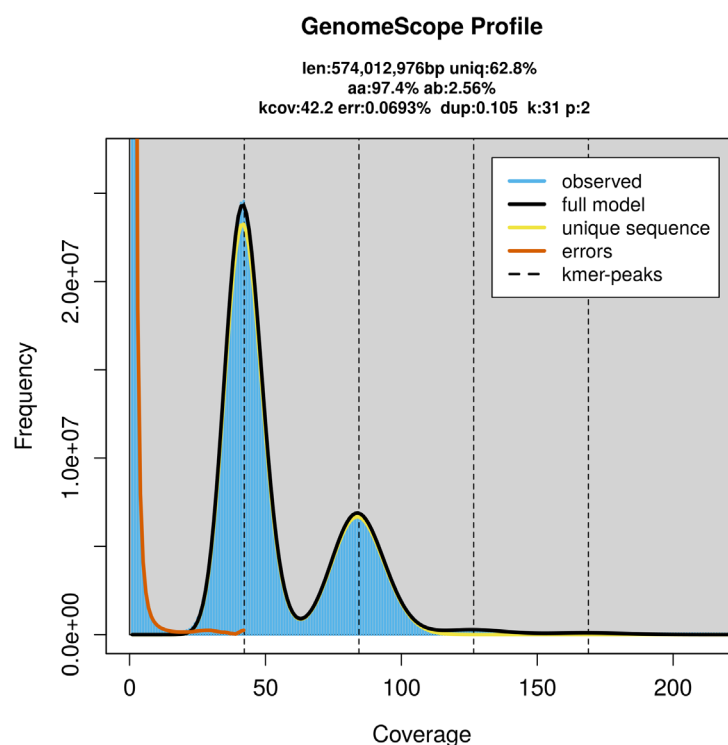kcov:42.2 err:0.0693% dup:0.105 k:31 p:2

**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 PRJEB86626.**

| Platform | PacBio HiFi | Hi-C |
|---|---|---|
| **ToLID** | ilAcrRumi2 | ilAcrRumi1 |
| **Specimen ID** | NHMUK013805967 | Ox000692 |
| **BioSample (source individual)** | SAMEA111458015 | SAMEA7701553 |
| **BioSample (tissue)** | SAMEA111458106 | SAMEA7701743 |
| **Tissue** | thorax | head \| thorax |
| **Instrument** | Revio | Illumina NovaSeq 6000 |
| **Run accessions** | ERR14777920 | ERR14782848 |
| **Read count total** | 5.23 million | 725.44 million |
| **Base count total** | 49.66 Gb | 109.54 Gb |

**Table 2. Genome assembly statistics.**

| Assembly name | ilAcrRumi2.hap1.1 | ilAcrRumi2.hap2.1 |
|---|---|---|
| **Assembly accession** | GCA_965234005.1 | GCA_965233995.1 |
| **Assembly level** | chromosome | scaffold |
| **Span (Mb)** | 582.86 | 528.05 |
| **Number of chromosomes** | 32 | N/A |
| **Number of contigs** | 73 | 512 |
| **Contig N50** | 18.93 Mb | 8.12 Mb |
| **Number of scaffolds** | 53 | 393 |
| **Scaffold N50** | 19.31 Mb | 18.34 Mb |
| **Longest scaffold length (Mb)** | 30.06 | N/A |
| **Sex chromosomes** | W and Z | N/A |
| **Organelles** | Mitochondrion: 15.39 kb | N/A |

For haplotype 1, the estimated QV is 67.5, and for haplotype 2, 66.1. When the two haplotypes are combined, the assembly achieves an estimated QV of 66.8. The $k$-mer completeness is 67.04% for haplotype 1, 62.36% for haplotype 2, and 99.58% for the combined haplotypes (Figure 4).

BUSCO analysis using the lepidoptera_odb10 reference set ($n$ = 5 286) identified 98.8% of the expected gene set (single = 98.4%, duplicated = 0.4%) 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 **7.C.Q67**, meeting the recommended reference standard.

## Wellcome Sanger Institute – Legal and Governance

The materials that have contributed to this genome note have been supplied by a Darwin Tree of Life Partner. The submission of materials by a Darwin Tree of Life Partner is subject to the **'Darwin Tree of Life Project Sampling Code of Practice'**, which can be found in full on the Darwin Tree of Life website. By agreeing with and signing up to the Sampling Code of Practice, the Darwin Tree of Life Partner agrees they will meet the legal and ethical requirements and standards set out within this document in respect of all samples acquired for, and supplied to, the Darwin Tree of Life Project. Further, the Wellcome Sanger Institute employs a process whereby due diligence is carried out proportionate to the nature of the materials themselves, and the circumstances under which they
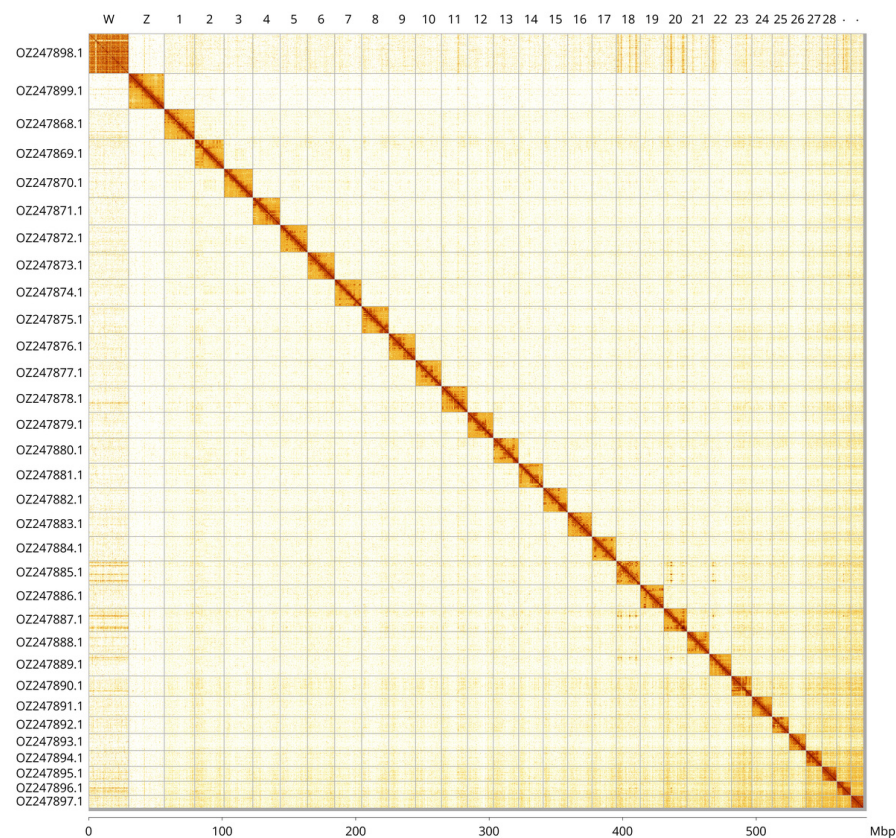
**Figure 3. Hi-C contact map of the *Acronicta rumicis* genome assembly.** Assembled chromosomes are shown in order of size and labelled along the axes, with a megabase scale shown below. The plot was generated using PretextSnapshot.

**Table 3. Chromosomal pseudomolecules in the haplotype 1 genome assembly of *Acronicta rumicis* ilAcrRumi2.**

| INSDC accession | Molecule | Length (Mb) | GC% |
| --- | --- | --- | --- |
| OZ247868.1 | 1 | 22.77 | 37 |
| OZ247869.1 | 2 | 22 | 37 |
| OZ247870.1 | 3 | 21.44 | 37 |
| OZ247871.1 | 4 | 20.51 | 36.50 |
| OZ247872.1 | 5 | 20.48 | 36.50 |
| OZ247873.1 | 6 | 20.33 | 37 |
| OZ247874.1 | 7 | 20.31 | 37 |
| OZ247875.1 | 8 | 20.26 | 37 |
| OZ247876.1 | 9 | 20.02 | 37 |
| OZ247877.1 | 10 | 19.54 | 37 |
| OZ247878.1 | 11 | 19.44 | 37 |
| OZ247879.1 | 12 | 19.31 | 36.50 |
| OZ247880.1 | 13 | 18.93 | 37 |
| OZ247881.1 | 14 | 18.45 | 37 |
| OZ247882.1 | 15 | 18.31 | 37 |

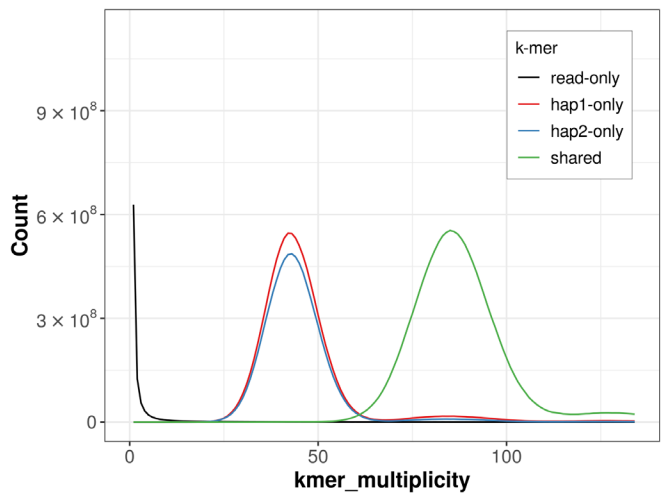| INSDC accession | Molecule | Length (Mb) | GC% |
| --- | --- | --- | --- |
| OZ247883.1 | 16 | 18.23 | 37 |
| OZ247884.1 | 17 | 18.20 | 37 |
| OZ247885.1 | 18 | 17.92 | 37.50 |
| OZ247886.1 | 19 | 17.60 | 37.50 |
| OZ247887.1 | 20 | 17.50 | 37.50 |
| OZ247888.1 | 21 | 16.72 | 37.50 |
| OZ247889.1 | 22 | 16.47 | 37 |
| OZ247890.1 | 23 | 15.43 | 38 |
| OZ247891.1 | 24 | 15.06 | 37.50 |
| OZ247892.1 | 25 | 12.73 | 37.50 |
| OZ247893.1 | 26 | 12.63 | 37.50 |
| OZ247894.1 | 27 | 12.09 | 38 |
| OZ247895.1 | 28 | 11.08 | 38 |
| OZ247896.1 | 29 | 10.59 | 38.50 |
| OZ247897.1 | 30 | 9.43 | 38.50 |
| OZ247898.1 | W | 30.06 | 39.50 |
| OZ247899.1 | Z | 26.71 | 36.50 |

**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.
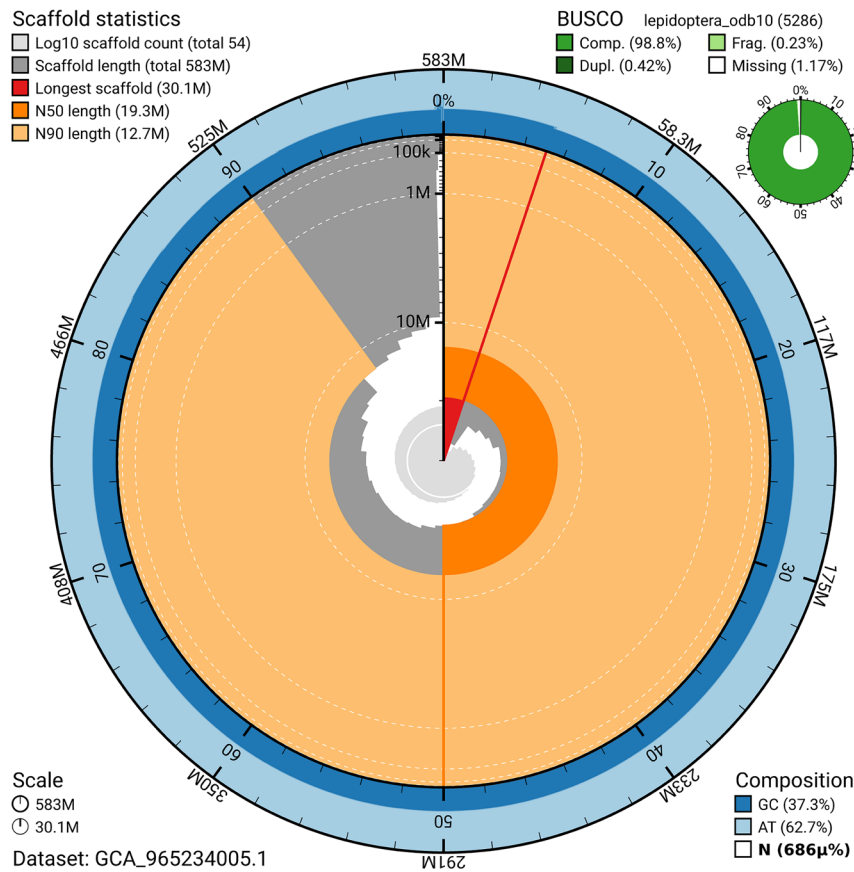


**Figure 5. Assembly metrics for ilAcrRumi2.hap1.1.** The BlobToolKit snail plot provides an overview of assembly metrics and BUSCO gene completeness. The circumference represents the length of the whole genome sequence, and the main plot is divided into 1 000 bins around the circumference. The outermost 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.
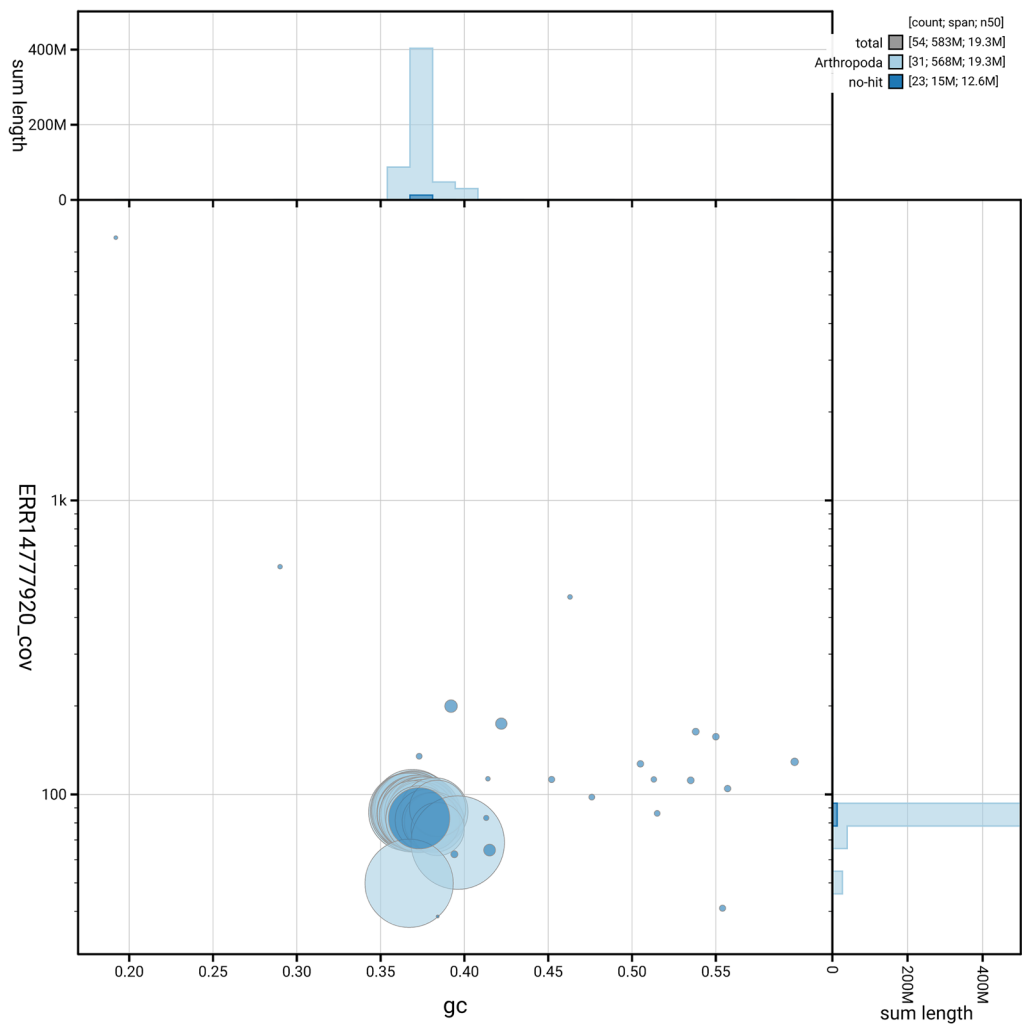
**Figure 6. BlobToolKit GC-coverage plot for ilAcrRumi2.hap1.1.** Blob plot showing sequence coverage (vertical axis) and GC content (horizontal axis). The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. The histograms along the 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.

**Table 4. Earth Biogenome Project summary metrics for the *Acronicta rumicis* assembly.**

| Measure | Value | Benchmark |
|---|---|---|
| EBP summary (haplotype 1) | 7.C.Q67 | 6.C.Q40 |
| Contig N50 length | 18.93 Mb | ≥1 Mb |
| Scaffold N50 length | 19.31 Mb | = chromosome N50 |
| Consensus quality (QV) | Haplotype 1: 67.5; haplotype 2: 66.1; combined: 66.8 | ≥40 |
| *k*-mer completeness | Haplotype 1: 67.04%; Haplotype 2: 62.36%; combined: 99.58% | ≥95% |
| BUSCO | C:98.8% [S:98.4%; D:0.4%]; F:0.2%; M:0.9%; n:5 286 | S > 90%; D < 5% |
| Percentage of assembly assigned to chromosomes | 99.60% | ≥90% |

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: Acronicta rumicis (knot grass). Accession number PRJEB86626. The genome sequence is released openly for reuse. The *Acronicta rumicis* genome sequencing initiative is part of the Darwin Tree of Life Project (PRJEB40665), the Sanger Institute Tree of Life Programme (PRJEB43745) and Project Psyche (PRJEB71705). All raw sequence data and the assembly have been deposited in INSDC databases. The genome will be annotated using available RNA-Seq data and presented through the Ensembl pipeline at the European Bioinformatics Institute. Raw data and assembly accession identifiers are reported in Table 1 and Table 2.

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

**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 |
| 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 |

| Software | Version | Source |
|----------|---------|--------|
| sanger-tol/blobtoolkit | v0.7.1 | https://github.com/sanger-tol/blobtoolkit |
| Seqtk | 1.3 | https://github.com/lh3/seqtk |
| Singularity | 3.9.0 | https://github.com/sylabs/singularity |
| TreeVal | 1.2.0 | 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 Natural History Museum Genome Acquisition Lab

- 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**

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**

Butterfly Conservation, Rothamsted Research: **The state of Britain's larger moths: summary data and analysis.** 2006.
**Reference Source**

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): 1332–1335.
**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**

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**

Fox R, Conrad KF, Parsons MS, *et al.*: **The state of Britain's larger moths.** 2006.
**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**

Joint Nature Conservation Committee (JNCC): **UK biodiversity action plan priority species: Terrestrial invertebrates.** 2007.
**Reference Source**

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**

Kimber I: **Knot grass (*Acronicta rumicis*).** 2025.
**Reference Source**

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**

National Museums Northern Ireland: **Northern Ireland priority species:** *Acronicta rumicis***.** 2025.
**Reference Source**

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**

The Wildlife Trusts: **Knot grass moth (*Acronicta rumicis*).** 2025.
**Reference Source**

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**

Waring P, Townsend M, Lewington R: **Field guide to the moths of Great Britain and Ireland, Third Edition.** Bloomsbury Publishing, 2017.
**Reference Source**

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**