DATA NOTE

# The genome sequence of the Violet Fritillary, *Boloria dia* (Linnaeus, 1767) (Lepidoptera: Nymphalidae)

[version 1; peer review: awaiting peer review]

Yannick Chittaro [ID][1], Mathieu Joron [ID][2], Kay Lucek [ID][3], Charlotte J. Wright [ID][4],
Joana I. Meier[4], Mark L. Blaxter [ID][4],
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, Project Psyche Community

[1]Info fauna, Neuchâtel, Switzerland
[2]CEFE CNRS, Montpellier, France
[3]University of Neuchâtel, Neuchâtel, Switzerland
[4]Tree of Life, Wellcome Sanger Institute, Hinxton, England, UK

**Open Peer Review**

**Approval Status** *AWAITING PEER REVIEW*

Any reports and responses or comments on the article can be found at the end of the article.

## Abstract

We present a genome assembly from a female specimen of *Boloria dia* (Violet Fritillary; Arthropoda; Insecta; Lepidoptera; Nymphalidae). The assembly contains two haplotypes with total lengths of 366.13 megabases and 333.50 megabases. Most of haplotype 1 (95.67%) 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 15.16 kilobases.

## Keywords

Boloria dia, Violet Fritillary, genome sequence, chromosomal, Lepidoptera

This article is included in the Tree of Life gateway.

**Corresponding author:** Charlotte J. Wright (cw22@sanger.ac.uk)

**Author roles: Chittaro Y**: Investigation, Resources; **Joron M**: Writing – Original Draft Preparation; **Lucek K**: Resources, Supervision; **Wright CJ**: Funding Acquisition, Project Administration, Supervision; **Meier JI**: Funding Acquisition, Project Administration, Supervision; **Blaxter ML**: Funding Acquisition, Project Administration, Supervision;

## 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; Papilionoidea; Nymphalidae; Heliconiinae; Argynnini; *Boloria*; *Clossiana*; *Boloria dia* (Linnaeus, 1767) (NCBI:txid596532)

## Background

*Clossiana dia* (Linnaeus, 1767), the Violet Fritillary, is a small nymphalid butterfly from a large genus occupying the Holarctic area (Simonsen *et al.*, 2010; Tuzov & Bozano, 2006). *Clossiana* is a monophyletic group now most often considered a subgenus of *Boloria* (Simonsen *et al.*, 2010) (see the nomenclature adopted by GBIF and the IUCN, Van Swaay *et al.*, 2010). As with most members of *Clossiana*, *C. dia* has a zonal Palaearctic distribution, ranging from southwestern Europe to eastern Russia (Tuzov & Bozano, 2006). In Europe, the Violet Fritillary is found in most temperate areas, but absent from northwestern regions (including the British and Irish Isles where it is an occasional visitor), and localised in Mediterranean regions. This species is listed as Least Concern by the IUCN (Van Swaay *et al.*, 2010).

*Clossiana dia* inhabits dry, sunny meadows and is often seen cruising at low height with a rapid, steady flight. The species exhibits a multivoltine life cycle with one to three generations per year depending on local conditions, but is rarely found above 1 200–1 500 metres, depending on latitude. Its larvae feed on various Viola species, and eggs are sometimes laid away from the host. *C. dia* displays variation in local wing patterns, notably in melanism, but has a distinctive pattern of purple and silver spots on its hindwing underside (Tuzov & Bozano, 2006).

This common species is part of genus widely used as a model to understand dispersal ecology in metapopulations, and in conservation (Baguette, 2003; Vandewoestijne & Baguette, 2002). Fritillaries have colonised many different host plant families, making the genus an excellent system to study how diet broadening contributes to niche and range expansions, especially into boreal habitats (Simonsen *et al.*, 2010). Other *Boloria* species inhabit high altitudes or circumpolar regions, making the group an interesting case to study contrasted responses to global warming via habitat loss. Finally, *Boloria* s.l. is a complex genus, with many systematic uncertainties within subclades.

The reference genome presented here constitutes a useful tool for better understanding butterfly voltinism, community dynamics, evolution, and biogeography in temperate and boreal biomes. The sequence data was derived from a female specimen (Figure 1) collected from Euseigne VS, Switzerland.

## Methods
### Sample acquisition

The specimen used for genome sequencing was an adult female *Boloria dia* (specimen ID SAN28000131, ToLID ilCloDiax1;



**Figure 1. Voucher photograph of the *Boloria dia* (ilCloDiax1) specimen used for genome sequencing.**

Figure 1), collected from Euseigne VS, Switzerland (latitude 46.1773, longitude 7.4159; elevation 800 m) on 26/04/2023. The specimen was collected and identified by Yannick Chittaro (Info Fauna, Neuchâtel, Switzerland).

### 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 ilCloDiax1 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 in the WSI Scientific Operations core using the Automated MagAttract v2 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 36.02 ng/μL and a yield of 1 692.94 ng, with a fragment size of 15.5 kb. The 260/280 spectrophotometric ratio was 1.94, and the 260/230 ratio was 2.94.

### PacBio HiFi library preparation and sequencing

Library preparation and sequencing were performed at the WSI Scientific Operations core. Samples with an average fragment size greater than 8 kb and total mass exceeding 400 ng were eligible for the low-input SMRTbell Prep Kit 3.0 protocol (Pacific Biosciences, California, USA), depending on genome size and required sequencing depth. Libraries were prepared

using the SMRTbell Prep Kit 3.0 according to 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.

Specimen details, sequencing platforms, and data yields are summarised in Table 1.

## Hi-C
### Sample preparation and crosslinking
The Hi-C sample was prepared from 20–50 mg of frozen head tissue of the ilCloDiax1 sample using the Arima-HiC v2 kit (Arima Genomics). Following the manufacturer's instructions, tissue was fixed and DNA crosslinked using TC buffer to a final formaldehyde concentration of 2%. The tissue was homogenised using the Diagnocine Power Masher-II. Crosslinked DNA was digested with a restriction enzyme master mix, biotinylated, and ligated. Clean-up was performed with SPRISelect beads before library preparation. DNA concentration was measured with the Qubit Fluorometer (Thermo Fisher Scientific) and Qubit HS Assay Kit. The biotinylation percentage was estimated using the Arima-HiC v2 QC beads.

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

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

Specimen details, sequencing platforms, and data yields are summarised in Table 1.

## 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.

**Table 1. Specimen and sequencing data for BioProject PRJEB78764.**

| Platform | PacBio HiFi | Hi-C |
|---|---|---|
| ToLID | ilCloDiax1 | ilCloDiax1 |
| Specimen ID | SAN28000131 | SAN28000131 |
| BioSample (source individual) | SAMEA115110009 | SAMEA115110009 |
| BioSample (tissue) | SAMEA115110041 | SAMEA115110039 |
| Tissue | thorax | head |
| Sequencing platform and model | Revio | Illumina NovaSeq X |
| Run accessions | ERR13485728 | ERR13493985 |
| Read count total | 2.13 million | 1 016.22 million |
| Base count total | 23.29 Gb | 153.45 Gb |

## 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 18 breaks and 31 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

Chromosomal painting was performed using lep_busco_painter using Merian elements, which represent the 32 ancestral linkage groups in Lepidoptera (Wright *et al.*, 2024). Painting was based on gene locations from the lepidoptera_odb10 BUSCO analysis and chromosome lengths from the genome index produced using SAMtools faidx (Danecek *et al.*, 2021). Each complete BUSCO (including both single-copy and duplicated BUSCOs) was assigned to a Merian element using a reference database, and coloured positions were plotted along chromosomes drawn to scale.

The Merqury.FK tool (Rhie *et al.*, 2020), run in a Singularity container (Kurtzer *et al.*, 2017), was used to evaluate $k$-mer completeness and assembly quality for 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 BlobToolKit pipeline (Challis *et al.*, 2020). The pipeline aligns PacBio reads using minimap2 (Li, 2018) and SAMtools (Danecek *et al.*, 2021) to generate coverage tracks. Simultaneously, it queries the GoaT database (Challis *et al.*, 2023) to identify relevant BUSCO lineages and runs BUSCO (Manni *et al.*, 2021). For the three domain-level BUSCO 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 package management via Conda and Bioconda (Grüning *et al.*, 2018), and containerisation through Docker (Merkel, 2014) and Singularity (Kurtzer *et al.*, 2017).

## Genome sequence report

### Sequence data

The genome of a specimen of *Boloria dia* was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating 23.29 Gb (gigabases) from 2.13 million reads, which were used to assemble the genome. GenomeScope2.0 analysis estimated the haploid genome size at 347.14 Mb, with a heterozygosity of 1.50% and repeat content of 20.18%. These estimates guided expectations for the assembly. Based on the estimated genome size, the sequencing data provided approximately 63× coverage. Hi-C sequencing produced 153.45 Gb from 1 016.22 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 366.13 Mb in 121 scaffolds, with 76 gaps, and a scaffold N50 of 12.04 Mb (Table 2).

Most of the assembly sequence (95.67%) was assigned to 32 chromosomal-level scaffolds, representing 30 autosomes and the W and Z sex chromosomes. The chromosome-level scaffolds, confirmed by Hi-C data, are named according to size (Figure 2; Table 3). The W chromosome is highly fragmented and repetitive and so most sequence has been assigned as unlocalised. Chromosome painting with Merian elements illustrates the distribution of orthologues along chromosomes and highlights patterns of chromosomal evolution relative to Lepidopteran ancestral linkage groups (Figure 3).

### Assembly quality metrics

For haplotype 1, the estimated QV is 60.3, and for haplotype 2, 61.2. When the two haplotypes are combined, the assembly achieves an estimated QV of 60.7. The $k$-mer completeness is

**Table 2. Genome assembly statistics.**

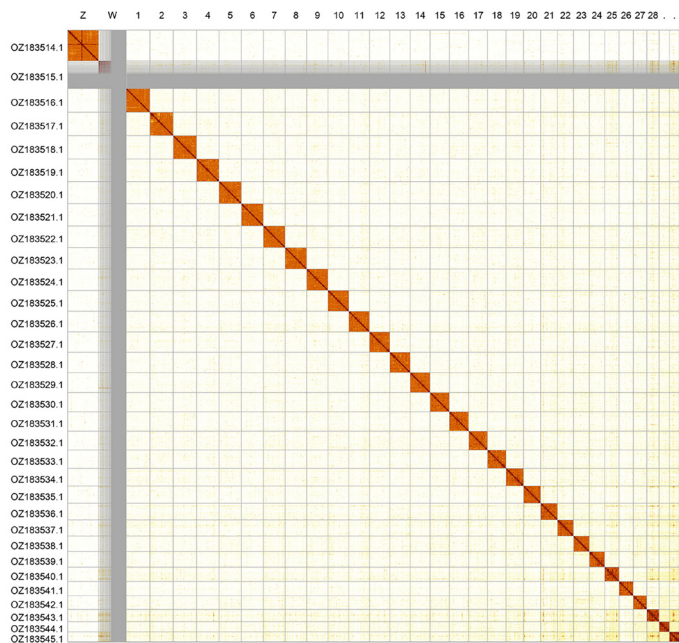| Assembly name | ilCloDiax1.hap1.1 | ilCloDiax1.hap2.1 |
|---|---|---|
| **Assembly accession** | GCA_964264805.1 | GCA_964264855.1 |
| **Assembly level** | chromosome | scaffold |
| **Span (Mb)** | 366.13 | 333.50 |
| **Number of chromosomes** | 32 | N/A |
| **Number of contigs** | 197 | 159 |
| **Contig N50** | 5.87 Mb | 5.63 Mb |
| **Number of scaffolds** | 121 | 106 |
| **Scaffold N50** | 12.04 Mb | 12.04 Mb |
| **Longest scaffold length (Mb)** | 18.69 | N/A |
| **Sex chromosomes** | W and Z | N/A |
| **Organelles** | Mitochondrial genome: 15.16 kb | N/A |

**Figure 2. Hi-C contact map of the *Boloria dia* 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 *Boloria dia* ilCloDiax1.**

| INSDC accession | Molecule | Length (Mb) | GC% | Assigned Merian elements |
|---|---|---|---|---|
| OZ183516.1 | 1 | 14.11 | 33 | M2 |
| OZ183517.1 | 2 | 14.03 | 33 | M17;M20 |
| OZ183518.1 | 3 | 13.96 | 33 | M1 |
| OZ183519.1 | 4 | 13.41 | 33 | M3 |
| OZ183520.1 | 5 | 13.37 | 32.50 | M8 |
| OZ183521.1 | 6 | 13.18 | 32.50 | M9 |
| OZ183522.1 | 7 | 12.89 | 33 | M7 |
| OZ183523.1 | 8 | 12.89 | 32.50 | M12 |
| OZ183524.1 | 9 | 12.87 | 33 | M5 |
| OZ183525.1 | 10 | 12.42 | 32.50 | M16 |
| OZ183526.1 | 11 | 12.38 | 32.50 | M18 |
| OZ183527.1 | 12 | 12.19 | 33 | M4 |
| OZ183528.1 | 13 | 12.04 | 33 | M21 |
| OZ183529.1 | 14 | 11.94 | 33 | M6 |
| OZ183530.1 | 15 | 11.58 | 33 | M22 |
| OZ183531.1 | 16 | 11.50 | 33 | M15 |

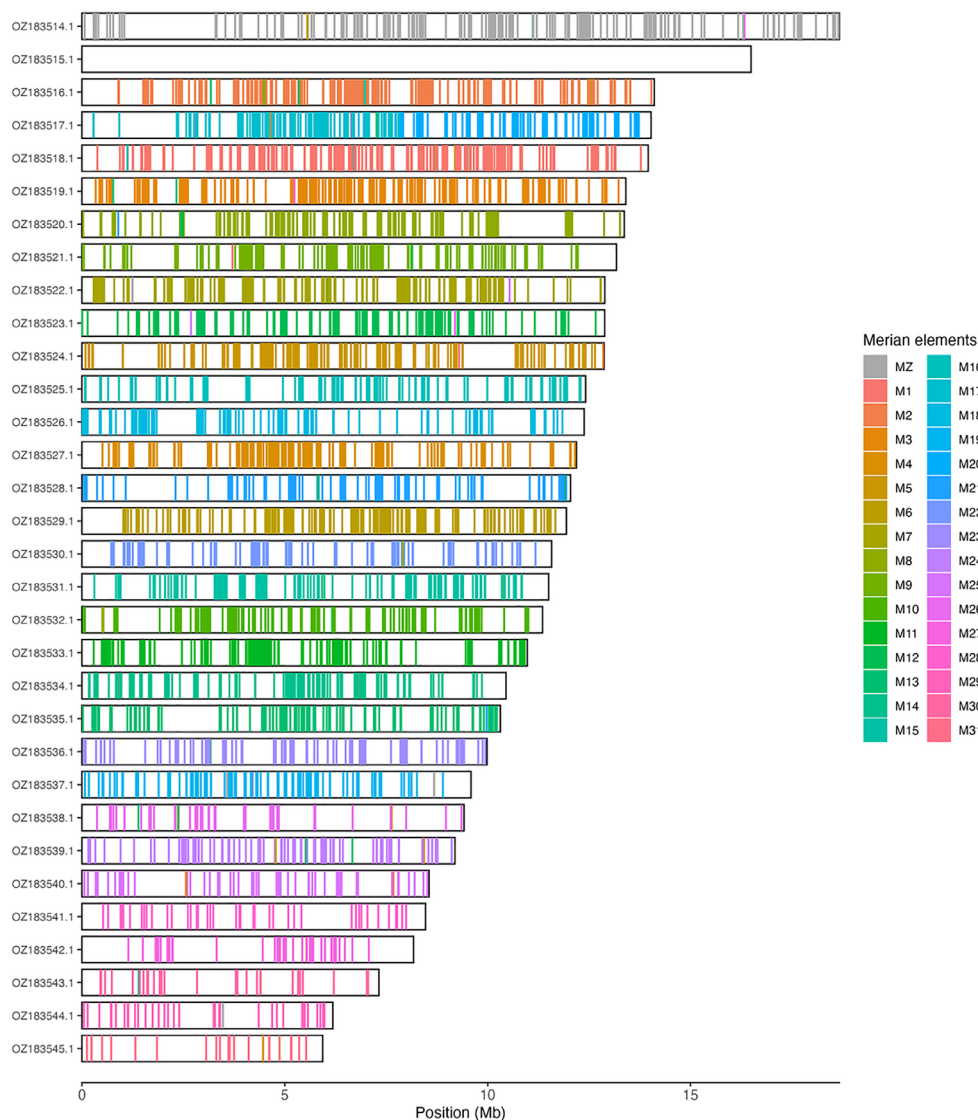| INSDC accession | Molecule | Length (Mb) | GC% | Assigned Merian elements |
|---|---|---|---|---|
| OZ183532.1 | 17 | 11.35 | 33 | M10 |
| OZ183533.1 | 18 | 10.97 | 33.50 | M11 |
| OZ183534.1 | 19 | 10.45 | 33.50 | M14 |
| OZ183535.1 | 20 | 10.31 | 33 | M13 |
| OZ183536.1 | 21 | 9.98 | 33.50 | M23 |
| OZ183537.1 | 22 | 9.59 | 34 | M19 |
| OZ183538.1 | 23 | 9.42 | 33.50 | M26 |
| OZ183539.1 | 24 | 9.19 | 33.50 | M24 |
| OZ183540.1 | 25 | 8.56 | 37 | M25 |
| OZ183541.1 | 26 | 8.47 | 33.50 | M28 |
| OZ183542.1 | 27 | 8.17 | 33.50 | M27 |
| OZ183543.1 | 28 | 7.32 | 36.50 | M30 |
| OZ183544.1 | 29 | 6.18 | 35.50 | M29 |
| OZ183545.1 | 30 | 5.93 | 36 | M31 |
| OZ183515.1 | W | 0.98 | 34 | N/A |
| OZ183514.1 | Z | 18.69 | 32.50 | MZ |
| OZ183546.1 | MT | 0.02 | 19 | N/A |

**Figure 3. Merian elements painted across chromosomes in the ilCloDiax1.hap1.1 assembly of *Boloria dia*.** Chromosomes are drawn to scale, with the positions of orthologues shown as coloured bars. Each orthologue is coloured by the Merian element that it belongs to. All orthologues which could be assigned to Merian elements are shown.

74.65% for haplotype 1, 70.59% for haplotype 2, and 99.49% for the combined haplotypes (Figure 4). BUSCO analysis using the lepidoptera_odb10 reference set ($n$ = 5 286) (Kriventseva *et al.*, 2019) identified 98.6% of the expected gene set (single = 98.1%, duplicated = 0.5%) 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.

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.

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 **6.C.Q60**, 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 Tree of Life collaborator. 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
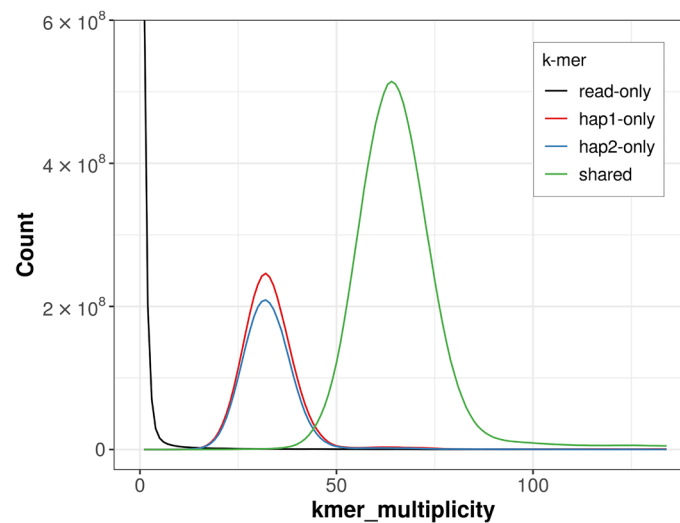
**Figure 4. Evaluation of *k*-mer completeness using MerquryFK.** This plot illustrates the recovery of *k*-mers from the original read data in the final assemblies. The horizontal axis represents *k*-mer multiplicity, and the vertical axis shows the number of *k*-mers. The black curve represents *k*-mers that appear in the reads but are not assembled. The green curve (the homozygous peak) corresponds to *k*-mers shared by both haplotypes and the red and blue curves (the heterozygous peaks) show *k*-mers found only in one of the haplotypes.
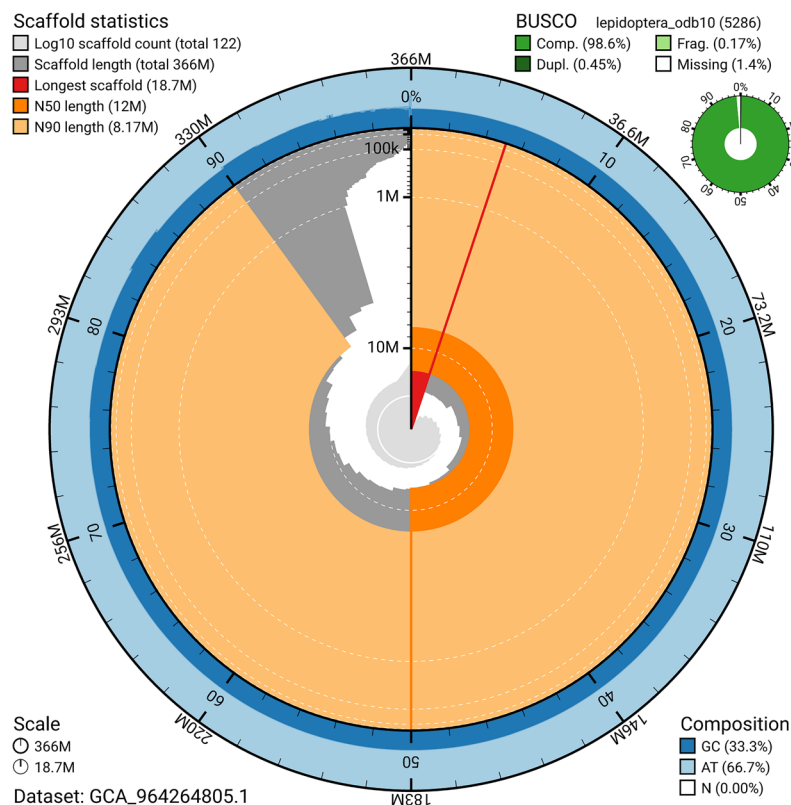


**Figure 5. Assembly metrics for ilCloDiax1.hap1.1.** The BlobToolKit snail plot provides an overview of assembly metrics and BUSCO gene completeness. The circumference represents the length of the whole genome sequence, and the main plot is divided into 1,000 bins around the circumference. The outermost 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.
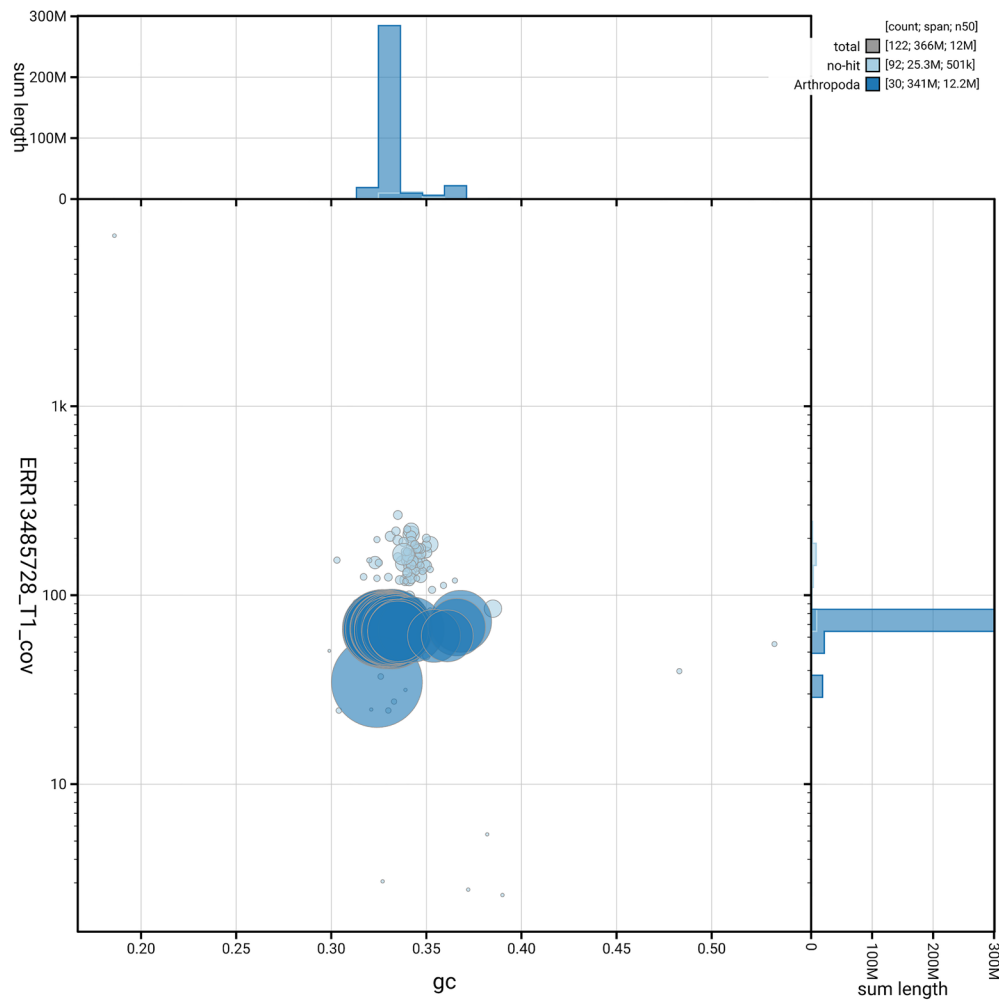
**Figure 6. BlobToolKit GC-coverage plot for ilCloDiax1.hap1.1.** Blob plot showing sequence coverage (vertical axis) and GC content (horizontal axis). The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. The histograms along the 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 *Boloria dia* assembly.**

| Measure (Benchmark) | Value |
| --- | --- |
| EBP summary (haplotype 1) | 6.C.Q60 |
| Contig N50 length (≥ 1 Mb) | 5.87 Mb |
| Scaffold N50 length (= chromosome N50) | 12.04 Mb |
| Consensus quality (QV) (≥ 40) | Haplotype 1: 60.3; haplotype 2: 61.2; combined: 60.7 |
| *k*-mer completeness (≥ 95%) | Haplotype 1: 74.65%; Haplotype 2: 70.59%; combined: 99.49% |
| BUSCO (S > 90%; D < 5%) | C:98.6%[S:98.1%,D:0.5%],F:0.2%,M:1.2%,n:5286 |
| Percentage of assembly assigned to chromosomes (≥ 90%) | 95.67% |

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 undertaken according to a Research Collaboration Agreement or Material Transfer Agreement entered into by the Tree of Life collaborator, Genome Research Limited (operating as the Wellcome Sanger Institute), and in some circumstances, other Tree of Life collaborators.

## Data availability

European Nucleotide Archive: Clossiana dia (violet fritillary). Accession number PRJEB78764. The genome sequence is released openly for reuse. The *Boloria dia* genome sequencing initiative is part of 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 Ensembl at the European Bioinformatics Institute. Raw data and assembly accession identifiers are reported in Table 1 and Table 2.

Pipelines used for genome assembly at the WSI Tree of Life are available at https://pipelines.tol.sanger.ac.uk/pipelines. Table 5 lists software versions used in this study.

## Author information
Contributors are listed at the following links:

- 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

- Project Psyche Community.

**Table 5. Software versions and sources.**

| Software | Version | Source |
|---|---|---|
| BEDTools | 2.30.0 | https://github.com/arq5x/bedtools2 |
| BLAST | 2.14.0 | ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/ |
| BlobToolKit | 4.3.9 | https://github.com/blobtoolkit/blobtoolkit |
| BUSCO | 5.5.0 | https://gitlab.com/ezlab/busco |
| bwa-mem2 | 2.2.1 | https://github.com/bwa-mem2/bwa-mem2 |
| Cooler | 0.8.11 | https://github.com/open2c/cooler |
| DIAMOND | 2.1.8 | https://github.com/bbuchfink/diamond |
| fasta_windows | 0.2.4 | https://github.com/tolkit/fasta_windows |
| FastK | 1.1 | https://github.com/thegenemyers/FASTK |
| GenomeScope2.0 | 2.0.1 | https://github.com/tbenavi1/genomescope2.0 |
| Gfastats | 1.3.6 | https://github.com/vgl-hub/gfastats |
| GoaT CLI | 0.2.5 | https://github.com/genomehubs/goat-cli |
| Hifiasm | 0.19.8-r603 | https://github.com/chhylp123/hifiasm |
| HiGlass | 1.13.4 | https://github.com/higlass/higlass |
| lep_busco_painter | 1.0.0 | https://github.com/charlottewright/lep_busco_painter |
| MerquryFK | 1.1.1 | https://github.com/thegenemyers/MERQURY.FK |
| Minimap2 | 2.24-r1122 | https://github.com/lh3/minimap2 |
| MitoHiFi | 3 | https://github.com/marcelauliano/MitoHiFi |

| Software | Version | Source |
|---|---|---|
| MultiQC | 1.14; 1.17 and 1.18 | https://github.com/MultiQC/MultiQC |
| Nextflow | 23.10.0 | https://github.com/nextflow-io/nextflow |
| PretextSnapshot | N/A | https://github.com/sanger-tol/PretextSnapshot |
| PretextView | 0.2.5 | https://github.com/sanger-tol/PretextView |
| samtools | 1.19.2 | https://github.com/samtools/samtools |
| sanger-tol/ascc | 0.1.0 | https://github.com/sanger-tol/ascc |
| sanger-tol/blobtoolkit | 0.6.0 | https://github.com/sanger-tol/blobtoolkit |
| 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 | 1.2a.2 | https://github.com/c-zhou/yahs |

# References

Allio R, Schomaker-Bastos A, Romiguier J, *et al.*: **MitoFinder: efficient automated large-scale extraction of mitogenomic data in target enrichment phylogenomics.** *Mol Ecol Resour.* 2020; **20**(4): 892–905.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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

Baguette M: **Long distance dispersal and landscape occupancy in a metapopulation of the cranberry fritillary butterfly.** *Ecography.* 2003; **26**(2): 153–60.
**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**

Challis R, Kumar S, Sotero-Caio C, *et al.*: **Genomes on a Tree (GoaT): a versatile, scalable search engine for genomic and sequencing project metadata across the eukaryotic Tree of Life [version 1; peer review: 2 approved].** *Wellcome Open Res.* 2023; **8**: 24.
**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): 1332–1335.
**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**

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**

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**

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

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

Kerpedjiev P, Abdennur N, Lekschas F, *et al.*: **HiGlass: web-based visual exploration and analysis of genome interaction maps.** *Genome Biol.* 2018; **19**(1): 125.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Kriventseva EV, Kuznetsov D, Tegenfeldt F, *et al.*: **OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs.** *Nucleic Acids Res.* 2019; **47**(D1): D807–D811.
**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**

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**

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**

Simonsen TJ, Wahlberg N, Warren AD, *et al.*: **The evolutionary history of** *Boloria* **(Lepidoptera: Nymphalidae): phylogeny, zoogeography and larval-foodplant relationships.** *Syst Biodivers.* 2010; **8**(4): 513–29.
**Publisher Full Text**

Tuzov VK, Bozano GC: **Guide to the Butterflies of the Palearctic Region. Nymphalidae, Part II. Tribe Argynnini, Genera Boloria, Proclossiana and Clossiana.** Milano: Omnes Partes, 2006.

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**

Van Swaay CAM, Wynhoff I, Verovnik R, *et al.*: *Boloria dia* **(Europe assessment).** 2010.

Vandewoestijne S, Baguette M: **The genetic structure of endangered populations in the Cranberry Fritillary,** *Boloria aquilonaris* **(Lepidoptera, Nymphalidae): RAPDs *vs* allozymes.** *Heredity (Edinb).* 2002; **89**(6): 439–45.
**PubMed Abstract** | **Publisher Full Text**

Vasimuddin M, Misra S, Li H, *et al.*: **Efficient architecture-aware acceleration of BWA-MEM for multicore systems.** In: *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS).* IEEE, 2019; 314–324.
**Publisher Full Text**

Wright CJ, Stevens L, Mackintosh A, *et al.*: **Comparative genomics reveals the dynamics of chromosome evolution in Lepidoptera.** *Nat Ecol Evol.* 2024; **8**(4): 777–790.
**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**