Check for updates

DATA NOTE

# The genome sequence of the black scavenger fly, *Nemopoda nitidula* (Fallén, 1820)

[version 1; peer review: 2 approved]

Steven Falk[1], Liam M. Crowley [2], Ruth Y. Akinmusola [3],
University of Oxford and Wytham Woods Genome Acquisition Lab,
Darwin Tree of Life Barcoding collective,
Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory team,
Wellcome Sanger Institute Scientific Operations: Sequencing Operations,
Wellcome Sanger Institute Tree of Life Core Informatics team,
Tree of Life Core Informatics collective, Darwin Tree of Life Consortium

[1]Independent researcher, Kenilworth, Warwickshire, England, UK
[2]Department of Biology, University of Oxford, Oxford, England, UK
[3]University of Bath, Bath, England, UK

## Abstract

We present a genome assembly from an individual specimen of *Nemopoda nitidula* (Arthropoda; Insecta; Diptera; Sepsidae). The genome sequence has a total length of 178.40 megabases. Most of the assembly (99.2%) is scaffolded into 6 chromosomal pseudomolecules. The mitochondrial genome has also been assembled and is 15.98 kilobases in length.

## Keywords

Nemopoda nitidula, black scavenger fly, genome sequence, chromosomal, Diptera

This article is included in the Tree of Life gateway.

**Open Peer Review**

**Approval Status** ✓ ✓

|  | 1 | 2 |
|---|---|---|
| **version 1**<br>07 Feb 2025 | ✓<br>view | ✓<br>view |

1. **Terrence Sylvester** , University of memphis, memphis, USA

2. **Thomas Brown** , Leibniz IZW (Berlin) & ERGA/BGE, Leibniz, Germany

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

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

**Author roles: Falk S**: Investigation, Resources; **Crowley LM**: Investigation, Resources; **Akinmusola RY**: Writing – Original Draft Preparation;

**How to cite this article:** Falk S, Crowley LM, Akinmusola RY *et al.* **The genome sequence of the black scavenger fly,** *Nemopoda nitidula* **(Fallén, 1820) [version 1; peer review: 2 approved]** Wellcome Open Research 2025, **10**:49 https://doi.org/10.12688/wellcomeopenres.23664.1

**First published:** 07 Feb 2025, **10**:49 https://doi.org/10.12688/wellcomeopenres.23664.1

## Species taxonomy

Eukaryota; Opisthokonta; Metazoa; Eumetazoa; Bilateria; Protostomia; Ecdysozoa; Panarthropoda; Arthropoda; Mandibulata; Pancrustacea; Hexapoda; Insecta; Dicondylia; Pterygota; Neoptera; Endopterygota; Diptera; Brachycera; Muscomorpha; Eremoneura; Cyclorrhapha; Schizophora; Acalyptratae; Sciomyzoidea; Sepsidae; *Nemopoda*; *Nemopoda nitidula* (Fallén, 1820) (NCBI:txid292387)

## Background

*Nemopoda nitidula* (Fallén, 1820), a black scavenger fly, is a prominent member of the Sepsidae family of flies. These flies have a spherical head, a petiolate abdomen, and an ant-like glossy black body (Frantsevich & Gladun, 2002; Pont & Meier, 2002). A typical *N. nitidula* adult is primarily black, except for red eyes and pale-yellow front legs. *N. nitidula* males possess sexually dimorphic forelegs, which are useful for successful positioning, mounting, clamping and copulation (Puniamoorthy *et al.*, 2008). Notwithstanding, the females can dislodge the males or halt the copulation attempt via vigorous shaking and other strategies (Hosken *et al.*, 2003).

They are scavengers that subsist on animal dung and decomposing plant and animal remains (Oleksakova *et al.*, 2016; Prado e Castro *et al.*, 2012; Rohner *et al.*, 2019). *N. nitidula* can thrive on various substrates derived from decaying faeces, carrion, vegetation and fungi (Meier, 1996; Oleksakova *et al.*, 2016; Pont & Meier, 2002; Rohner *et al.*, 2019). They also host *Vickermania spadyakhi*, a parasitic flagellate of the Trypanosomatidae family, usually found in their midgut (Kostygov *et al.*, 2020). Thus, they are ecologically relevant in microbial ecology and nutrient cycling.

*N. nitidula* is generally distributed in tropical and temperate climes (Khaghaninia *et al.*, 2014). It prefers damp and cool locations such as woodlands and areas with vegetation cover (Meier, 1996; Pont & Meier, 2002). However, its country-specific distribution varies based on environmental conditions such as altitude, humidity, temperature, wind velocity and seasonal variations (Prado e Castro *et al.*, 2012; Rohner *et al.*, 2019). They are widespread in Europe, with the highest records from the United Kingdom and the Netherlands (GBIF Secretariat, 2024).

Here, we describe the first reference genome for *Nemopoda nitidula*. This genome assembly will promote future investigations on the phylogeny and evolution of different traits in sepsid flies.

## Genome sequence report

### Sequencing data

The genome of a specimen of *Nemopoda nitidula* (Figure 1) was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating 19.16 Gb from 2.08 million reads. GenomeScope analysis of the PacBio HiFi data estimated the haploid genome size at 160.31 Mb, with a heterozygosity of 1.32% and repeat content of 16.08%. These values provide an initial assessment of genome complexity and the challenges anticipated during assembly. Based on this estimated genome size, the sequencing data provided approximately 113.0x



**Figure 1.** Photograph of the *Nemopoda nitidula* (idNemNiti2) specimen used for genome sequencing.

coverage of the genome. Hi-C data produced 93.02 Gb from 616.05 million reads. Table 1 summarises the specimen and sequencing information, including the BioProject, study name, BioSample numbers, and sequencing data for each technology.

### Assembly statistics

The primary haplotype was assembled, and contigs corresponding to an alternate haplotype were also deposited in INSDC databases. The assembly was improved by manual curation, which corrected 50 misjoins or missing joins and removed 6 haplotypic duplications. These interventions reduced the total assembly length by 0.81%, decreased the scaffold count by 47.37%, and increased the scaffold N50 by 1.46%. The final assembly has a total length of 178.45 Mb in 29 scaffolds, with 169 gaps, and a scaffold N50 of 29.83 Mb (Table 2).

The snail plot in Figure 2 provides a summary of the assembly statistics, indicating the distribution of scaffold lengths and other assembly metrics. Figure 3 shows the distribution of scaffolds by GC proportion and coverage. Figure 4 presents a cumulative assembly plot, with separate curves representing different scaffold subsets assigned to various phyla, illustrating the completeness of the assembly.

Most of the assembly sequence (99.22%) was assigned to 6 chromosomal-level scaffolds. These chromosome-level scaffolds, confirmed by Hi-C data, are named according to size (Figure 5; Table 3). Sex chromosome(s) were not identified during curation, as sequence data from the heterogametic sex was not available and homology is unreliable for sex chromosome identification in Diptera due to frequent sex chromosome turnover (Vicoso & Bachtrog, 2015).

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 in GenBank.

### Assembly quality metrics

The estimated Quality Value (QV) and *k*-mer completeness metrics, along with BUSCO completeness scores, were calculated for each haplotype and the combined assembly. The QV

**Table 1. Specimen and sequencing data for *Nemopoda nitidula*.**

| Project information | | | |
|---|---|---|---|
| **Study title** | Nemopoda nitidula | | |
| **Umbrella BioProject** | PRJEB74593 | | |
| **Species** | *Nemopoda nitidula* | | |
| **BioSample** | SAMEA113425719 | | |
| **NCBI taxonomy ID** | 292387 | | |
| **Specimen information** | | | |
| **Technology** | **ToLID** | **BioSample accession** | **Organism part** |
| **PacBio long read sequencing** | idNemNiti2 | SAMEA113425904 | Whole organism |
| **Hi-C sequencing** | idNemNiti1 | SAMEA10979530 | Whole organism |
| **Sequencing information** | | | |
| **Platform** | **Run accession** | **Read count** | **Base count (Gb)** |
| **Hi-C Illumina NovaSeq 6000** | ERR13302608 | 6.16e+08 | 93.02 |
| **PacBio Revio** | ERR12875148 | 2.08e+06 | 19.16 |

**Table 2. Genome assembly data for *Nemopoda nitidula*, idNemNiti2.1.**

| Genome assembly | | |
|---|---|---|
| Assembly name | idNemNiti2.1 | |
| Assembly accession | GCA_964194425.1 | |
| *Accession of alternate haplotype* | GCA_964194455.1 | |
| Span (Mb) | 178.40 | |
| Number of contigs | 199 | |
| Number of scaffolds | 29 | |
| Longest scaffold (Mb) | 41.4 | |
| **Assembly metrics*** | | ***Benchmark*** |
| Contig N50 length (Mb) | 1.7 | *≥ 1 Mb* |
| Scaffold N50 length (Mb) | 29.8 | *= chromosome N50* |
| Consensus quality (QV) | 60.3 | *≥ 40* |
| *k*-mer completeness | 100.0% | *≥ 95%* |
| BUSCO** | C:97.9%[S:97.4%,D:0.5%], F:0.6%,M:1.5%,n:3,285 | *S > 90%, D < 5%* |
| Percentage of assembly mapped to chromosomes | 99.2% | *≥ 90%* |
| Sex chromosomes | Not identified | *localised homologous pairs* |
| Organelles | Mitochondrial genome: 15.98 kb | *complete single alleles* |

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

** BUSCO scores based on the diptera_odb10 BUSCO set using version 5.4.3. C = complete [S = single copy, D = duplicated], F = fragmented, M = missing, n = number of orthologues in comparison. A full set of BUSCO scores is available at https://blobtoolkit.genomehubs.org/view/Nemopoda_nitidula/dataset/GCA_964194425.1/busco.
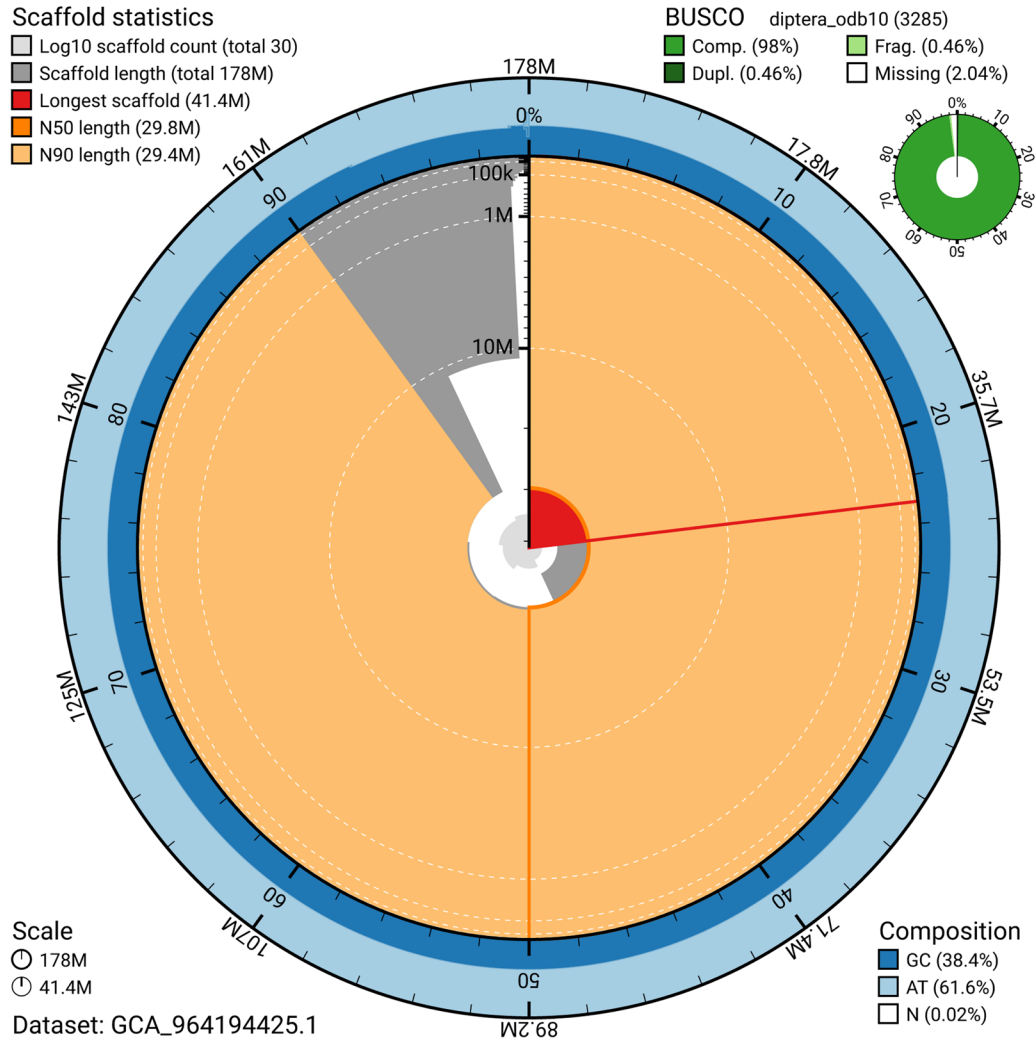
**Figure 2. Genome assembly of *Nemopoda nitidula*, idNemNiti2.1: metrics.** The BlobToolKit snail plot provides an overview of assembly metrics and BUSCO gene completeness. The circumference represents the length of the whole genome sequence, and the main plot is divided into 1,000 bins around the circumference. 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 diptera_odb10 set is presented at the top right. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/GCA_964194425.1/dataset/GCA_964194425.1/snail.

reflects the base-level accuracy of the assembly, while *k*-mer completeness indicates the proportion of expected *k*-mers identified in the assembly. BUSCO scores provide a measure of completeness based on benchmarking universal single-copy orthologues.

The primary haplotype has a QV of 60.3, and the combined primary and alternate assemblies achieve an estimated QV of 60.0. The *k*-mer completeness for the primary haplotype is 75.75%, and for the alternate haplotype it is 76.37%. The combined primary and alternate assemblies achieve a *k*-mer completeness of 99.73%. BUSCO analysis using the

diptera_odb10 reference set ($n = 3,285$) indicated a completeness score of 98.0% (single = 97.5%, duplicated = 0.5%).

Table 2 provides assembly metric benchmarks adapted from Rhie *et al.* (2021) and the Earth BioGenome Project Report on Assembly Standards September 2024. The primary assembly achieves the EBP reference standard of 6.C.60.

## Methods
### Sample acquisition and DNA barcoding
An adult *Nemopoda nitidula* (specimen ID Ox003142, ToLID idNemNiti2) was collected from Wytham Woods, Oxfordshire,
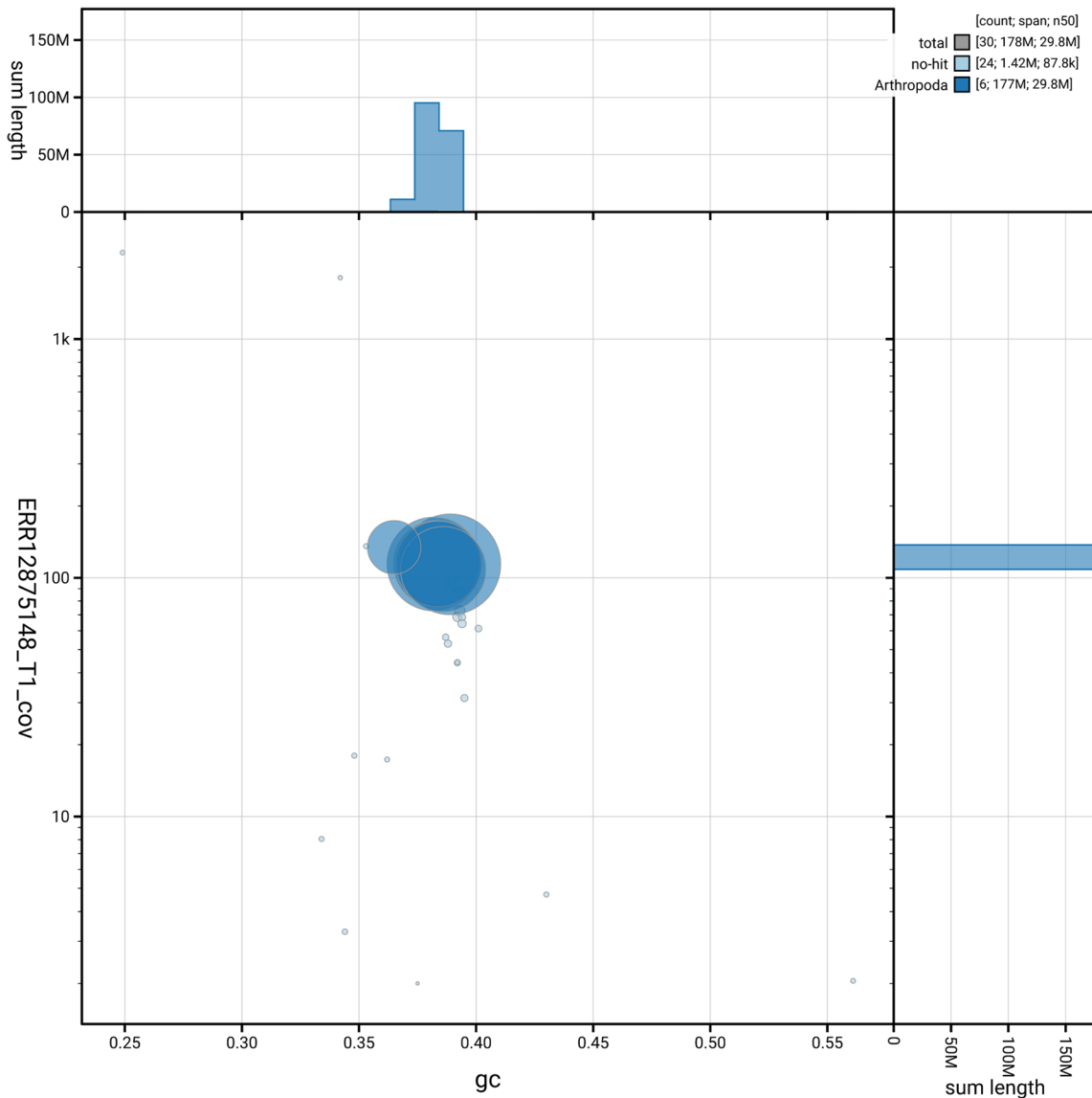
**Figure 3. Genome assembly of *Nemopoda nitidula*, idNemNiti2.1: BlobToolKit GC-coverage plot showing sequence coverage (vertical axis) and GC content (horizontal axis).** The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. The histograms along the axes display the total length of sequences distributed across different levels of coverage and GC content. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/GCA_964194425.1/dataset/GCA_964194425.1/blob.

United Kingdom (latitude 51.76, longitude –1.34) on 2022-10-03 by net. The specimen was collected by Steven Falk (independent researcher) and Liam Crowley (University of Oxford) and identified by Steven Falk. The specimen used for Hi-C sequencing (specimen ID Ox001862, ToLID idNemNiti1) was collected by Liam Crowley from the same location on 2021-09-03 by net. Both specimens were flash frozen on dry ice.

The initial identification was verified by an additional DNA barcoding process according to the framework developed by Twyford *et al.* (2024). A small sample was dissected from the specimens and stored in ethanol, while the remaining parts were shipped on dry ice to the Wellcome Sanger Institute (WSI). The tissue was lysed, the COI marker region was amplified by PCR, and amplicons were sequenced and compared to the BOLD database, confirming the species identification (Crowley *et al.*, 2023). Following whole genome sequence generation, the relevant DNA barcode region was also used alongside the initial barcoding data for sample tracking at the WSI (Twyford *et al.*, 2024). The standard operating procedures for Darwin Tree of Life barcoding have been deposited on protocols.io (Beasley *et al.*, 2023).
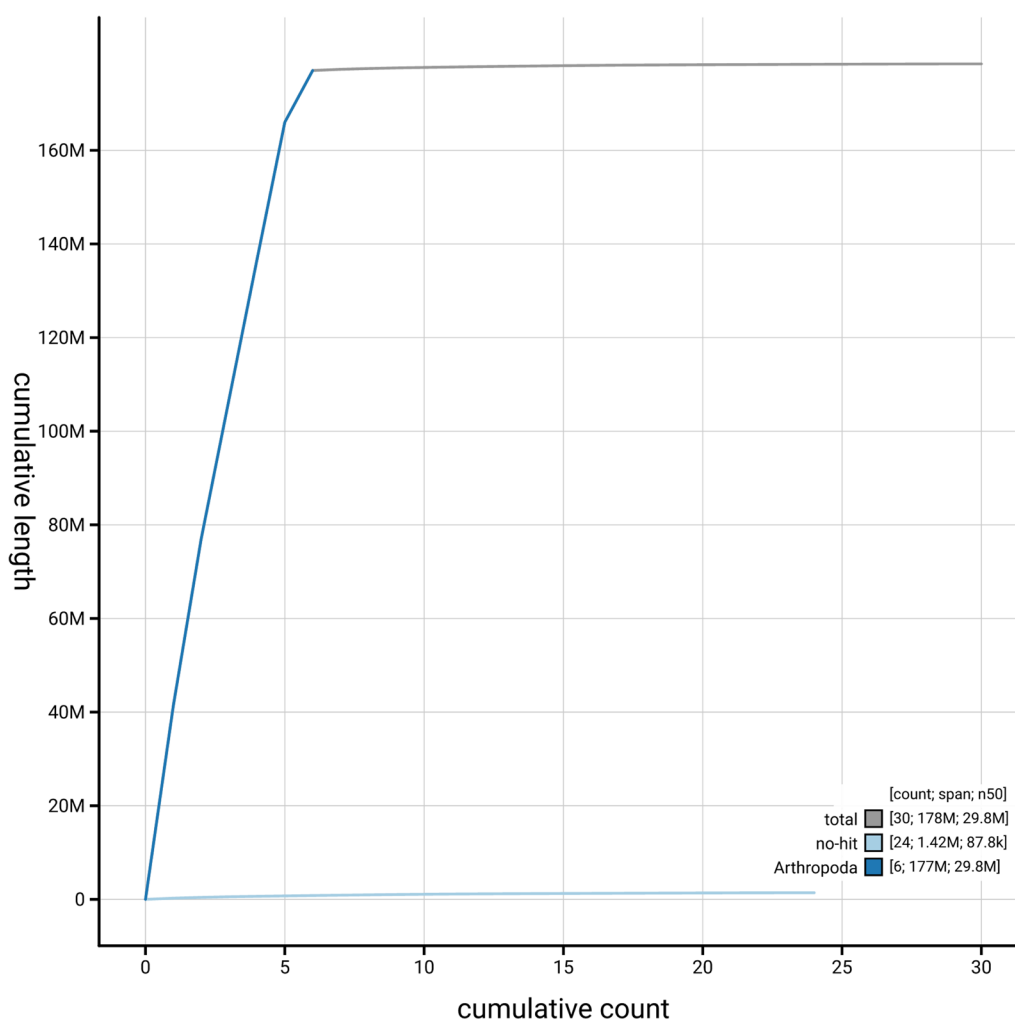
**Figure 4. Genome assembly of *Nemopoda nitidula* idNemNiti2.1: BlobToolKit cumulative sequence plot.** The grey line shows cumulative length for all scaffolds. Coloured lines show cumulative lengths of scaffolds assigned to each phylum using the buscogenes taxrule. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/GCA_964194425.1/dataset/GCA_964194425.1/cumulative.

## Nucleic acid extraction

The workflow for high molecular weight (HMW) DNA extraction at the Wellcome Sanger Institute (WSI) Tree of Life Core Laboratory includes a sequence of procedures: sample preparation and homogenisation, DNA extraction, fragmentation and purification. Detailed protocols are available on protocols.io (Denton *et al.*, 2023b). The idNemNiti2 sample was prepared for DNA extraction by weighing and dissecting it on dry ice (Jay *et al.*, 2023). Tissue from the whole organism was homogenised using a PowerMasher II tissue disruptor (Denton *et al.*, 2023a).

HMW DNA was extracted in the WSI Scientific Operations core using the Automated MagAttract v2 protocol (Oatley *et al.*, 2023). The DNA was sheared into an average fragment size of 12–20 kb in a Megaruptor 3 system (Bates *et al.*, 2023).

Sheared DNA was purified by solid-phase reversible immobilisation, using AMPure PB beads to eliminate shorter fragments and concentrate the DNA (Strickland *et al.*, 2023). The concentration of the sheared and purified DNA was assessed using a Nanodrop spectrophotometer and Qubit Fluorometer using the Qubit dsDNA High Sensitivity Assay kit. Fragment size distribution was evaluated by running the sample on the FemtoPulse system.

## Hi-C sample preparation

Tissue of the idNemNiti1 sample was processed at the WSI Scientific Operations core, using the Arima-HiC v2 kit. Tissue (stored at –80 °C) was fixed, and the DNA crosslinked using a TC buffer with 22% formaldehyde. After crosslinking, the tissue was homogenised using the Diagnocine Power Masher-II and BioMasher-II tubes and pestles. Following the
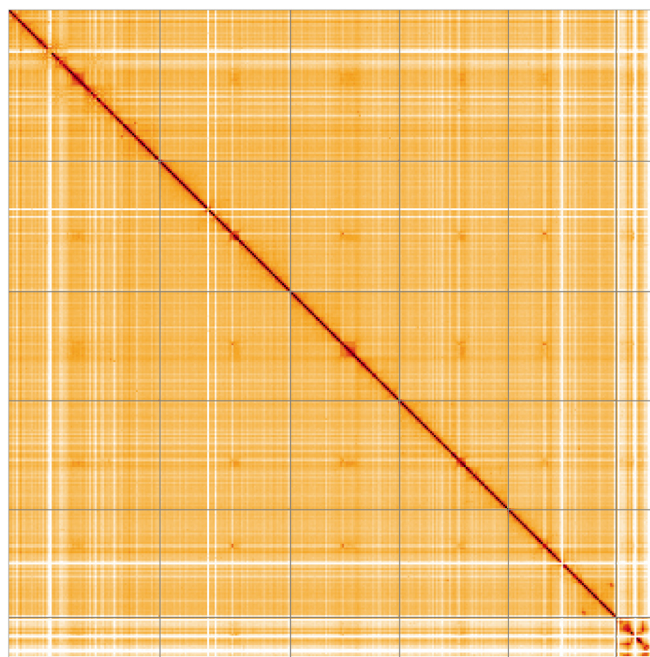
**Figure 5. Genome assembly of *Nemopoda nitidula* idNemNiti2.1: Hi-C contact map of the idNemNiti2.1 assembly, visualised using HiGlass.** Chromosomes are shown in order of size from left to right and top to bottom. An interactive version of this figure may be viewed at https://genome-note-higlass.tol.sanger.ac.uk/l/?d=KuhezG4LTbaUcf8MG8ql2Q.

**Table 3. Chromosomal pseudomolecules in the genome assembly of *Nemopoda nitidula*, idNemNiti2.**

| INSDC accession | Name | Length (Mb) | GC% |
|---|---|---|---|
| OZ077518.1 | 1 | 41.4 | 39.0 |
| OZ077519.1 | 2 | 35.64 | 38.0 |
| OZ077520.1 | 3 | 29.83 | 38.5 |
| OZ077521.1 | 4 | 29.73 | 38.5 |
| OZ077522.1 | 5 | 29.4 | 38.5 |
| OZ077523.1 | 6 | 11.05 | 36.5 |
| OZ077524.1 | MT | 0.02 | 25.0 |

kit manufacturer's instructions, crosslinked DNA was digested using a restriction enzyme master mix. The 5'-overhangs were then filled in and labelled with biotinylated nucleotides and proximally ligated. An overnight incubation was carried out for enzymes to digest remaining proteins and for crosslinks to reverse. A clean up was performed with SPRIselect beads prior to library preparation.

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

*PacBio HiFi*

The sample requires Covaris g-TUBE shearing to approximately 10 kb prior to library preparation. Ultra-low input libraries were prepared using PacBio SMRTbell® Express Template Prep Kit 2.0 and PacBio SMRTbell® gDNA Sample Amplification Kit. To begin, samples were normalised to 20 ng of DNA. Initial removal of single-strand overhangs, DNA damage repair, and end repair/A-tailing were performed per manufacturer's instructions. From the SMRTbell® gDNA Sample Amplification Kit, amplification adapters were then ligated. A 0.85X pre-PCR clean-up was performed with Promega ProNex beads and the sample was then divided into two for a dual PCR. PCR reactions A and B each followed the PCR programs as described in the manufacturer's protocol. A 0.85X post-PCR clean-up was performed with ProNex beads for PCR reactions A and B and DNA concentration was quantified using the Qubit Fluorometer v4.0 (Thermo Fisher Scientific) and Qubit HS Assay Kit and fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) and gDNA 55kb BAC analysis kit. PCR reactions A and B were then pooled, ensuring the total mass was ≥500 ng in 47.4 µl. The pooled sample then repeated the process for DNA damage repair, end repair/A-tailing and additional hairpin adapter ligation. A 1X clean-up was performed with ProNex beads and DNA concentration was quantified using the Qubit and fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies). Size selection was performed using

Sage Sciences' PippinHT system with target fragment size determined by analysis from the Femto Pulse, usually a value between 4000 and 9000 bp. Size selected libraries were then cleaned-up using1.0X ProNex beads and normalised to 2 nM before proceeding to sequencing.

Samples were sequenced on a Revio instrument (Pacific Biosciences, California, USA). Prepared libraries were normalised to 2 nM, and 15 μL was used for making complexes. Primers were annealed and polymerases were hybridised to create circularised complexes according to manufacturer's instructions. The complexes were purified with the 1.2X clean up with SMRTbell beads. The purified complexes were then diluted to the Revio loading concentration (in the range 200–300 pM), and spiked with a Revio sequencing internal control. Samples were sequenced on Revio 25M SMRT cells (Pacific Biosciences, California, USA). The SMRT link software, a PacBio web-based end-to-end workflow manager, was used to set-up and monitor the run, as well as perform primary and secondary analysis of the data upon completion.

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

## Genome assembly, curation and evaluation
### Assembly
The HiFi reads were first assembled using Hifiasm (Cheng et al., 2021) with the --primary option. Haplotypic duplications were identified and removed using purge_dups (Guan et al., 2020). The Hi-C reads were mapped to the primary contigs using bwa-mem2 (Vasimuddin et al., 2019). The contigs were further scaffolded using the provided Hi-C data (Rao et al., 2014) in YaHS (Zhou et al., 2023) using the --break option for handling potential misassemblies. The scaffolded assemblies were evaluated using Gfastats (Formenti et al., 2022), BUSCO (Manni et al., 2021) and MERQURY.FK (Rhie et al., 2020).

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

### Assembly curation
The assembly was decontaminated using the Assembly Screen for Cobionts and Contaminants (ASCC) pipeline (article in preparation). Flat files and maps used in curation were generated

in TreeVal (Pointon et al., 2023). Manual curation was primarily conducted using PretextView (Harry, 2022), with additional insights provided by JBrowse2 (Diesh et al., 2023) and HiGlass (Kerpedjiev et al., 2018). Scaffolds were visually inspected and corrected as described by Howe et al. (2021). Any identified contamination, missed joins, and mis-joins were corrected, and duplicate sequences were tagged and removed. The curation process is documented at https://gitlab.com/wtsi-grit/rapid-curation (article in preparation).

### Assembly quality assessment
The Merqury.FK tool (Rhie et al., 2020), run in a Singularity container (Kurtzer et al., 2017), was used to evaluate k-mer completeness and assembly quality for the primary and alternate haplotypes using the k-mer databases ($k = 31$) that were computed prior to genome assembly. The analysis outputs included assembly QV scores and completeness statistics.

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

The blobtoolkit pipeline is a Nextflow port of the previous Snakemake Blobtoolkit pipeline (Challis et al., 2020). It aligns the PacBio reads in SAMtools and minimap2 (Li, 2018) and generates coverage tracks for regions of fixed size. In parallel, it queries the GoaT database (Challis et al., 2023) to identify all matching BUSCO lineages to run BUSCO (Manni et al., 2021). For the three domain-level BUSCO lineages, the pipeline aligns the BUSCO genes to the UniProt Reference Proteomes database (Bateman et al., 2023) with DIAMOND blastp (Buchfink et al., 2021). The genome is also divided into chunks according to the density of the BUSCO genes from the closest taxonomic lineage, and each chunk is aligned to the UniProt Reference Proteomes database using DIAMOND blastx. Genome sequences without a hit are chunked using seqtk and aligned to the NT database with blastn (Altschul et al., 1990). The blobtools suite combines all these 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), relying on the Conda package manager, the Bioconda initiative (Grüning et al., 2018), the Biocontainers infrastructure (da Veiga Leprevost et al., 2017), as well as the Docker (Merkel, 2014) and Singularity (Kurtzer et al., 2017) containerisation solutions.

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

## 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

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

| Software tool | 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 | 4.3.7 | https://github.com/blobtoolkit/blobtoolkit |
| BUSCO | 5.4.3 and 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 | 427104ea91c78c3b8b8b49f1a7d6bbeaa869ba1c | https://github.com/thegenemyers/FASTK |
| 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-r587 | https://github.com/chhylp123/hifiasm |
| HiGlass | 44086069ee7d4d3f6f3f0012569789ec138f42b84aa44357826c0b6753eb28de | https://github.com/higlass/higlass |
| Merqury.FK | d00d98157618f4e8d1a9190026b19b471055b22e | https://github.com/thegenemyers/MERQURY.FK |
| MitoHiFi | 3 | https://github.com/marcelauliano/MitoHiFi |
| MultiQC | 1.14, 1.17, and 1.18 | https://github.com/MultiQC/MultiQC |
| NCBI Datasets | 15.12.0 | https://github.com/ncbi/datasets |
| Nextflow | 23.04.0-5857 | https://github.com/nextflow-io/nextflow |
| PretextView | 0.2.5 | https://github.com/sanger-tol/PretextView |
| purge_dups | 1.2.5 | https://github.com/dfguan/purge_dups |
| samtools | 1.16.1, 1.17, and 1.18 | https://github.com/samtools/samtools |
| sanger-tol/ascc | - | https://github.com/sanger-tol/ascc |
| Seqtk | 1.3 | https://github.com/lh3/seqtk |
| Singularity | 3.9.0 | https://github.com/sylabs/singularity |
| TreeVal | 1.0.0 | https://github.com/sanger-tol/treeval |
| YaHS | 1.2a.2 | https://github.com/c-zhou/yahs |

of materials by a Darwin Tree of Life Partner is subject to the **'Darwin Tree of Life Project Sampling Code of Practice'**, which can be found in full on the Darwin Tree of Life website here. By agreeing with and signing up to the Sampling 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: Nemopoda nitidula. Accession number PRJEB74593; https://identifiers.org/ena.embl/PRJEB74593. The genome sequence is released openly for reuse. The *Nemopoda nitidula* genome sequencing initiative is part of the Darwin Tree of Life (DToL) project. All raw sequence data and the assembly have been deposited in INSDC databases. 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.

## Author information

Members of the University of Oxford and Wytham Woods Genome Acquisition Lab are listed here: https://doi.org/10.5281/zenodo.12157525.

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

Members of the Wellcome Sanger Institute Tree of Life Management, Samples and Laboratory team are listed here: https://doi.org/10.5281/zenodo.12162482.

Members of Wellcome Sanger Institute Scientific Operations: Sequencing Operations are listed here: https://doi.org/10.5281/zenodo.12165051.

Members of the Wellcome Sanger Institute Tree of Life Core Informatics team are listed here: https://doi.org/10.5281/zenodo.12160324.

Members of the Tree of Life Core Informatics collective are listed here: https://doi.org/10.5281/zenodo.12205391.

Members of the Darwin Tree of Life Consortium are listed here: https://doi.org/10.5281/zenodo.4783558.

## References

Abdennur N, Mirny LA: **Cooler: scalable storage for Hi-C data and other genomically labeled arrays.** *Bioinformatics.* 2020; **36**(1): 311–316.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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**

Bates A, Clayton-Lucey I, Howard C: **Sanger Tree of Life HMW DNA fragmentation: diagenode Megaruptor®3 for LI PacBio.** *protocols.io.* 2023.
**Publisher Full Text**

Beasley J, Uhl R, Forrest LL, *et al.*: **DNA barcoding SOPs for the Darwin Tree of Life project.** *protocols.io.* 2023; [Accessed 25 June 2024].
**Publisher Full Text**

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**

Crowley L, Allen H, Barnes I, *et al.*: **A sampling strategy for genome sequencing the British terrestrial arthropod fauna [version 1; peer review: 2 approved].** *Wellcome Open Res.* 2023; **8**: 123.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

da Veiga Leprevost F, Grüning BA, Alves Aflitos S, *et al.*: **BioContainers: an open-source and community-driven framework for software standardization.** *Bioinformatics.* 2017; **33**(16): 2580–2582.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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**

Denton A, Oatley G, Cornwell C, *et al.*: **Sanger Tree of Life sample homogenisation: PowerMash.** *protocols.io.* 2023a.
**Publisher Full Text**

Denton A, Yatsenko H, Jay J, *et al.*: **Sanger Tree of Life wet laboratory protocol collection V.1.** *protocols.io.* 2023b.
**Publisher Full Text**

Diesh C, Stevens GJ, Xie P, *et al.*: **JBrowse 2: a modular genome browser with views of synteny and structural variation.** *Genome Biol.* 2023; **24**(1): 74.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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**

Frantsevich L, Gladun D: **Evolution of the middle leg basal articulations in flies (Diptera).** *Acta Zool.* 2002; **83**(2): 125–147.
**Publisher Full Text**

GBIF Secretariat: ***Nemopoda nitidula* (Fallén, 1820).** GBIF Backbone Taxonomy, 2024.
**Reference Source**

Grüning B, Dale R, Sjödin A, *et al.*: **Bioconda: sustainable and comprehensive software distribution for the life sciences.** *Nat Methods.* 2018; **15**(7): 475–476.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Guan D, McCarthy SA, Wood J, *et al.*: **Identifying and removing haplotypic duplication in primary genome assemblies.** *Bioinformatics.* 2020; **36**(9): 2896–2898.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Harry E: **PretextView (Paired REad TEXTure Viewer): a desktop application for viewing pretext contact maps.** 2022.
**Reference Source**

Hosken DJ, Martin OY, Born J, *et al.*: **Sexual conflict in *Sepsis cynipsea*: female reluctance, fertility and mate choice.** *J Evol Biol.* 2003; **16**(3): 485–490.
**PubMed Abstract** | **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**

Jay J, Yatsenko H, Narváez-Gómez JP, *et al.*: **Sanger Tree of Life sample preparation: triage and dissection.** *protocols.io.* 2023.
**Publisher Full Text**

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

Khaghaninia S, Zarghani E, Gharajedaghi Y: **A recent contribution to the black scavenger flies (Diptera: Sepsidae) in Iran.** *Biharean Biol.* 2014; **8**(1): 21–23.
**Reference Source**

Kostygov AY, Frolov AO, Malysheva MN, *et al.*: *Vickermania* **gen. nov., trypanosomatids that use two joined flagella to resist midgut peristaltic flow within the fly host.** *BMC Biol.* 2020; **18**(1): 187.
**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**

Meier R: **Larval morphology of the Sepsidae (Diptera: Sciomyzoidea), with a cladistic analysis using adult and larval characters.** *Bull Am Mus Nat Hist.* 1996; **228**: 1–147.
**Reference Source**

Merkel D: **Docker: lightweight Linux containers for consistent development and deployment.** *Linux J.* 2014; **2014**(239): 2, [Accessed 2 April 2024].
**Reference Source**

Oatley G, Denton A, Howard C: **Sanger Tree of Life HMW DNA extraction: automated MagAttract v.2.** *protocols.io.* 2023.
**Publisher Full Text**

Oleksakova T, Klimesova V, Bartak M, *et al.*: **Sepsidae (Diptera) associated with animal and human decomposition in the Czech Republic.** *Mendelnet.* 2016; **23**: 795–800.
**Reference Source**

Pointon DL, Eagles W, Sims Y, *et al.*: **sanger-tol/treeval v1.0.0 – Ancient Atlantis.** 2023.
**Publisher Full Text**

Pont AC, Meier R: **The Sepsidae (Diptera) of Europe.** Brill, **2002; 37**.
**Reference Source**

Prado e Castro C, Serrano A, da Silva P, *et al.*: **Carrion flies of forensic interest: a study of seasonal community composition and succession in Lisbon, Portugal.** *Med Vet Entomol.* 2012; **26**(4): 417–431.
**PubMed Abstract** | **Publisher Full Text**

Puniamoorthy N, Su KF, Meier R: **Bending for love: losses and gains of sexual dimorphisms are strictly correlated with changes in the mounting position of sepsid flies (Sepsidae: Diptera).** *BMC Evol Biol.* 2008; **8**: 155.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Quinlan AR, Hall IM: **BEDTools: a flexible suite of utilities for comparing genomic features.** *Bioinformatics.* 2010; **26**(6): 841–842.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Rao SSP, Huntley MH, Durand NC, *et al.*: **A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping.** *Cell.* 2014; **159**(7): 1665–1680.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Rhie A, McCarthy SA, Fedrigo O, *et al.*: **Towards complete and error-free genome assemblies of all vertebrate species.** *Nature.* 2021; **592**(7856): 737–746.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Rhie A, Walenz BP, Koren S, *et al.*: **Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies.** *Genome Biol.* 2020; **21**(1): 245.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Rohner PT, Haenni JP, Giesen A, *et al.*: **Temporal niche partitioning of Swiss black scavenger flies in relation to season and substrate age (Diptera, Sepsidae).** *Alpine Entomology.* 2019; **3**(1–2): 1–10.
**Publisher Full Text**

Strickland M, Cornwell C, Howard C: **Sanger Tree of Life fragmented DNA clean up: manual SPRI.** *protocols.io.* 2023.
**Publisher Full Text**

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**

Vicoso B, Bachtrog D: **Numerous transitions of sex chromosomes in Diptera.** *PLoS Biol.* 2015; **13**(4): e1002078.
**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: ✔ ✔

---

✔ **Thomas Brown** (ID)

Leibniz IZW (Berlin) & ERGA/BGE, Leibniz, Germany

Falk et al present the genome assembly for the black scavenger fly Nemopoda nitudula as part of the DToL program. The methods used are state-of-the-art and well explained in general. A little more explanation could be given as to what features about the phylogeny and evolution of "different traits" that will be uncovered with the genome would give the background section more weight.

I only have a couple of minor points:

The abstract states the assembly is from a single specimen, yet the PacBio and Hi-C reads are stated as coming from different specimens in the methods. I find this a little misleading.

The kmer completeness is stated as 100% complete in the table. I think it should be made clearer that this corresponds to the combination of primary and alternate assemblies, as these are each stated as having completenesses of 76% in the text. The contig and scaffold N50s and BUSCO scores most likely correspond to the primary assembly, I find this misleading

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

Partly

**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:* Genome assembly, Genome annotation, Workflows, FAIR research

**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 06 June 2025

https://doi.org/10.21956/wellcomeopenres.26106.r122339

**Terrence Sylvester** iD

University of memphis, memphis, USA

The authors provide a high-quality genome assembly of the black scavenger fly, *Nemopoda nitidula.* Genome quality matrices indicate that the presented genome was assembled and curated to minimise assembly errors. This genome would be a valuable community resource for other researchers. If the authors can include the metrics of the successful extraction that was used for PacBio sequencing (e.g., Nanodrop graph, femtopulse fragment size distribution, etc.), that would be something of value for other researchers, giving a point of reference to compare their projects with. Other than that, I do not have significant remarks on the paper.

**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:* 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.**