Check for updates

DATA NOTE

# The genome sequence of the Heath Bumblebee, *Bombus jonellus* (Kirby, 1802)

[version 1; peer review: 2 approved]

Gavin R. Broad [1], Inez Januszczak [1], Chris Fletcher [1],
Natural History Museum 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]Natural History Museum, London, England, UK

## Abstract

We present a genome assembly from a female specimen of *Bombus jonellus* (Heath Bumblebee; Arthropoda; Insecta; Hymenoptera; Apidae). The genome sequence has a total length of 357.90 megabases. Most of the assembly (78.06%) is scaffolded into 18 chromosomal pseudomolecules. The mitochondrial genome has also been assembled, with a length of 24.83 kilobases.

## Keywords

Bombus jonellus, Heath Bumblebee, genome sequence, chromosomal, Hymenoptera

This article is included in the Tree of Life gateway.

**Open Peer Review**

**Approval Status** ✓ ✓

|  | 1 | 2 |
| --- | --- | --- |
| **version 1**<br>23 May 2025 | ✓<br>view | ✓<br>view |

1. **Jason Charamis** [iD], Foundation for Research and Technology - Hellas, Irákleion, Greece

2. **Erich D Jarvis**, Rockefeller University, Millbrook, USA

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: Broad GR**: Investigation, Resources, Writing – Original Draft Preparation, Writing – Review & Editing; **Januszczak I**: Investigation, Resources; **Fletcher C**: Investigation, Resources;

**How to cite this article:** Broad GR, Januszczak I, Fletcher C *et al.* **The genome sequence of the Heath Bumblebee, *Bombus jonellus* (Kirby, 1802) [version 1; peer review: 2 approved]** Wellcome Open Research 2025, **10**:269 https://doi.org/10.12688/wellcomeopenres.24256.1

**First published:** 23 May 2025, **10**:269 https://doi.org/10.12688/wellcomeopenres.24256.1

## Species taxonomy

Eukaryota; Opisthokonta; Metazoa; Eumetazoa; Bilateria; Protostomia; Ecdysozoa; Panarthropoda; Arthropoda; Mandibulata; Pancrustacea; Hexapoda; Insecta; Dicondylia; Pterygota; Neoptera; Endopterygota; Hymenoptera; Apocrita; Aculeata; Apoidea; Anthophila; Apidae; Apinae; Bombini; *Bombus*; *Pyrobombus*; *Bombus jonellus* (Kirby, 1802) (NCBI: txid85663)

## Background

*Bombus jonellus*, the Heath Bumblebee, lives up to its name by inhabiting mainly heaths in southern England, but across Britain has a wider range of open habitats, often inhabiting moorland and machair in Scotland and various coastal habitats. It is one of the short-faced species with two yellow bands on the thorax and a white 'tail' (tergites four and five), often of similar size to *B. pratorum* but in *B. pratorum* the 'tail' is orange. Note that in the northern Isles of Scotland, the 'tail' is yellow and in males from Orkney and the western Isles, orange, thus easily confused with *B. pratorum*, should that species colonise the northern Isles. The ubiquitous *Bombus hortorum* is superficially similar but usually larger, brighter yellow, shorter-haired and with a longer face. As with many bumblebees, the males of *B. jonellus* have a yellow-haired face while females are black-haired. There are several recent identification keys which include *B. jonellus* (e.g., Benton, 2009; Else & Edwards, 2018; Falk, 2015). Benton (2009) and Else & Edwards (2018) summarise the ecology of *B. jonellus*.

This is typically a bivoltine species but in parts of its range there may only be one generation per year. Nests seem to be established in a variety of locations, on, below and sometimes far above the ground and a wide range of flowers are visited, with Fabaceae reported to be the most important for pollen (Else & Edwards, 2018). Colonies are usually small, from 30 up to 120 workers (Benton, 2009; Else & Edwards, 2018) and queens can be found from March through to September.

Within Britain, Williams (Williams, 1989a) classified *B. jonellus* as a 'widespread' local species, although Benton (2009) suggests that it is less local, more widespread than usually reported. This might not be true of southern England, as this species is disappearing from Dungeness, Kent, which has held on to a rich range of *Bombus* species compared to the increasingly impoverished surrounding areas, probably because of the continued high density of appropriate flowers (Williams, 1989b). Else & Edwards (2018) map a distribution heavily skewed towards northwest/northern Scotland and southeast and southwest England. Throughout central England, *B. jonellus* is virtually absent and is of very scattered distribution in Wales and northern England. GRB was pleased, very early in his entomological career, to find the first *B. jonellus* for Cheshire. The specimen which had its genome sequenced (Figure 1) was collected in Heather-rich moorland in Beinn Eighe NNR, where the species was common.



**Figure 1. Photograph of the *Bombus jonellus* (iyBomJone1) specimen used for genome sequencing.**

## Genome sequence report

### Sequencing data

The genome of a specimen of *Bombus jonellus* (Figure 1) was sequenced using Pacific Biosciences single-molecule HiFi long reads, generating 27.20 Gb from 2.58 million reads, which were used to assemble the genome. GenomeScope analysis estimated the haploid genome size at 401.80 Mb, with a heterozygosity of 0.38% and repeat content of 42.03%. These estimates guided expectations for the assembly. Based on the estimated genome size, the sequencing data provided approximately 65 coverage. Hi-C sequencing produced 139.13 Gb from 921.40 million reads, and was used to scaffold the assembly. Table 1 summarises the specimen and sequencing details.

### 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 30 misjoins or missing joins. These interventions decreased the scaffold count by 3.46% and increased the scaffold N50 by 4.58%. The final assembly has a total length of 357.90 Mb in 250 scaffolds, with 118 gaps, and a scaffold N50 of 14.15 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 (78.06%) was assigned to 18 chromosomal-level scaffolds. These chromosome-level scaffolds,

**Table 1. Specimen and sequencing data for *Bombus jonellus*.**

| Project information | | | |
|---|---|---|---|
| **Study title** | Bombus jonellus | | |
| **Umbrella BioProject** | PRJEB64937 | | |
| **Species** | *Bombus jonellus* | | |
| **BioSpecimen** | SAMEA14448317 | | |
| **NCBI taxonomy ID** | 85663 | | |
| **Specimen information** | | | |
| **Technology** | **ToLID** | **BioSample accession** | **Organism part** |
| **PacBio long read sequencing** | iyBomJone1 | SAMEA14448516 | thorax |
| **Hi-C sequencing** | iyBomJone1 | SAMEA14448517 | head |
| **Sequencing information** | | | |
| **Platform** | **Run accession** | **Read count** | **Base count (Gb)** |
| **Hi-C Illumina NovaSeq 6000** | ERR11837502 | 9.21e+08 | 139.13 |
| **PacBio Sequel IIe** | ERR11843415 | 2.58e+06 | 27.2 |

**Table 2. Genome assembly data for *Bombus jonellus*.**

| Genome assembly | | |
|---|---|---|
| Assembly name | iyBomJone1.1 | |
| Assembly accession | GCA_964197665.1 | |
| *Alternate haplotype accession* | *GCA_964197655.1* | |
| Assembly level for primary assembly | chromosome | |
| Span (Mb) | 357.90 | |
| Number of contigs | 368 | |
| Number of scaffolds | 250 | |
| Longest scaffold (Mb) | 22.69 | |
| **Assembly metric** | **Measure** | ***Benchmark*** |
| Contig N50 length | 2.76 Mb | *≥ 1 Mb* |
| Scaffold N50 length | 14.15 Mb | *= chromosome N50* |
| Consensus quality (QV) | Primary: 62.7; alternate: 63.9; combined: 63.5 | *≥ 40* |
| *k*-mer completeness | Primary: 88.02%; alternate: 87.17%; combined: 97.56% | *≥ 95%* |
| BUSCO* | C:97.6%[S:97.3%,D:0.3%], F:0.4%,M:2.1%,n:5,991 | *S > 90%; D < 5%* |
| Percentage of assembly assigned to chromosomes | 78.06% | *≥ 90%* |
| Sex chromosomes | None | *localised homologous pairs* |
| Organelles | Mitochondrial genome: 24.83 kb | *complete single alleles* |

* BUSCO scores based on the hymenoptera_odb10 BUSCO set using version 5.5.0. C = complete [S = single copy, D = duplicated], F = fragmented, M = missing, n = number of orthologues in comparison.
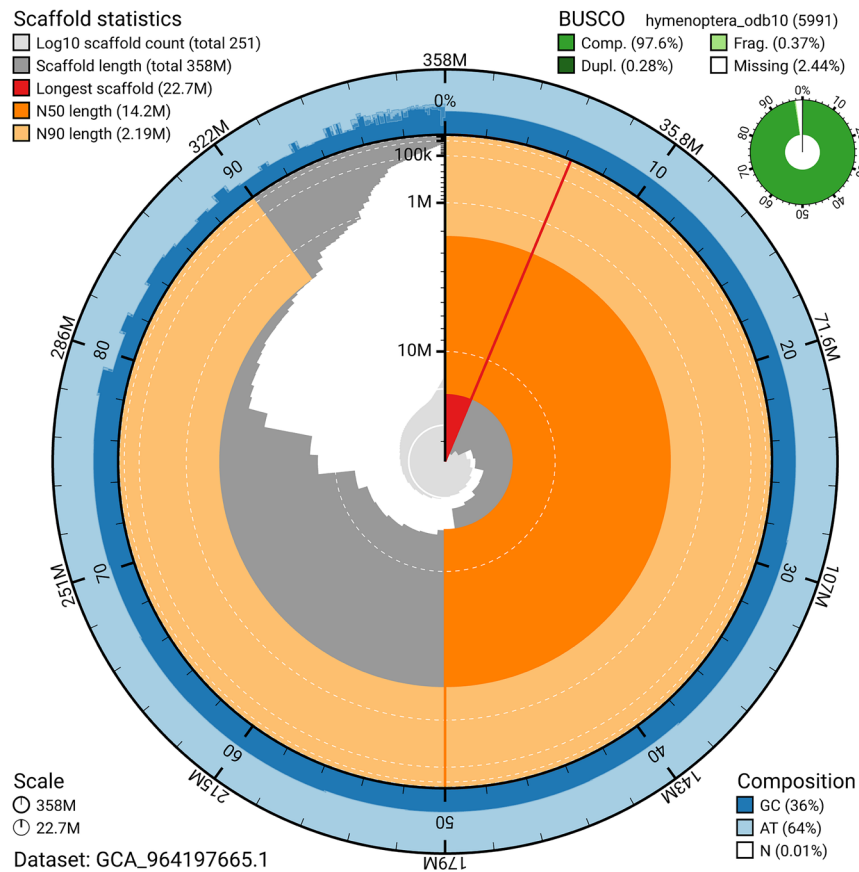
**Figure 2. Genome assembly of *Bombus jonellus*, iyBomJone1.1: metrics.** The BlobToolKit snail plot provides an overview of assembly metrics and BUSCO gene completeness. The circumference represents the length of the whole genome sequence, and the main plot is divided into 1,000 bins around the circumference. 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 set is presented at the top right. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/GCA_964197665.1/dataset/GCA_964197665.1/snail.

confirmed by Hi-C data, are named according to size (Figure 5; Table 3).

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.

## 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 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 combined primary and alternate assemblies achieve an estimated QV of 63.5. The *k*-mer recovery for the primary

haplotype is 88.02%, and for the alternate haplotype 87.17%; the combined primary and alternate assemblies have a *k*-mer recovery of 97.56%. BUSCO v.5.5.0 analysis using the hymenoptera_odb10 reference set (*n* = 5,991) identified 97.6% of the expected gene set (single = 97.3%, duplicated = 0.3%).

Table 2 provides assembly metric benchmarks adapted from Rhie *et al.* (2021) and the Earth BioGenome Project Report on Assembly Standards September 2024. The assembly achieves the EBP reference standard of **6.7.Q62**.

## Methods

### Sample acquisition and DNA barcoding

The specimen used for genome sequencing was an adult female *Bombus jonellus* (specimen ID NHMUK014451607, ToLID iyBomJone1), collected from Beinn Eighe National Nature Reserve National Nature Reserve, Scotland, United Kingdom, Scotland, United Kingdom (latitude 57.63,
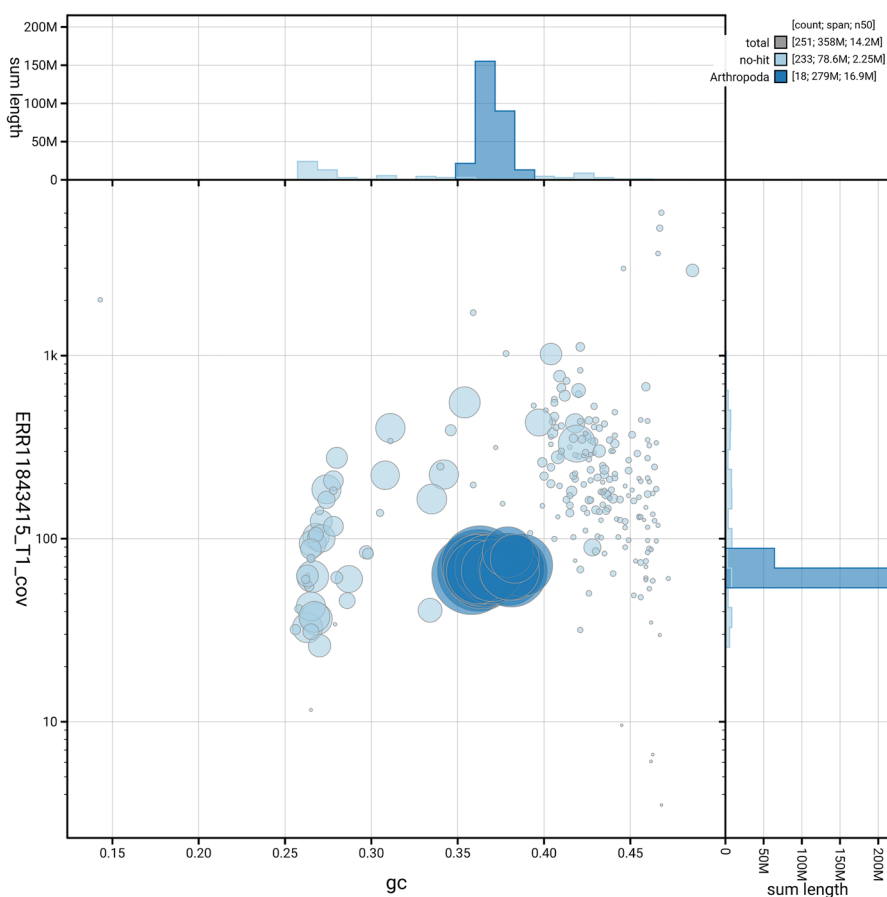
**Figure 3. Genome assembly of *Bombus jonellus*, iyBomJone1.1: BlobToolKit GC-coverage plot.** Blob plot showing sequence coverage (vertical axis) and GC content (horizontal axis). The circles represent scaffolds, with the size proportional to scaffold length and the colour representing phylum membership. 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 at https://blobtoolkit.genomehubs.org/view/GCA_964197665.1/dataset/GCA_964197665.1/blob.

longitude −5.35) on 2021-09-10, using an aerial net. The specimen was collected by Gavin Broad, David Lees, Inez Januszczak and Chris Fletcher (Natural History Museum), identified by Gavin Broad (Natural History Museum) and preserved by dry freezing (−80 °C).

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) (Pereira *et al.*, 2022). 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).

Metadata collection for samples adhered to the Darwin Tree of Life project standards described by Lawniczak *et al.* (2022).

## 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 iyBomJone1 sample was prepared for DNA extraction by weighing and dissecting it on dry ice (Jay *et al.*, 2023). Tissue from the thorax was homogenised using a PowerMasher II tissue disruptor (Denton *et al.*, 2023a).
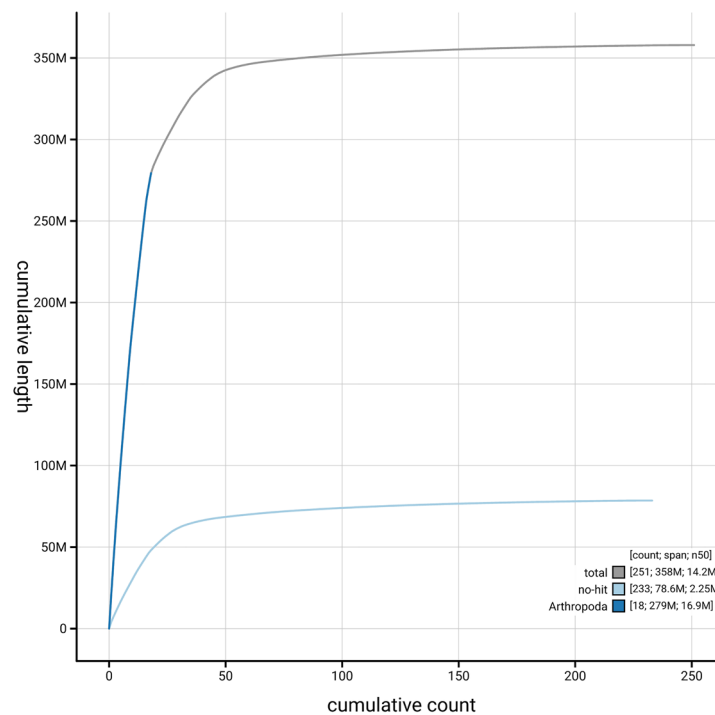
**Figure 4. Genome assembly of *Bombus jonellus*, iyBomJone1.1: BlobToolKit cumulative sequence plot.** The grey line shows cumulative length for all scaffolds. Coloured lines show cumulative lengths of scaffolds assigned to each phylum using the buscogenes taxrule. An interactive version of this figure is available at https://blobtoolkit.genomehubs.org/view/GCA_964197665.1/dataset/GCA_964197665.1/cumulative.
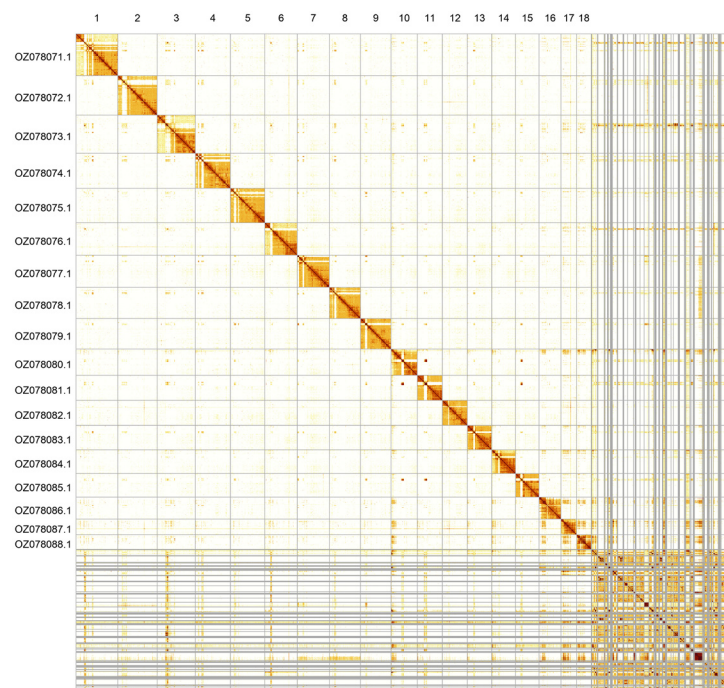


**Figure 5. Genome assembly of *Bombus jonellus*: Hi-C contact map of the iyBomJone1.1 assembly, visualised using PretextSnapshot.** Chromosomes are shown in order of size from left to right and top to bottom. An interactive version of this figure may be viewed at https://genome-note-higlass.tol.sanger.ac.uk/l/?d=PUlRBbxQR3COirbclnKaPA.

**Table 3. Chromosomal pseudomolecules in the genome assembly of *Bombus jonellus*, iyBomJone1.**

| INSDC accession | Name | Length (Mb) | GC% |
|---|---|---|---|
| OZ078071.1 | 1 | 22.69 | 36 |
| OZ078072.1 | 2 | 21.47 | 36 |
| OZ078073.1 | 3 | 20.71 | 36.5 |
| OZ078074.1 | 4 | 18.96 | 36.5 |
| OZ078075.1 | 5 | 18.68 | 36.5 |
| OZ078076.1 | 6 | 17.54 | 38 |
| OZ078077.1 | 7 | 17.45 | 38 |
| OZ078078.1 | 8 | 16.88 | 36 |
| OZ078079.1 | 9 | 16.57 | 37 |
| OZ078080.1 | 10 | 14.15 | 36.5 |
| OZ078081.1 | 11 | 13.64 | 36 |
| OZ078082.1 | 12 | 13.53 | 38 |
| OZ078083.1 | 13 | 13.19 | 38 |
| OZ078084.1 | 14 | 12.93 | 38.5 |
| OZ078085.1 | 15 | 12.79 | 37 |
| OZ078086.1 | 16 | 11.87 | 38 |
| OZ078087.1 | 17 | 8.49 | 38 |
| OZ078088.1 | 18 | 7.82 | 38.5 |
| OZ078089.1 | MT | 0.02 | 14.5 |

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 and crosslinking

Hi-C data were generated from the head of the iyBomJone1 sample using the Arima-HiC v2 kit (Arima Genomics) with 20–50 mg of frozen tissue (stored at –80 °C). As per manufacturer's instructions, tissue was fixed, and the DNA crosslinked using a TC buffer with 22% formaldehyde concentration, and a final formaldehyde concentration of 2%. The tissue was then homogenised using the Diagnocine Power Masher-II. The crosslinked DNA was digested using a restriction enzyme master mix, then biotinylated and ligated. A clean up was performed with SPRIselect beads prior to library preparation. DNA concentration was quantified using the Qubit Fluorometer v4.0 (Thermo Fisher Scientific) and Qubit HS Assay Kit, and sample biotinylation percentage was estimated using the Arima-HiC v2 QC beads.

## Library preparation and sequencing

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

### PacBio HiFi

Samples need to have an average fragment size exceeding 8 kb and a total mass over 400 ng to proceed to the low-input SMRTbell Prep Kit 3.0 protocol (Pacific Biosciences), depending on genome size and sequencing depth required. Libraries were prepared using the SMRTbell Prep Kit 3.0 as per the manufacturer's instructions. The kit includes the reagents required for end repair/A-tailing, adapter ligation, post-ligation SMRTbell bead cleanup, and nuclease treatment. Size-selection and clean-up were carried out using diluted AMPure PB beads (Pacific Biosciences). DNA concentration was quantified using the Qubit Fluorometer v4.0 (ThermoFisher Scientific) with Qubit 1X dsDNA HS assay kit and the final library fragment size analysis was carried out using the Agilent Femto Pulse Automated Pulsed Field CE Instrument (Agilent Technologies) and the gDNA 55kb BAC analysis kit.

Samples were sequenced using the Sequel IIe system (Pacific Biosciences, California, USA). The concentration of the library loaded onto the Sequel IIe was in the range 40–135 pM. 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, the biotinylated DNA constructs were fragmented using a Covaris E220 sonicator and size-selected to 400–600 bp using SPRISelect beads. DNA was then enriched using Arima-HiC v2 Enrichment beads. The NEBNext Ultra II DNA Library Prep Kit (New England Biolabs) was used for end repair, A-tailing, and adapter ligation, following a modified protocol in which library preparation is carried out while the DNA remains bound to the enrichment beads. PCR amplification was performed using KAPA HiFi HotStart mix and custom dual-indexed adapters (Integrated DNA Technologies) in a 96-well plate format. Depending on sample concentration and biotinylation percentage determined at the crosslinking stage, samples were amplified for 10–16 PCR cycles. Post-PCR clean-up was carried out using SPRISelect beads. The libraries were quantified using the Accuclear Ultra High Sensitivity dsDNA Standards Assay kit (Biotium) and normalised to 10 ng/μL before sequencing. Hi-C sequencing was performed on the Illumina NovaSeq 6000 instrument using 150 bp paired-end reads.

## Genome assembly, curation and evaluation
### 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 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 (Rao *et al.*, 2014) were mapped to the primary contigs using bwa-mem2 (Vasimuddin *et al.*, 2019), and the contigs were scaffolded using 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. Flat files and maps used in curation were generated via the TreeVal pipeline (Pointon *et al.*, 2023). Manual curation was conducted primarily in PretextView (Harry, 2022) and HiGlass (Kerpedjiev *et al.*, 2018), with additional insights provided by JBrowse2 (Diesh *et al.*, 2023). Scaffolds were visually inspected and corrected as described by Howe *et al.* (2021). Any identified contamination, missed joins, and mis-joins were amended, and duplicate sequences were tagged and removed. 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), 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) computed prior to genome assembly. The analysis outputs included assembly QV scores and completeness statistics.

The genome was analysed in the blobtoolkit pipeline, a Nextflow (Di Tommaso *et al.*, 2017) port of the previous Snakemake Blobtoolkit pipeline (Challis *et al.*, 2020). It aligns the PacBio reads in SAMtools (Danecek *et al.*, 2021) 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 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.

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

| Software tool | Version | Source |
|---|---|---|
| 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 |
| DIAMOND | 2.1.8 | https://github.com/bbuchfink/diamond |
| fasta_windows | 0.2.4 | https://github.com/tolkit/fasta_windows |
| FastK | 666652151335353eef2fcd58880bcef5bc2928e1 | https://github.com/thegenemyers/FASTK |
| GenomeScope2.0 | 2.0.1 | https://github.com/tbenavi1/genomescope2.0 |

| Software tool | Version | Source |
|---|---|---|
| Gfastats | 1.3.6 | https://github.com/vgl-hub/gfastats |
| GoaT CLI | 0.2.5 | https://github.com/genomehubs/goat-cli |
| Hifiasm | 0.16.1 | https://github.com/chhylp123/hifiasm |
| HiGlass | 44086069ee7d4d3f6f3f0012569789ec138f42b84aa44357826c0b6753eb28de | https://github.com/higlass/higlass |
| MerquryFK | d00d98157618f4e8d1a9190026b19b471055b22e | https://github.com/thegenemyers/MERQURY.FK |
| Minimap2 | 2.24-r1122 | https://github.com/lh3/minimap2 |
| MitoHiFi | None | https://github.com/marcelauliano/MitoHiFi |
| MultiQC | 1.14, 1.17, and 1.18 | https://github.com/MultiQC/MultiQC |
| Nextflow | 23.10.0 | https://github.com/nextflow-io/nextflow |
| PretextView | 0.2.5 | https://github.com/sanger-tol/PretextView |
| purge_dups | 1.2.3 | https://github.com/dfguan/purge_dups |
| 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 | 1.1a.2 | https://github.com/c-zhou/yahs |

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: Bombus jonellus. Accession number PRJEB64937; https://identifiers.org/ena.embl/PRJEB64937. The genome sequence is released openly for reuse. The *Bombus jonellus* genome sequencing initiative is part of the Darwin Tree of Life Project (PRJEB40665) and the Sanger Institute Tree of Life Programme (PRJEB43745). 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 Natural History Museum Genome Acquisition Lab are listed here: https://doi.org/10.5281/zenodo.12159242.

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

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**

Benton T: **Bumblebees**. HarperCollins, 2009.

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**

Di Tommaso P, Chatzou M, Floden EW, *et al.*: **Nextflow enables reproducible computational workflows.** *Nat Biotechnol.* 2017; **35**(4): 316–319.
**PubMed Abstract** | **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**

Else GR, Edwards M: **Handbook of the bees of the British Isles.** London: The Ray Society, 2018; **2**.
**Reference Source**

Ewels P, Magnusson M, Lundin S, *et al.*: **MultiQC: summarize analysis results for multiple tools and samples in a single report.** *Bioinformatics.* 2016; **32**(19): 3047–3048.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

Ewels PA, Peltzer A, Fillinger S, *et al.*: **The nf-core framework for community-curated bioinformatics pipelines.** *Nat Biotechnol.* 2020; **38**(3): 276–278.
**PubMed Abstract** | **Publisher Full Text**

Falk S: **Field guide to the bees of Britain and Ireland.** Bloomsbury Wildlife Guides, 2015.
**Reference Source**

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**

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

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

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

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

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

Kurtzer GM, Sochat V, Bauer MW: **Singularity: scientific containers for mobility of compute.** *PLoS One.* 2017; **12**(5): e0177459.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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, [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**

Pereira L, Sivell O, Sivess L, *et al.*: **DToL Taxon-specific Standard Operating Procedure for the terrestrial and freshwater arthropods working group**. 2022.
**Publisher Full Text**

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

Ranallo-Benavidez TR, Jaron KS, Schatz MC: **GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes.** *Nat Commun.* 2020; **11**(1): 1432.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

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

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

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

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**

Williams PH: **Bumble bees – and their decline in Britain.** Ilford: Central Association of Bee-Keepers, 1989a; [Accessed 14 April 2025].
**Reference Source**

Williams PH: **Why are there so many species of bumble bees at Dungeness?** *Bot J Linn Soc.* 1989b; **101**(1): 31–44.
**Publisher Full Text**

Zhou C, McCarthy SA, Durbin R: **YaHS: yet another Hi-C scaffolding tool.** *Bioinformatics.* 2023; **39**(1): btac808.
**PubMed Abstract** | **Publisher Full Text** | **Free Full Text**

# Open Peer Review

## Current Peer Review Status: ✓ ✓

---

**Version 1**

✓ **Erich D Jarvis**

Rockefeller University, Millbrook, USA

This is a good genomic data note. It is well written, explained well, with good figures. The genome assembly is good. Not sure why the authors did not use the Hi-C phasing mode to get two relatively complete haplotypes. In future assemblies, it would be good to do so. No changes though are necessary for the current Data Note

**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 and Neuroscience

**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.**

✓ **Jason Charamis** [iD]

Foundation for Research and Technology - Hellas, Irákleion, Greece

This data note presents a high-quality genome assembly of the Heath Bumblebee (Bombus jonellus), contributing valuable genomic resources to the Darwin Tree of Life project. The technical work is solid and the manuscript is generally well-written.

The assembly achieves good quality metrics with 357.90 Mb total length and 78.06% scaffolded into 18 chromosomes, while it also has a high BUSCO completeness score (97.6%) with low duplication rate.

I have only two comments to make:

1. The introduction would benefit from brief discussion of why the B. jonellus genome is scientifically important beyond its contribution to the Darwin Tree of Life project.

2. The manuscript does not discuss broader implications of this genomic resource in ecological studies.

**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:* Arthropod Comparative 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.**