

First chromosome-level genome assembly of a ribbon worm from the Hoplonemertea clade, *Emplectonema gracile*, and its structural annotation

Alberto Valero-Gracia^{1*}, Nickellaus G. Roberts², Meghan Yap-Chiongco², Ana Teresa Capucho¹, Kevin M. Kocot^{2,3}, Michael Matschiner¹, Torsten H. Struck¹.

¹ Natural History Museum, University of Oslo, Blindern, P.O. Box 1172, 0318 Oslo, Norway

² Department of Biological Sciences, University of Alabama, Tuscaloosa, AL 35487, USA

³ Alabama Museum of Natural History, University of Alabama, Tuscaloosa, AL 35487, USA

* Correspondence: Alberto Valero-Gracia, Natural History Museum, University of Oslo, Norway, alberto.valero-gracia@nhm.uio.no

Abstract

Genome-wide information has so far been unavailable for ribbon worms of the clade Hoplonemertea, the most species-rich class within the phylum Nemertea. While species within Pilidiophora, the sister clade of Hoplonemertea, possess a pilidium larval stage and lack stylets on their proboscis, Hoplonemertea species have a planuliform larva and are armed with stylets employed for the injection of toxins into their prey. To further compare these developmental, physiological, and behavioral differences from a genomic perspective, the availability of a reference genome of a Hoplonemertea species is crucial. To this end, we herein present the annotated chromosome-level genome assembly for *Emplectonema gracile* (Nemertea; Hoplonemertea; Monostilifera; Emplectonematidae), an easily collected nemertean well-suited for laboratory experimentation. The genome is 157.9 Mbp in span. Hi-C scaffolding yielded 15 putative chromosomes with a scaffold N50 of 10.0 Mbp and a BUSCO completeness score of

95.3%. Structural annotation predicted 20,684 protein-coding genes. The high-quality reference genome reaches an Earth BioGenome standard level of 7.C.Q50. These data will be highly useful for future investigations towards a better understanding of the evolution, development, morphology, and toxicology of Nemertea.

Keywords: de novo assembly, genome sequence, 3D genomics, Hi-C, HiFi

Significance

The genome of *Emplectonema gracile* is highly contiguous, well annotated, and shorter than those of the other two ribbon worm species sequenced to date. This genome is a valuable resource for studies on molecular ecology, venom evolution, and regeneration in marine invertebrates.

Introduction

Nemerteans, commonly known as ribbon worms, are a phylum of about 1,200 species of predatory worms that exhibit spiral cleavage and a variety of life histories, typically including pelagic and benthic stages (Gibson 1994, Maslakova and Hiebert 2015). While nemerteans are mainly marine, some species have entered freshwater habitats, and a few have colonized moist, terrestrial habitats (Gibson 1994). Phylogenetically, Nemertea is nested within Spiralia (*sensu* Giribet and Edgecombe 2020). However, their exact phylogenetic position is not well established (Struck and Fisse 2008; Struck et al. 2014; Andrade et al. 2014; Laumer et al. 2015; Kocot et al. 2017; Bleidorn 2019; Marlétaz et al. 2019; Drábková et al. 2022).

Nemertea is divided into three main clades: Paleonemertea, Pilidiophora, and Hoplonemertea (Figure 1A) (Andrade et al. 2014; Kvist et al. 2014). To date, only two nemertean nuclear genomes are available, both coming from species of the family Lineidae

(Pilidiophora). One of these nemertean genomes, *Lineus longissimus*, meets the current Earth BioGenome Project standards for reference genomes (Kwiatkowski et al. 2021), while the second, the genome of *Notospermus geniculatus*, was published only at the level of scaffolds (Luo et al. 2018) (Figure 1A).

In this study, we aimed to enrich the available genomic resources for Nemertea, a venomous animal group characterized by their regeneration capacities and obscure phylogenetic position (Stricker and Cloney 1983; Zattara et al. 2019). For it, we have sequenced, assembled, and annotated the genome of one representative of the Hoplonemertea clade, the species *Emplectonema gracile* (Figure 1A).

Emplectonema gracile inhabits the rocky shores of the North Atlantic Ocean and the Mediterranean Sea. This species has been selected for ease of collection, culturing, and spawning in the lab. Its slender, bi-toned body, armored with a venomous stylet used for capturing prey, can reach lengths of up to about 50 cm. The availability of the *E. gracile* genome will facilitate evolutionary, developmental, morphological, and toxicological studies within Nemertea, and will ultimately contribute to clarify the phylogenetic position of nemerteans within the Tree of life.

Results and Discussion

For this diploid genome, HiFi sequencing yielded 26.6 Gbp of information contained in a total of 1,779,646 reads. Analysis of the genomic data with GenomeScope (Vurture et al. 2017) inferred a genome size of 157.9 Mbp with a heterozygosity of 1.5%, a uniqueness of 76.8%, and an error rate of 0.1% (Figure 1B, Figure S1A). K-mer analysis indicates the k-mers with the highest percentage out of all heterozygous k-mer pairs are diploid (92%) with only a small percentage of heterozygous k-mers being triploid or tetraploid (4%); these values indicate that the genome is diploid (Figure S1B).

Assembly, purging redundant haplotigs, and filtering out contamination resulted in a final genome assembly of 157.9 Mbp consisting of 22 contigs with an N50 of 10.0 Mbp (Figure 1B, Table 1). After scaffolding with Hi-C reads, the longest scaffold was 17.8 Mbp and the scaffold L50 was 6. The overall quality of the *E. gracile* genome was established by means of Mercury as 66.0 (Rhie et al. 2020). The assembly is rather complete with 95.3% of BUSCO markers detected (95.3% complete including 0.2% duplicated markers plus 1.0% fragmented). 99.8% of the k-mers mapped to the combined primary and alternative assembly, while 79.1% mapped to the assembly using only the primary one. The selected values for the different assembly parameters for each step (i.e., the unpurged primary genome assembly, the primary assembly after purging haplotypic duplications, the decontaminated primary genome assembly, and the HiC scaffolded genome) are shown in Table 1.

For Hi-C, 52,305,505 reads were obtained. According to our Hi-C assembly and structural annotation, the *E. gracile* genome contains 15 putative chromosomes (Fig. 1C). After annotation, the resulting BUSCO protein score was 89.6% complete (84.6% single copy, 5% duplicated, 3.1% fragmented). This protein assessment was done employing the protein set BUSCO metazoan_odb10 dataset (Simão et al. 2015) combined with the selected nemertean transcriptomes shown in Table S1.

With values ranging between 157.9 and 161.8 Mbp, the *Emplectonema gracile* genome is shorter than the haploid size of 210 ± 5 Mbp previously estimated based on flow cytometry (Paule et al. 2021). Moreover, the *E. gracile* genome is substantially smaller than those of the two pilidiophoran species previously sequenced: *Lineus longissimus* (391 Mbp in 15 putative chromosomes distributed over 109 contigs) (Kwiatkowski et al. 2021), and *Notospermus geniculatus* (859 Mbp distributed over 60,228 contigs) (Luo et al. 2018).

Materials and Methods

Sampling

For this work, two specimens of *Emplectonema gracile* of unknown sex were selected. One was used for HiFi sequencing (specimen ID N59), while the other one was used for Hi-C sequencing (specimen ID N53) (Lieberman-Aiden et al. 2009; Hu et al. 2021). Both individuals were collected in Norway from the beach of Jeløya (Moss, Viken; N 59°25'23.6", E 010°37'05.9"; WGS84; ±2 m). The specimens, approximately 10 cm long, were cut. The anterior body part of N53 was preserved in 4% formaldehyde as a voucher (Natural History Museum, University of Oslo, Norway; catalog number NHMO C7190). Remaining pieces of N53, and all N59, were flash frozen in liquid nitrogen.

Identification

Morphological identification and DNA barcoding were performed for both specimens. Specimen N59 was barcoded using the DNA extracted for HiFi sequencing, while the DNA of N53 was extracted using the Qiagen DNeasy Blood and Tissue Kit following manufacturer's instructions. 16S and 18S rRNA gene sequences of N59 were PCR amplified using the following primers: forward (16S: 5' CCGGTCTGAACTCAGATCACGT 3'; 18S: 5' CCCCATAATTGGAATGAGTACA 3'), reverse (16S: 5' CGCCTGTTTATCAAAAACAT 3'; 18S: 5' AGCTCTCAATCTTGTCATCCT 3'); and the following settings: 1x 2 minutes at 94 °C, 40x [30 seconds at 94 °C, 60 seconds at 51 °C, 60 seconds at 72 °C], 1x 2 minutes at 72 °C. For N53, COI was PCR amplified using forward primer LCO1490-JJ (5' CHACWAAYCATAAAGATARYGG 3'), and reverse primer HCO2198 JJ (5' AWAATTCVGGRTGVCCAAARAARCA 3') (Astrin and Stüben 2008). Resulting PCR products were Sanger sequenced by Macrogen (Amsterdam). The 16S (XXX), 18S (XXX), and COI (XXX) sequences confirmed the morphological identification. The COI sequence for N53 was identical to the COI sequence of the publicly available mitochondrial genome of *E. gracile*

(NC_016952.1). The mitochondrial genome of the *Emplectonema gracile* specimen N59 determined by BLAST had a >99% similarity to the previously published mitochondrial genome for the same species (NCBI accession number JF727825).

Genome Sequencing

For HiFi sequencing, DNA was extracted from posterior parts of N59. Samples were weighed and minced on dry ice followed by tissue disruption using a TissueRuptor II (QIAGEN, Germany) on its maximum settings for 10 seconds. High molecular weight (HMW) DNA was extracted using the Nanobind Tissue Big DNA kit (Circulomics Inc, USA). DNA quality and concentration were determined with a Nanodrop UV/Vis spectrophotometer (Thermo Fisher Scientific, USA), a Qubit BR dsDNA assay (Thermo Fisher Scientific), a pulsed field gel, and a Fragment Analyzer (Agilent, USA) run. Low molecular weight DNA was removed using the BluePippin (Sage Science, USA) system with High Pass Plus Gel cassettes. DNA was further purified and concentrated using the AMPure XP purification kit (Beckman Coulter, USA). A final concentration of 143 ng/ μ L in a volume of 75 μ L was obtained. The library for HiFi circular consensus sequencing was constructed and sequenced on a SEQUEL II (Pacific Biosciences) platform at the Norwegian Sequencing Centre (Oslo, Norway).

For Hi-C sequencing, a library was prepared using the Arima High Coverage Hi-C+ Kit (Arima Genomics, USA). Specifically, the restriction enzymes used for Arima Hi-C 2.0 cut at the following recognition sites: \wedge GATC, G \wedge ANTC, C \wedge TNAG, T \wedge TAA. For this reaction, ~40 mg of disrupted tissue was used. Subsequently, a library was generated following the manufacturer's instructions. For quality control, a combination of Qubit (Thermo Fisher Scientific, USA) and Fragment Analyzer (Agilent, USA) measurements, as well as a Kapa library quantification kit for Illumina libraries (Roche, Switzerland), were used. The final barcoded library was pooled on a quarter S4 flow cell in 2x150 bp paired-end mode on an

Illumina NovaSeq sequencer (Illumina, USA). Hi-C library preparation and sequencing were done at the Norwegian Sequencing Centre (Oslo, Norway).

Genome Profiling

Genome profiling steps were followed to assess k-mer frequencies within raw sequencing reads, and to estimate major genome characteristics such as genome size, heterozygosity, and repetitiveness. The k-mer distribution, with a k-mer size of 21, was generated using Jellyfish 2.3.0 using default settings. Based on this k-mer distribution, SmudgePlot 0.2.4 was run to test the ploidy of the genome (Marcais et al. 2012; Ranallo-Benavidez et al. 2020). GenomeScope 2.0 (Ranallo-Benavidez et al. 2020) with a k-mer size of 21, diploid level and a high-bound value of 1 million, was used to calculate the genome size, repetitiveness, and heterozygosity for a diploid genome using a combinatorial approach fitting a mathematical model to the k-mer distribution.

De Novo Genome Assembly

Assembly of HiFi reads was carried out with Hifiasm 0.18.2 (Cheng et al. 2021), using default settings. Haplotypic duplications were purged with Purge_dups 1.4 (Guan et al. 2020). The primary assembly was checked for contamination and corrected using the BlobToolKit 3.1.4 software (Challis et al. 2020); all contigs not matching metazoan hits were excluded (Fig. S1C).

Computational Scaffolding

The Arima Mapping Pipeline (https://github.com/ArimaGenomics/mapping_pipeline) was used for mapping raw Hi-C reads to the purged and decontaminated assembly outlined above. Briefly, Hi-C paired reads are first aligned to the reference independently using BWA-MEM to identify potential chimeric reads to be filtered out. Filtered single-end Hi-C reads are then paired and sorted based on mapping quality to produce a quality filtered, paired-end BAM file. Picard Tools is then used to flag PCR duplicates which are then discarded using SAMtools

(Camacho et al. 2009). The quality filtered BAM file was then used as input for scaffolding. Scaffolding was performed in YaHS: yet another Hi-C scaffolding tool (<https://github.com/czhou/yahs>) (Zhou et al. 2022). The output of YaHS was then converted into .hic and .assembly files using Juicer tools 1.9.9_jcuda.0 (Durand et al. 2016) for manual curation in Juicebox Assembly Tools 1.9.1 (Dudchenko et al. 2018). The Hi-C contact map generation was done using Juicebox 1.9.1 (Robinson et al. 2018).

Quality Control Checks

Several quality control checks were conducted after each analytical step (i.e., unpurged assembly, assembly after purging, decontaminated assembly, and scaffolded genome). Quast 5.0.2 (Gurevich et al. 2013) was used to determine statistical parameters of the primary genome assembly. Using Merqury 1.3 (Rhie et al. 2020), a meryl database was generated, and quality statistics such as consensus quality and k-mer completeness retrieved. The different assemblies were benchmarked against the 954 universal single-copy orthologs of the metazoa_odb10 dataset using BUSCO+ 5.5.0 (Simão et al. 2015, Manni et al. 2021). In preparation for BlobToolKit, Blast+ 2.13.0 (Camacho and Madden 2013) was used to map each contig of the assembly against a local copy of the NCBI nucleotide (nt) database downloaded as part of the pipeline. Additionally, HiFi reads were mapped against the primary assembly with Minimap2 2.17 (Li 2018), and further prepared for BlobTools with SAMtools 1.10 (Camacho et al. 2009). The BUSCO scores, BLAST results, and read coverage were uploaded and further analyzed within BlobToolKit 3.1.4. After assembly, mitochondrial genome sequences were retrieved using Blast+ and the amino-acid sequences of all protein-coding genes of the mitochondrial genome NC_000931. These mitochondrial sequences were queried against the nt with Blast+ 2.13.0 to confirm species identification and possible sources of contamination.

Genome Annotation

For structural annotation, RepeatModeler 2.0.1 (Flynn et al. 2020) was used to model repeat content followed by soft repeat masking utilizing RepeatMasker 4.1.2 with default settings (Smith et al. 2015). As no transcriptome was available for this species, annotation was performed using the Braker 3 (Hoff et al. 2019) pipeline based on protein sequences from closely related species. 37 publicly available nemertean transcriptomes belonging to 25 species (Table S1) were downloaded from NCBI, assembled with Trinity (Grabher et al. 2011), and translated with TransDecoder (<https://github.com/TransDecoder/>). The translated transcriptomes, combined with the OrthoDB v10 metazoa dataset, were used to generate protein prediction hints with ProtHint 2.6.0 (Hoff et al. 2019). The ProtHint mapping pipeline was used by Braker 3 to produce protein hints to train the model. The soft-masked and decontaminated primary genome assembly and protein databases were used as input to Braker 3. Gene annotations were assessed for completeness using BUSCO+ 5.5.0 (Simão et al. 2015) metazoa odb_10 database.

Data Availability

European Nucleotide Archive: *Emplectonema gracile*. Accession number XXX. Unprocessed sequence data have been archived in the NCBI Sequence Read Archive under Bioproject XXX. The Refseq genome assembly can be found under accession XXX along with NCBI *E. gracile* Annotation release 1000. The genome sequence is released openly for reuse. All custom scripts are available at GitHub <https://github.com/torstenstruck/InvertOmics> and <https://github.com/mkyapchiongco/Hi-C-WorkFlow>.

Ethical approval

The Nagoya protocol does not apply to this work. Both sample collection and molecular work were done in Norway.

Acknowledgements

This work was funded by the Research Council of Norway project “InvertOmics – Phylogeny and evolution of lophotrochozoan invertebrates based on genomic data” (Project number: 300587 to THS). KMK, MKY, and NGR were funded by NSF DEB-1846174. We thank Matz Berggren (University of Gothenburg) for assisting and allowing AVG to use the camera setup used for taking the picture of a specimen of *E. gracile* included in Figure 1A, and Miguel Ángel Naranjo Ortiz (University of Oslo) and Emanuela Di Martino (University of Oslo) for discussion.

Author Contributions

AVG and THS collected and preserved the specimens of *Emplectonema gracile*, identified by AVG. AVG carried out the molecular parts of this work which were not conducted at the Norwegian Sequencing Center, except for barcoding the voucher specimen N53, done by ATC. AVG ran the different genome profiling, genome assembly, and quality control checks developed by THS. MY-C performed Hi-C scaffolding of the genome plotted by AVG. NR performed the genome structural annotation. AVG, KMK, MM, and THS conceived the study. AVG wrote the first draft of the manuscript, and all authors approved the submitted version.

References

Ament-Velásquez SL, Figuet E, Ballenghien M, Zattara EE, Norenburg JL, Fernández-Álvarez FA, Bierne J, Bierne N, Galtier N. Population genomics of sexual and asexual lineages

- in fissiparous ribbon worms (Lineus, Nemertea): hybridization, polyploidy and the Meselson effect. *Mol Ecol.* 2016;25(14):3356–3369. <https://doi.org/10.1111/mec.13717>.
- Andrade SC, Montenegro H, Strand M, Schwartz ML, Kajihara H, Norenburg JL, Turbeville JM, Sundberg P, Giribet G. A transcriptomic approach to ribbon worm systematics (Nemertea): resolving the Pilidiophora problem. *Mol Biol Evol.* 2014;31(12): 3206–3215. <https://doi.org/10.1093/molbev/msu253>.
- Astrin JJ, Stüben PE. Phylogeny in cryptic weevils: molecules, morphology and new genera of western Palaearctic Cryptorhynchinae (Coleoptera: Curculionidae). *Invertebr Syst.* 2008;22(5):503–522. <https://doi.org/10.1071/IS07057>.
- Bleidorn C. Recent progress in reconstructing lophotrochozoan (spiralian) phylogeny. *Org Divers Evol.* 2019;19(4):557–566. <https://doi.org/10.1007/s13127-019-00412-4>.
- Brůna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. BRAKER2: automatic eukaryotic genome annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genom Bioinform.* 2021;3(1):lqaa108. <https://doi.org/10.1093/nargab/lqaa108>.
- Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. BLAST+: architecture and applications. *BMC Bioinformatics.* 2009;10:1–9. <https://doi.org/10.1186/1471-2105-10-421>.
- Camacho C, Madden T. BLAST+ release notes. BLAST® Help. 2013 [Internet].
- Challis R, Richards E, Rajan J, Cochrane G, Blaxter M. BlobToolKit–interactive quality assessment of genome assemblies. *G3 (Bethesda).* 2020;10(4):1361–1374. <https://doi.org/10.1534/g3.119.400908>.

- Cheng H, Concepcion GT, Feng X, Zhang H, Li H. Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat Methods*. 2021;18(2):170–175. <https://doi.org/10.1038/s41592-020-01056-5>.
- Coe WR. *On the development of the pilidium of certain nemerteans* (Volume 10 of Transactions of the Connecticut Academy of Arts and Sciences). 1899. Connecticut Academy of Arts and Sciences, USA.
- Drábková M, Kocot KM, Halanych KM, Oakley TH, Moroz LL, Cannon JT, Kuris A, Garcia-Vedrenne AE, Pankey MS, Ellis EA, et al. Different phylogenomic methods support monophyly of enigmatic ‘Mesozoa’ (Dicyemida+Orthonectida, Lophotrochozoa). *Proc R Soc Lond B Biol Sci*. 2022;289(1978):20220683. <https://doi.org/10.1098/rspb.2022.0683>.
- Dudchenko O, Shamim MS, Batra SS, Durand NC, Musial NT, Mostofa R, Pham M, St Hilaire BG, Yao W, Stamenova E, et al. The Juicebox Assembly Tools module facilitates de novo assembly of mammalian genomes with chromosome-length scaffolds for under \$1000. *bioRxiv*. 2018: 254797. <https://doi.org/10.1101/254797>.
- Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, Aiden EL. Juicer provides a one-click system for analyzing loop-resolution Hi-C experiments. *Cell Syst*. 2016;3(1):95–98. <http://dx.doi.org/10.1016/j.cels.2016.07.002>.
- Egger B, Lapraz F, Tomiczek B, Müller S, Dessimoz C, Girstmair J, Škunca N, Rawlinson KA, Cameron CB, Beli E, et al. A transcriptomic-phylogenomic analysis of the evolutionary relationships of flatworms. *Curr Biol*. 2015;25(10):1347–1353. <http://dx.doi.org/10.1016/j.cub.2015.03.034>.
- Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. RepeatModeler2 for automated genomic discovery of transposable element families. *Proc Natl Acad Sci USA*. 2020;117(17): 9451–9457. <https://doi.org/10.1073/pnas.1921046117>.

- Gibson R. *Nemertean* 2nd ed. 1994: vol. 24. Field Studies Council, USA.
- Giribet G, Edgecombe GD. *The Invertebrate Tree of Life*. 2020. Princeton University Press, USA.
- Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, et al. Trinity: reconstructing a full-length transcriptome without a genome from RNA-Seq data. *Nat Biotechnol*. 2011;29(7):644. <http://dx.doi.org/10.1038/nbt.1883>
- Guan D, McCarthy SA, Wood J, Howe K, Wang Y, Durbin R. Identifying and removing haplotypic duplication in primary genome assemblies. *Bioinformatics*. 2020;36(9):2896–2898. <https://doi.org/10.1093/bioinformatics/btaa025>.
- Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUASt: quality assessment tool for genome assemblies. *Bioinformatics*. 2013;29(8):1072–1075. <https://doi.org/10.1093/bioinformatics/btt086>.
- Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: unsupervised RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics*. 2016;32(5):767–769. <https://doi.org/10.1093/bioinformatics/btv661>.
- Hoff KJ, Lomsadze A, Borodovsky M, Stanke M. Whole-genome annotation with BRAKER. *Gene Prediction: Methods and Protocols*. Humana, New York, NY. 2019, p. 65–95. https://doi.org/10.1007/978-1-4939-9173-0_5.
- Hoff KJ, Stanke M. Predicting genes in single genomes with AUGUSTUS. *Curr Protoc Bioinformatics*. 2019;65(1):e57. <https://doi.org/10.1002/cpbi.57>.
- Hu T, Chitnis N, Monos D, Dinh A. Next-generation sequencing technologies: An overview. *Hum Immunol*. 2021;82(11):801–811. <https://doi.org/10.1016/j.humimm.2021.02.012>.
- Jiang Z, Rokhsar DS, Harland RM. Old can be new again: HAPPY whole genome sequencing, mapping and assembly. *Int J Biol Sci*. 2009;5(4):298. <https://doi.org/10.7150/ijbs.5.298>.

- Kocot KM, Struck TH, Merkel J, Waits DS, Todt C, Brannock PM, Weese DA, Cannon JT, Moroz LL, Lieb B, et al. Phylogenomics of Lophotrochozoa with consideration of systematic error. *Syst Biol.* 2017;66(2):256–282. <https://doi.org/10.1093/sysbio/syw079>.
- Kvist S, Laumer CE, Junoy J, Giribet G. New insights into the phylogeny, systematics and DNA barcoding of Nemertea. *Invertebr Syst.* 2014;28(3):287–308. <https://doi.org/10.1071/IS13061>.
- Kwiatkowski D, Blaxter M, Darwin Tree of Life Barcoding collective, Wellcome Sanger Institute Tree of Life programme, Wellcome Sanger Institute Scientific Operations: DNA Pipelines collective, Tree of Life Core Informatics collective, Darwin Tree of Life Consortium. The genome sequence of the bootlace worm, *Lineus longissimus* (Gunnerus, 1770). *Wellcome Open Res.* 2021;6. <https://doi.org/10.12688/wellcomeopenres.17193.1>.
- Laumer CE, Bekkouche N, Kerbl A, Goetz F, Neves RC, Sørensen MV, Kristensen RM, Hejnol A, Dunn CW, Giribet G, et al. Spiralian phylogeny informs the evolution of microscopic lineages. *Current Biol.* 2015;25(15):2000–2006. <http://dx.doi.org/10.1016/j.cub.2015.06.068>.
- Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34(18):3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>.
- Lieberman-Aiden E, Van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ, Dorschner MO, et al. Comprehensive mapping of long-range interactions reveals folding principles of the human genome. *Science.* 2009;326(5950):289–293. <https://doi.org/10.1126/science.118136>.
- Luo YJ, Kanda M, Koyanagi R, Hisata K, Akiyama T, Sakamoto H, Sakamoto T, Satoh N. Nemertean and phoronid genomes reveal lophotrochozoan evolution and the origin of

- bilaterian heads. *Nat Ecol Evol.* 2018;2(1):141–151. <https://doi.org/10.1038/s41559-017-0389-y>.
- Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. BUSCO update: novel and streamlined workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic, and viral genomes. *MBE.* 2021;38(10):4647–4654. <https://doi.org/10.1093/molbev/msab199>.
- Marcais G, Kingsford C. Jellyfish: A fast k-mer counter. Version 1.1.4. *Tutorialis e Manuais.* 2012;1:1–8.
- Marlétaz F, Peijnenburg KT, Goto T, Satoh N, Rokhsar DS. A new spiralian phylogeny places the enigmatic arrow worms among gnathiferans. *Current Biol.* 2019;2012:29(2):312–318. <https://doi.org/10.1016/j.cub.2018.11.042>.
- Maslakova SA, Hiebert TC. From trochophore to pilidium and back again—a larva’s journey. *Int J Dev Biol.* 2015;58(6–8):585–591. <https://doi.org/10.1387/ijdb.140090sm>.
- Paule J, von Döhren J, Sagorny C, Nilsson MA. Genome Size Dynamics in Marine Ribbon Worms (Nemertea, Spiralia). *Genes.* 2021;12(9):1347. <https://doi.org/10.3390/genes12091347>.
- 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. <https://doi.org/10.1038/s41467-020-14998-3>.
- Rhie A, Walenz BP, Koren S, Phillippy AM. Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biol.* 2020;21(1):1–27. <https://doi.org/10.1186/s13059-020-02134-9>.
- Schmidt GA. Issledovania po embryologii nemertin. II. Pilidii *Cerebratulus pantherinus* i *marginatus*. *Russkii Zool Zh.* 1930;10:113–127.

- Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*. 2015;31(19):3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>.
- Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0. 2015:2013–2015.
- Smit AFA, Hubley S, Green P. RepeatMasker. 2021:Version 4.1. 2.
- Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*. 2008;24(5):637–644. <https://doi.org/10.1093/bioinformatics/btn013>.
- Stanke M, Schöffmann O, Morgenstern B, Waack S. Gene prediction in eukaryotes with a generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics*. 2006;7(1):1–11. <https://doi.org/10.1186/1471-2105-7-62>.
- Stricker SA, Cloney RA. The ultrastructure of venom-producing cells in *Paranemertes peregrina* (Nemertea, Hoplonemertea). *J Morphol*. 1983;177(1):89–107. <https://doi.org/10.1002/jmor.1051770108>.
- Struck TH, Fisse F. Phylogenetic position of Nemertea derived from phylogenomic data. *Mol Biol Evol*. 2008;25(4):728–736. <https://doi.org/10.1093/molbev/msn019>.
- Struck TH, Wey-Fabrizius AR, Golombek A, Hering L, Weigert A, Bleidorn C, Klebow S, Iakovenko N, Hausdorf B, Petersen M, et al. Platyzoan paraphyly based on phylogenomic data supports a noncoelomate ancestry of Spiralia. *Mol Biol Evol*. 2014;31(7):1833–1849. <https://doi.org/10.1093/molbev/msu143>.
- Uliano-Silva M, Ferreira JGR, Krasheninnikova K, Formenti G, Abueg L, Torrance J, Myers EW, Durbin R, Blaxter M, McCarthy SA. MitoHiFi: a python pipeline for mitochondrial genome assembly from PacBio high fidelity reads. *BMC Bioinformatics*. 2023;24(1):288. <https://doi.org/10.1186/s12859-023-05385-y>.

Vlasenko AE, Kuznetsov VG, Magarlamov TY. Investigation of Peptide Toxin Diversity in Ribbon Worms (Nemertea) Using a Transcriptomic Approach. *Toxins*. 2022;14(8):542. <https://doi.org/10.3390/toxins14080542>.

von Reumont BM, Lüddecke T, Timm T, Lochnit G, Vilcinskas A, von Döhren J, Nilsson MA. Proteo-transcriptomic analysis identifies potential novel toxins secreted by the predatory, prey-piercing ribbon worm *Amphiporus lactifloreus*. *Mar Drugs*. 2020;18(8):407. <https://doi.org/10.3390/md18080407>.

Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC. GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics*. 2017;33(14):2202–2204. <https://doi.org/10.1093/bioinformatics/btx153>.

Zattara EE, Fernández-Álvarez FA, Hiebert TC, Bely AE, Norenburg JL. A phylum-wide survey reveals multiple independent gains of head regeneration in Nemertea. *Proc R Soc Lond B Biol Sci*. 2019;286(1898):20182524. <https://doi.org/10.1098/rspb.2018.2524>.

Zhou C, McCarthy SA, Durbin R. YaHS: yet another Hi-C scaffolding tool. *Bioinformatics*. 2023;39(1):btac808. <https://doi.org/10.1093/bioinformatics/btac808>.

Figures and Tables

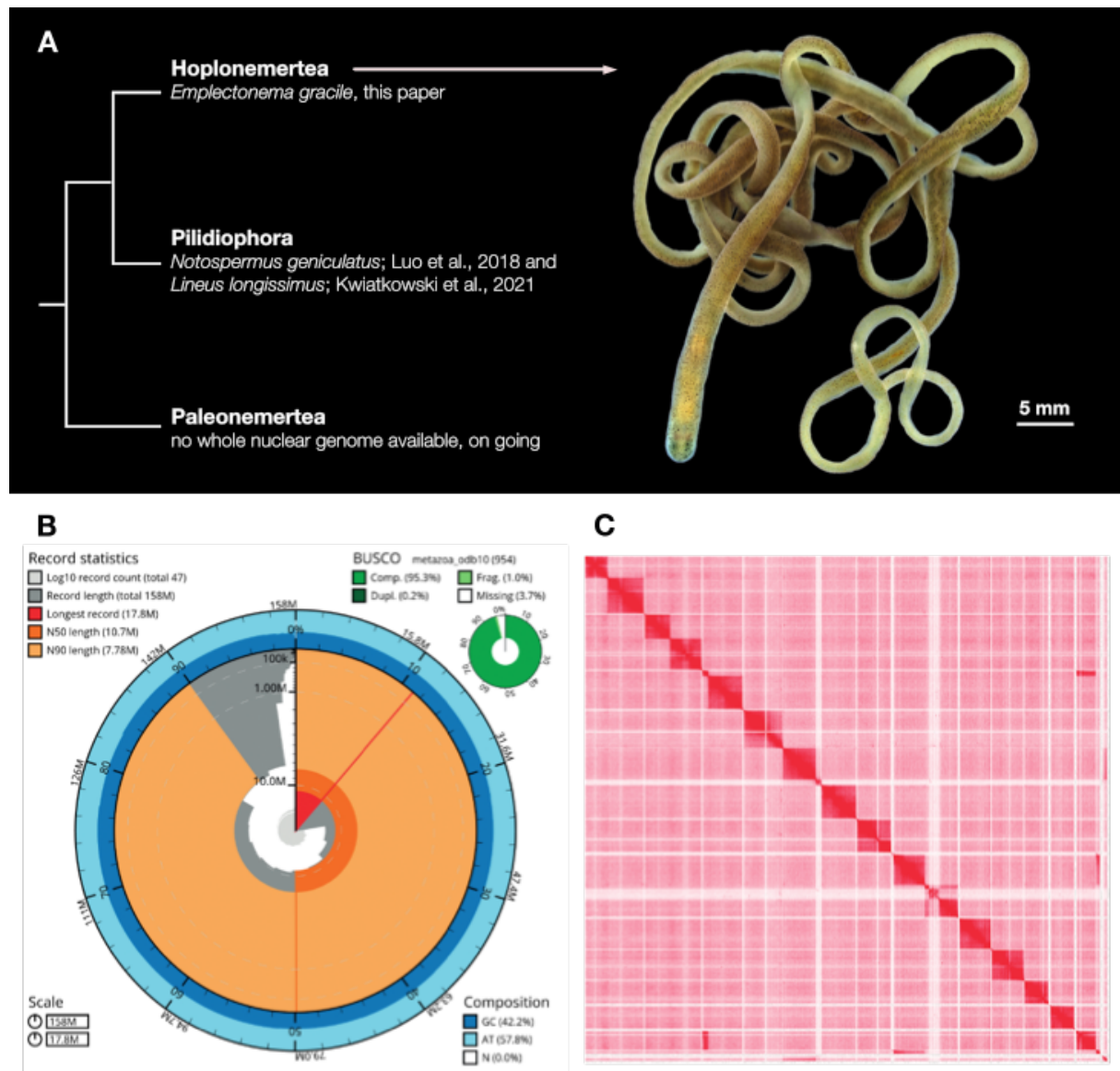


Fig. 1. –The genome of *Emplectonema gracile*. **(A)** Phylogeny of subgroups within the Nemertea phylum, with information about the species for which a whole genome is available (left), and adult specimen of *Emplectonema gracile* (right), photo taken by AVG. **(B)** “Snailplot” produced with BlobToolKit, illustrating N50 metrics and BUSCO gene completeness. **(C)** Hi-C contact map representing the final genome assembly of *E. gracile*, visualized with Juicebox 1.9.1.

PROJECT ACCESSION DATA					
Assembly Identifier	XXXX				
Species	<i>Emplectonema gracile</i>				
Specimens	tnEmpGrac1				
NCBI taxonomy ID	6230				
BioProject	XXXX				
BioSample ID	XXXX				
Isolate information	XXXX				
RAW DATA ACCESSION NUMBERS					
PacificBiosciences Sequel II	XXXX				
Hi-C Illumina	XXXX				
GENOME ASSEMBLY					
Assembly accession	XXXX				
ASSEMBLY METRICS	Unpurged primary genome assembly	Primary genome assembly after purging haplotypic duplications	Decontaminated primary genome assembly	HiC scaffolded genome	Benchmark
Span (Mb)	161.8	158.5	157.9	157.9	na
Number of contigs	135	49	22	na	na
Contig N50 length (Mb)	10.0	10.0	10.0	na	≥ 10
Longest contig (Mb)	13.1	13.1	13.1	na	na
Contig L50 length	8	8	8	na	na
Consensus quality (QV)	61.6 (primary only)/ 62.8 (primary & alternative)	65.0 (primary only)	67.0 (primary only)	67.0 (primary only)	≥ 50
k-mer completeness	79.1% (primary only)/ 99.8% (primary & alternative)	79.1% (primary only)	79.1% (primary only)	79.1% (primary only)	≥ 95%
BUSCO scores (n:954)	C:95.6%[S:95.3%,D:0.3%], F:0.7%,M:3.7%	C:95.6%[S:95.3%,D:0.3%], F:0.7%,M:3.7%	C:95.5%[S:95.2%,D:0.3%], F:0.8%,M:3.7%	C:95.3%[S:95.1%,D:0.2%], F:1.0%,M:3.7%	C ≥ 95%
BUSCO protein	na	na	na	C:89.6%[S:84.6%,D:5%], F:3.1%,M:7.3%	
Percentage of assembly mapped to chromosomes	na	na	na	99.1%	≥ 95%
Number of scaffolds		na	na	47	na
Scaffold N50 length (Mb)	na	na	na	10.7	≥
Longest scaffold (Mb)	na	na	na	17.8	na
Scaffold L50 length	na	na	na	6	na
Organelles Mitochondrial genome assembled	complete single contig	na	na	na	na

Table 1. – Project accession data and Assembly information for *E. gracile*. *BUSCO scored based on the metazoan_odb10 BUSCO set using v5.1.2. C=complete [S=single copy, D=duplicated], F=fragmented, M=missing, n=number of orthologues in comparison.