



OPEN ACCESS

EDITED BY

Iker Irisarri,
University of Göttingen, Germany

REVIEWED BY

Joel Vizueta,
University of Copenhagen, Denmark
Andreas Hejnol,
University of Bergen, Norway

*CORRESPONDENCE

Samuel Abalde,
✉ saabalde@gmail.com
Ulf Jondelius,
✉ ulf.jondelius@zoologi.su.se

RECEIVED 23 June 2023

ACCEPTED 12 September 2023

PUBLISHED 26 September 2023

CITATION

Abalde S, Tellgren-Roth C, Heintz J,
Vinnere Pettersson O and Jondelius U
(2023), The draft genome of the
microscopic *Nemertoderma westbladi*
sheds light on the evolution of
Acoelomorpha genomes.
Front. Genet. 14:1244493.
doi: 10.3389/fgene.2023.1244493

COPYRIGHT

© 2023 Abalde, Tellgren-Roth, Heintz,
Vinnere Pettersson and Jondelius. This is
an open-access article distributed under
the terms of the [Creative Commons
Attribution License \(CC BY\)](#). The use,
distribution or reproduction in other
forums is permitted, provided the original
author(s) and the copyright owner(s) are
credited and that the original publication
in this journal is cited, in accordance with
accepted academic practice. No use,
distribution or reproduction is permitted
which does not comply with these terms.

The draft genome of the microscopic *Nemertoderma westbladi* sheds light on the evolution of Acoelomorpha genomes

Samuel Abalde^{1*}, Christian Tellgren-Roth², Julia Heintz²,
Olga Vinnere Pettersson² and Ulf Jondelius^{1,3*}

¹Department of Zoology, Swedish Museum of Natural History, Stockholm, Sweden, ²Department of Immunology, Genetics and Pathology, SciLifeLab, Uppsala University, Uppsala, Sweden, ³Department of Zoology, Stockholm University, Stockholm, Sweden

Background: Xenacoelomorpha is a marine clade of microscopic worms that is an important model system for understanding the evolution of key bilaterian novelties, such as the excretory system. Nevertheless, Xenacoelomorpha genomics has been restricted to a few species that either can be cultured in the lab or are centimetres long. Thus far, no genomes are available for Nemertodermatida, one of the group's main clades and whose origin has been dated more than 400 million years ago.

Methods: DNA was extracted from a single specimen and sequenced with HiFi following the PacBio Ultra-Low DNA Input protocol. After genome assembly, decontamination, and annotation, the genome quality was benchmarked using two acoel genomes and one Illumina genome as reference. The gene content of three cnidarians, three acoelomorphs, four deuterostomes, and eight protostomes was clustered in orthogroups to make inferences of gene content evolution. Finally, we focused on the genes related to the ultrafiltration excretory system to compare patterns of presence/absence and gene architecture among these clades.

Results: We present the first nemertodermatid genome sequenced from a single specimen of *Nemertoderma westbladi*. Although genome contiguity remains challenging (N50: 60 kb), it is very complete (BUSCO: 80.2%, Metazoa; 88.6%, Eukaryota) and the quality of the annotation allows fine-detail analyses of genome evolution. Acoelomorph genomes seem to be relatively conserved in terms of the percentage of repeats, number of genes, number of exons per gene and intron size. In addition, a high fraction of genes present in both protostomes and deuterostomes are absent in Acoelomorpha. Interestingly, we show that all genes related to the excretory system are present in Xenacoelomorpha except *Osr*, a key element in the development of these organs and whose acquisition seems to be interconnected with the origin of the specialised excretory system.

Conclusion: Overall, these analyses highlight the potential of the Ultra-Low Input DNA protocol and HiFi to generate high-quality genomes from single animals, even for relatively large genomes, making it a feasible option for sequencing challenging taxa, which will be an exciting resource for comparative genomics analyses.

KEYWORDS

Ultra-Low DNA Input, Xenacoelomorpha, HiFi, gene content, excretory system, *Osr*

1 Introduction

Access to a growing number of high-quality genomes from non-model animal species has helped us understand the origin of key evolutionary novelties (Albertin et al., 2015; Dunwell et al., 2017; Rubin et al., 2019). However, small yields of extracted DNA is a limiting factor in genome sequencing of small animals, also when using whole-body extractions. In this regard, the recent development of Ultra-Low DNA Input protocols has significantly reduced the amount of input DNA, enabling the sequencing of high-quality genomes from millimetric animals (Kingan et al., 2019; Korlach, 2020; Schneider et al., 2021). Yet, this approach is recommended for genomes smaller than 500 Mb (PacBio's library prep protocol), and it is unclear how well it performs beyond that limit, which is not a minor detail. Despite the general trend that miniaturised animals tend to have smaller genomes (Liu et al., 2012; Gross et al., 2019; Xu et al., 2021), there are several animals, such as Xenacoelomorpha, whose genome size is comparable to that of larger animals (Arimoto et al., 2019; Gehrke et al., 2019; Martinez et al., 2022).

Xenacoelomorpha is a clade of mostly marine, microscopic worms consisting of the clades Acoela, Nemertodermatida, and their sister taxon *Xenoturbella*. Early molecular phylogenetic studies placed Xenacoelomorpha as the sister group of all other Bilateria (Ruiz-Trillo et al., 1999; Jondelius et al., 2002). This hypothesis received support from the simple morphology of Xenacoelomorpha, which lacks typical bilaterian structures such as excretory organs, through-gut and circulatory system (Haszprunar, 2016) and the name Nephrozoa was introduced for its sister group under this hypothesis (Jondelius et al., 2002). The Nephrozoa hypothesis was further supported by analyses of gene content and phylogenomic inference (Cannon et al., 2016; Juravel et al., 2023). However, an alternative hypothesis based on analyses of nucleotide sequence data places Xenacoelomorpha as sister group to Ambulacraria (echinoderms and hemichordates) within the deuterostomes (Philippe et al., 2019; Kapli and Telford, 2020). In either case, xenacoelomorphs offer a good opportunity for studying the origin of important animal novelties. Due to their lack of specialised excretory organs, xenacoelomorphs make a good comparison reference to better understand the evolution of this system. A recent study based on spatial transcriptomics has shown the expression in Xenacoelomorpha of several genes involved in the excretory process in other bilaterians, as well as several genes specifically related to the ultrafiltration excretory system (*Nephrin*, *Kirrel*, and *ZOI*; (Andrikou et al., 2019)), although their expression was observed throughout the body, unlike in other organisms with specialised excretory organs (Gašiorowski et al., 2021). In addition to analysing their expression, the comparison of high-quality genomes from xenacoelomorphs, protostomes, and deuterostomes would offer a better understanding of the evolution of these genes, thanks to a more accurate assessment of gene presence/absence, the annotation of all gene copies in the genome, information about their distribution in the genomes, or comparisons of gene architecture, among other analyses. However, the set of available xenacoelomorph genomes is still limited.

Several xenacoelomorph species have drawn interest as model systems to study the evolution of body regeneration, the nervous system, and endosymbiosis (Martín-Durán et al., 2018; Andrikou et al., 2019; Gehrke et al., 2019), resulting in the generation of genomes from *Xenoturbella* [*Xenoturbella bocki*; (Schiffer et al., 2023)] and Acoela [*Hofstenia miamia* and the closely related acoel species *Praesagittifera naikaiensis* and *Symsagittifera roscoffensis*; (Arimoto et al., 2019; Gehrke et al., 2019; Martinez et al., 2022)]. Thus, to fully capture the diversity of Xenacoelomorpha it is necessary to generate new genomes from Nemertodermatida, the sister group of Acoela from which it diverged more than 400 MYBP (Dos Reis et al., 2015). This, however, is challenging due to their microscopic size. The four available xenacoelomorph genomes were sequenced from species that can be either cultured in the lab and/or are relatively big (*X. bocki* and *H. miamia* can reach four and 2 cm in body length, respectively), but that is not the case for the vast majority of xenacoelomorphs, requiring more sophisticated methods. Despite their small size, the acoel genomes sequenced so far range between 700 and 1000 Mb, two to three times larger than any other published genome sequenced with the Ultra-Low Input protocol (Kingan et al., 2019; Korlach, 2020; Schneider et al., 2021), and thus represent a good opportunity for testing its performance in a challenging animal group. Here, we applied the PacBio Ultra-Low DNA Input protocol to sequence the genome of *Nemertoderma westbladi* from a single, microscopic worm, the first nemertodermatid and the longest genome sequenced with this protocol. We demonstrate the potential of this approach to generate relatively good-quality genomes through comparisons with other xenacoelomorphs. In addition, we explore the evolution of acoelomorph genomes, analyze the evolution of gene content in Bilateria and provide insights into the evolution of the genes related to the excretory system.

2 Materials and methods

2.1 DNA extractions, library preparation, and sequencing

High molecular weight DNA was extracted from single individuals of the nemertodermatid *N. westbladi* stored in either ethanol, RNAlater, or RNA Shield using two different methods: the salting-out protocol (Guimaraes, 2018) and the QIAamp Micro DNA kit (Qiagen). The Qubit dsDNA HS kit, a 2% agarose gel, and a Femto Pulse system were used to ensure the extraction met the minimum requirements for DNA yield and fragment size (the majority of gDNA over 20 kb). These worms were isolated in dishes of seawater before their long-term preservation.

Library preparation and sequencing followed the PacBio Ultra-Low DNA Input protocol with small modifications. This protocol has been heavily tested at the Wellcome Sanger Institute by Dr. Laumer, who kindly shared his experience and recommendations with us (Laumer, pers. comm.). Briefly, DNA was sheared to 10 kb using Megaruptor 3 instead of Covaris g-TUBE (Covaris). After removing single-strand overhangs and repairing the fragment ends (SMRTbell Express Template Prep Kit 2.0; PacBio), DNA fragments were ligated to the amplification adapter and PCR

amplified in two independent reactions (Reaction Mix 5A and 5B) for 15 cycles each (SMRTbell® gDNA Sample Amplification Kit; PacBio). Amplified DNA was purified using ProNex Beads (Promega), pooled in a single sample, damage-repaired for the second time, and ligated to the hairpin adapters. Size selection of the prepared SMRTbell library was done using a 35% dilution of AMPure PB beads (PacBio), which removed all fragments shorter than 3 kb, instead of the BluePippin system (Sage Science). Finally, the library was sequenced in one SMRT cell on the Sequel IIe platform.

2.2 Data filtering, assembly, and decontamination

The ‘Trim gDNA Amplification Adapters’ pipeline from SMRT Link v11 was used to remove the adapter sequences. Three genome assembly strategies were attempted and compared: the IPA HiFi Genome Assembler included in SMRT Link v11 (PacBio), Hifiasm v.0.7 (Cheng et al., 2021), and Flye v.2.8.3 (Kolmogorov et al., 2019). Default parameters were used in the three approaches, but with the “--meta” option activated in Flye. Based on genome length and completeness (measured with BUSCO and the metazoa odb10 database; Supplementary Table S1), the Flye assembly was selected for downstream analyses, which included two additional scaffolding approaches. First, the two *N. westbladi* transcriptomes were mapped to the genome using HISAT2 v.2.0.5 (Kim et al., 2019) and fed to P_RNA_SCAFFOLDER (Zhu et al., 2018). Second, the genome of *S. roscoffensis* was used as a reference to map the assembled genome with RagTag v.2.0.1 (Alonge et al., 2022). Unfortunately, none of these attempts improved the genome contiguity any further. Finally, assembly redundancy was removed using the kmerDedup pipeline (<https://github.com/xiekunwhy/kmerDedup>). First, the Kmer spectra ($k = 21$) from the trimmed reads and the assembly were extracted and then merged with Jellyfish v.2.3.0 (Marçais and Kingsford, 2011). The list of kmers was converted to a fasta file and mapped to the assembly with Bowtie2 v.2.5.1 (Langmead and Salzberg, 2012), using the very-sensitive mode and the end-to-end approach, with a minimum alignment score of “L, -0.6, -0.2”. Then, kmerDedup was run with different minimum kmer coverage (mcv) and maximum duplication percentage (mpr), but the best combination according to BUSCO was to set both to 30%.

The raw assembly was decontaminated following the BlobTools2 pipeline (Challis et al., 2020). Coverage data was calculated by mapping the filtered HiFi reads to the assembled genome using Minimap2 with default parameters (Li, 2021), genome completeness inferred with BUSCO v.5.2.2 (Seppey et al., 2019) and the Metazoa odb10 database, and taxonomic information was identified through BLAST searches of the contigs versus the UniProt database (Release 2022_05) using diamond v.0.9.26.127 (Buchfink et al., 2021). Only the contigs identified as “Metazoa” were kept at this stage. Additionally, a BLAST search and a custom database were used to remove mitochondrial contigs. Finally, Minimap2 was used to map the reads back to the decontaminated genome to separate the nemertodermatid reads. The k-mer approaches GenomeScope v.2.0 (assuming a diploid genome and a 35x coverage) and SmudgePlot (Ranallo-Benavidez

et al., 2020) were used to calculate the genome heterozygosity and ploidy before and after the decontamination step with a Kmer length of 21. To identify the contaminant contigs, the diamond output was used to extract the *Taxid* information of the hits, which is associated with a unique taxonomic category on the NCBI database.

2.3 Genome annotation

RepeatMasker v.4.1.2-p1 (Smit et al., 2015) was used to soft mask the repeats in the decontaminated genome with the rmblast engine, for which a custom repeat database was generated with RepeatModeler v.2.0.1 (Smit and Hubley, 2015) and the -LTRStruct option activated. Afterwards, the genome was annotated with BRAKER2 (Brůna et al., 2021) using transcriptomic and proteomic evidence. The two available transcriptomes for *N. westbladi*, sequenced from a whole-body extraction, were downloaded and quality filtered in a two-step approach. Adapters removal and the first trimming were performed with Trimmomatic v.0.36 (as implemented in Trinity v2.6.6 with default parameters, (Grabherr et al., 2011)), followed by a more thorough cleaning with PRINSEQ v.0.20.3 (Schmieder and Edwards, 2011): trim all terminal bases with a quality below 30 and filter out reads whose mean quality is below 25, low complexity sequences (minimum entropy 50), and reads shorter than 75 bp. Clean reads were mapped to the soft-masked genome with STAR v.2.7.9 (Dobin et al., 2013) and the options “--sjdbOverhang 100 --genomeSAindexNbases 13 --genomeChrBinNbits 15” and “--chimSegmentMin 40 --twopassMode Basic”. For the proteomes, the gene models from the coel *P. naikaiensis* (Arimoto et al., 2019), the BUSCO Metazoa odb10 database, and a custom set of single-copy orthogroups, inferred from published transcriptomes with OrthoFinder v.2.4.1 (Emms and Kelly, 2019), were concatenated and mapped to the *N. westbladi* genome using ProtHint v.2.6 (Brůna et al., 2020). The inferred gene models were functionally annotated by pfam_scan v.1.6 (Mistry et al., 2007) and the PFAM 31.0 database.

2.4 Quality control

The quality of the decontaminated genome was assessed using QUAST v.5.2.0 (Gurevich et al., 2013) and the completeness of the genome and the annotation with BUSCO v.5.2.2 using the Metazoa and Eukaryota odb10 databases. Since all the metazoan contigs were kept during the decontamination step, two approaches were followed to ensure they belong to the nemertodermatid genome. First, a distance tree was inferred with FastMe v.2.1.5 (Lefort et al., 2015) based on a distance matrix calculated with Skmer (Sarmashghi et al., 2017), an alignment-free method designed to estimate genomic distances, over the *N. westbladi* genome and 18 metazoan genomes downloaded from GenBank (Supplementary Table S2). Second, a phylogenetic tree was inferred from these genomes except for three for which the annotated proteome was not available. Briefly, orthogroups were inferred with OrthoFinder v.2.4.1 (Emms and Kelly, 2019) and cleaned from paralogs with PhyloPyPruner v.1.2.3 (Thalén et al., 2021) using the “Largest Subtree” method, collapsing nodes with bootstrap support lower than 60, and pruning branches more than

five times longer than the standard deviation of all branch lengths in the tree. Then, orthogroups were aligned with MAFFT v.7.475 using the L-INS-i algorithm (Katoh and Standley, 2013), cleaned from poorly aligned sites with BMGE v.1.12 (Criscuolo and Gribaldo, 2010), tested for stationarity and homogeneity (symmetry tests) with IQ-TREE2 v.2.1.3 (Minh et al., 2020), and concatenated with FASconCAT v.1.05 (Kück and Longo, 2014). Finally, a phylogenetic tree was inferred using coalescence [ASTRAL; (Zhang et al., 2017)] and concatenation by maximum likelihood with a site-specific heterogeneous model (assuming 20 amino acid categories, C20) with IQ-TREE v.1.6.12 (Nguyen et al., 2015).

All the genome metrics, including length, contiguity, number of genes, and completeness, among others, were compared to the acael genomes from *P. naikaiensis* (Arimoto et al., 2019) and *S. roscoffensis* (Martinez et al., 2022), which were also tested for contaminants using BlobTools2, following the same pipeline and with the same filtering criteria. The genomes of *H. miamia* and *X. bocki* (Gehrke et al., 2019; Schiffer et al., 2023) were not considered because an annotation file with details of protein structure is not available for any of them. Additionally, a second *N. westbladi* genome sequenced in an Illumina HiSeq2500 platform was also included in the comparisons to estimate the improvement in genome quality with HiFi data relative to a short-read approach. Briefly, DNA was extracted from a pool of 12 individuals, collected in the same location at the same time, the sequencing library was prepared with a Rubicon kit, and the sequencing generated more than 385 million reads. The Illumina reads were assembled with SPAdes v.3.14.1 (Bankev et al., 2012), with four Kmer lengths (21, 33, 55, 75) and error correction activated. Finally, this genome was analysed with the same parameters as the HiFi genome to eliminate contamination contigs (but not redundancy), produce completeness stats, and annotate gene models.

2.5 Analysis of gene content

To analyse the evolution of gene content in Acoelomorpha, the annotated genomes of 18 metazoans were compared, including *N. westbladi* (Nemertodermatida) and *P. naikaiensis* and *S. symsagittifera* (Acoela) as representatives of Acoelomorpha, eight protostome genomes, four deuterostomes, and three cnidarians as the outgroup to Bilateria (Supplementary Table S2). Redundancies in the gene models of all genomes were removed with CD-HIT (Fu et al., 2012), clustering all sequences more than 95% identical, and then functionally annotated with pfam_scan v.1.6 (Mistry et al., 2007) and the PFAM 31.0 database. The annotated proteins were clustered using OrthoFinder v.2.4.1 (Emms and Kelly, 2019) and used to calculate the number of genes specific to or shared among the four main clades of interest: Cnidaria, Acoelomorpha, Deuterostomia, and Protostomia. Genes present in at least one cnidarian and one bilaterian were considered to be shared across Metazoa, whereas genes present in at least two of Acoelomorpha, Deuterostomia, and Protostomia were considered to be shared across Bilateria. The proportion of “metazoan” and “bilaterian” genes absent from each of the three bilaterian clades was calculated based on these two datasets.

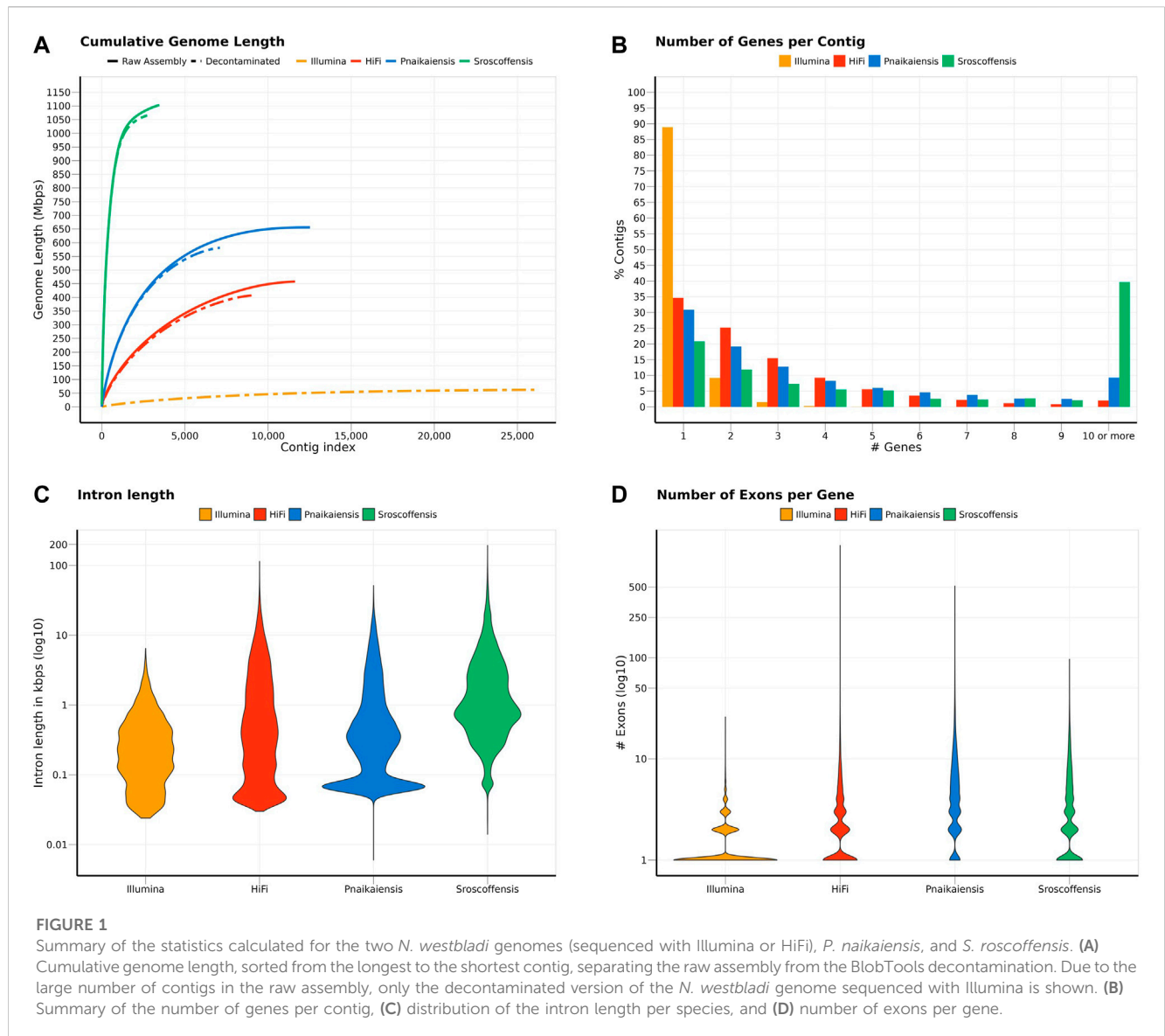
2.6 Annotation and comparison of the genes related to the ultrafiltration excretory system

This analysis was based on the results of Gąsiorowski et al. (Gąsiorowski et al., 2021), who used spatial transcriptomics to identify the genes involved in the development of the ultrafiltration excretory system in several protostomes and one hemichordate species. All the protein sequences annotated in that study were downloaded from GenBank except *Hunchback*, as they found no evidence of this gene being involved in nephridiogenesis, for a total of three structural proteins: *Nephrin*, *Kirrel*, and *ZO1*; and six transcription factors: *Eya*, *Lhx1/5*, *Osr*, *POU3*, *Sall*, and *Six1*. These genes were annotated in the same genomes used to analyse gene content evolution through BLAST searches with diamond v0.9.26.127 (Buchfink et al., 2021). The correct identification of these genes was later confirmed through phylogenetic analyses with IQ-TREE v.1.6.12 (Nguyen et al., 2015) and manual BLAST searches on the NCBI web server. The identification of the *Lhx1/5* and *Six1* transcription factors was not always straightforward, as they are thoroughly mixed in the phylogenetic tree with many other gene variants and sometimes different isoform names were proposed in the BLAST searches for the same sequence, and thus they represent a mixture of isoforms of the same gene. All genes were later confirmed to be correctly annotated by mapping them to the same proteins annotated in the *N. westbladi* transcriptome. A custom R script was written to locate the filtered genes in the GFF files and extract three metrics related to gene architecture: protein length, number of exons per protein, and average exon length per gene. Unfortunately, the GFF annotation file was not available for all these genomes, so not all of them could be included in this analysis (Supplementary Table S2). To ameliorate the misleading effect of highly fragmented genes we filtered out all proteins shorter than half of the average protein length of the respective gene (a total of 10 proteins). To test if the observed differences in the three gene metrics were statistically significant, the Shapiro-Wilk’s method and the Barlett test were used to check if they follow a normal distribution and the homogeneity of their variances, respectively. For each gene, the differences among clades were tested with either an ANOVA or a Kruskal-Wallis test, depending on the result of the normality and homoscedasticity tests. Finally, the Bonferroni correction (ANOVA) and the Dunn test (Kruskal-Wallis) were selected to run pairwise comparisons in all cases identified as statistically different. In all cases, a *p*-value of 0.05 was set as the significance threshold.

3 Results

3.1 The *Nemertoderma westbladi* genome

The best extraction was produced from a sample stored in RNAlater using the QIAamp Micro kit, obtaining a fragment size over 20 kb and ca. 20 ng of total DNA, which yielded 990 ng after DNA shearing and whole-genome amplification. About half of this DNA was selected for sequencing. A total of 2,313,071 reads were produced during HiFi sequencing, later reduced to 2,297,478 after quality filtering with an average length of 6.6 kb.



Flye produced the best assembly (Supplementary Table S1), which was 678.9 Mb long and contained 26,880 contigs, and later reduced to 458 Mb and 11,625 contigs when assembly redundancies were removed (Figure 1A). The longest contig was 2 Mb long, with an N50 of 57.1 kb and contained 82.6% of the BUSCO Metazoa odb10. The assembly contained two repeats of 507 and 531 bp with 70,000 and 79,000 copies, respectively, corresponding to 11% of the assembled (raw) genome. BlobTools2 revealed the presence of many contaminants, with only 78.9% of the contigs identified as metazoan (Supplementary Table S3). Thus, the decontaminated assembly was only 407.7 Mb, split into 9,167 contigs with an N50 of 60 kb (Figure 1A; Table 1), but 80.2% of the Metazoa and 88.6% of the Eukaryota BUSCO genes were still present (Supplementary Figure S1). The two databases reported a relatively high proportion of duplicated BUSCO genes, 7.8% and 11% in the Metazoa and Eukaryota databases, respectively (Supplementary Figure S1). The smudgeplot was markedly different before and after the decontamination step, as the inferred ploidy went from triploid to diploid after the decontamination (Supplementary Figure S2).

The genome size estimated by GenomeScope was 237.1 Mb, with an average coverage of 39.6, and high heterozygosity (3.99%), although these numbers must be taken cautiously given the poor fit of the model (39.3%; Supplementary Figure S3).

The decontaminated Illumina genome was also relatively complete, with 76.8% of the metazoan BUSCO genes present in the assembly (with barely any duplicates, 0.6%), but much shorter (62.2 Mb) and much more fragmented (49,310 contigs; N50: 4 kb) (Figure 1A). Despite being sequenced from cultured, starved and symbiont-free populations, BlobTools also identified some contaminants in the published genomes of *P. naikaiensis* and *S. roscoffensis*. The former went from 656.1 Mb and 12,525 contigs to 581.4 Mb and 7,104 contigs, whereas the latter went from 1103 Mb to 3,460 contigs to 1064.9 Mb and 2,730 contigs (Figure 1A). The N50 of the two genomes rose from 127 to 130 kb in *P. naikaiensis*, and from 1.04 to 1.08 Mb in *S. roscoffensis*. Despite the observed differences in genome size and contiguity, the four genomes show very similar completeness results. Around 90% of the Eukaryota BUSCO genes were identified in the decontaminated genomes of all

TABLE 1 Statistics of the four genomes analysed in this study after the decontamination step. The *N. westbladi* genomes are presented as “HiFi” and “Illumina” to differentiate the two sequencing approaches.

Parameter	Illumina	HiFi	Pnaikaiensis	Sroscoffensis
Length after BlobTools (Mb)	62.229	407.663	581.371	1064.926
N's (count)	49,310	12,700	7,367,142	1,589,933
N's (%)	0.079	0.003	1.267	0.149
Number of contigs	26,021	9,167	7,104	2,730
Longest contig (Kb)	65.353	601.587	702.461	8,003.794
Average contig length (Kb)	2.391	44.471	81.837	390.083
N50 (Kb)	3.996	59.976	129.752	1077.644
Number of gene models	23,120	22,578	20,303	28,513
Functionally annotated proteins	14,486	10,074	13,708	17,717
Max. number of genes per contig	33	92	37	280
Average number of genes per contig	0.876	2.352	2.858	12.281
Max. number of exons per gene	26	349 *	512	97
Average number of exons per gene	1.531	3.284	6.386	4.244

*One gene with >1200 exons, thought to be an annotation artefact, was removed for this calculation

species, with *P. naikaiensis* presenting the fewest genes (14.9% of missing genes) (Supplementary Figure S1A). Differences among genomes were slightly higher with the Metazoa database, with almost a 10% difference between the most (*S. roscoffensis*; 18.5% missing genes) and the least (*P. naikaiensis*; 27%) complete genomes. In *N. westbladi*, the HiFi genome was almost as complete as *S. roscoffensis* (19.8% missing genes), whereas the Illumina genome was in an intermediate position (23.1%) (Supplementary Figure S1B).

The number of gene models in the four genomes ranged from 20,303 (*P. naikaiensis*) to 28,513 (*S. roscoffensis*), although the differences were reduced when only functionally annotated genes were considered: 10,074 (*N. westbladi*, HiFi), 13,708 (*P. naikaiensis*), 14,486 (*N. westbladi*, Illumina), and 17,717 (*S. roscoffensis*) (Table 1). The organisation of these genes in the genome somehow reflected the differences observed in genome contiguity. In the *N. westbladi* genome sequenced with Illumina, the average number of genes per contig was just 0.876, with a single gene in almost 90% of the contigs (Figure 1B), and the contig with the highest number of genes presented 33 gene models (Table 1). In the HiFi-sequenced *N. westbladi* genome, up to 92 genes were found in a single contig, with an average of 2.4 genes per contig. Similarly, an average of 2.9 genes per contig were annotated in the *P. naikaiensis* genome, but in this case, the maximum number of genes in one contig was only 37. The *S. roscoffensis* genome stands out, with a maximum of 280 genes in a single contig and more than 10 genes in almost 40% of the contigs (Figure 1B; Table 1). This trend, however, was not observed in gene architecture. The gene models in *P. naikaiensis*, *S. roscoffensis*, and the HiFi genome of *N. westbladi* were similar, ranging between an average of 3.2–6.3 exons per gene, whereas almost all the genes presented a single exon in the Illumina

genome (average 1.5) (Figure 1D). One gene with over 1200 exons was annotated in the *N. westbladi* genome. This is likely an annotation error and it was not included in the calculation of these metrics, but it is shown in Figure 1D. The intron size was very variable in all genomes, ranging from 6 (*P. naikaiensis*) to 193,733 (*S. roscoffensis*) bp. The intron size distribution was similar between *N. westbladi* and *P. naikaiensis*, but with generally longer introns in *S. roscoffensis* (Figure 1C). Nevertheless, the intron size range was similar in the three genomes, but visibly smaller in the *N. westbladi* Illumina genome.

According to RepeatMasker, the *N. westbladi* genome is very repetitive, masking up to 59.66% of the genome (Supplementary Table S4). The majority of these repeats are interspersed throughout the genome (57.84%) and more than a fifth (22.61%) were not classified into any known repeat family. Among the classified repeats, the most common ones are retroelements (31.19%), particularly the long terminal repeats (LTR, 19.32%) and long interspersed nuclear elements (LINEs, 11.53%). The Illumina genome presents a sharp contrast, with just 16.40% of the genome masked as repetitive, although LINEs (4.28%) and LTR (3.43%) are still the most abundant repeat elements (Supplementary Table S4).

3.2 Identification of the contaminant contigs

More than half of the taxonomic groups identified within the set of contaminant contigs were bacteria, including several of the major taxonomic groups: Bacteroidetes, Tectomicrobia, Proteobacteria (including Alpha-, Beta-, Delta/Epsilon-, and Gammaproteobacteria), Planctomycetes, Actinobacteria, Cyanobacteria, and Firmicutes. None of the “Candidate Phyla Radiation” phyla (Brown et al., 2015) were

important sources of contamination besides bacteria are algae (Chlorophyta, Rhodophyta, and Streptophyta), and fungi (Ascomycota, Basidiomycota, Microsporidia, Mucoromycota, and Zoopagomycota). These groups accumulate 87% of the taxonomic diversity within the contaminants. In addition, we also found Protista (Amoebozoa, Perkinsozoa, Endomyxa, and Oomycota), and Virus (Uroviricota). A complete description of these results is provided in [Supplementary Table S5](#).

3.3 Gene content evolution

The comparison of 18 animal genomes, representing Acoelomorpha, Cnidaria, Deuterostomia, and Protostomia revealed a high degree of specificity in gene content: 17.2% of all orthogroups present in Cnidaria are exclusive to this phylum (1754 out of 10,172), 23.1% in Acoelomorpha (2,101 out of 9,080), 45.3% in Deuterostomia (7,976 out of 17,593), and 48.7% in Protostomia (9,573 out of 19,669; [Figure 2A](#)). Hence, only 33.4% (10,736 out of 32,141) of all orthogroups were annotated in at least two of the four groups. Among these, more than half (53.4%) were present in at least one species of each clade, whereas only 3.8% were present in all bilaterian clades but Cnidaria. A total of 8,418 genes were identified as shared across Metazoa (present in Cnidaria and at least one Bilateria), and 2,318 for Bilateria (present in at least two bilaterian clades). Acoelomorpha had 71.5% of the metazoan genes and 41.5% of the bilaterian ones, contrasting with deuterostomes (91.3% and 83.1%) and protostomes (94.4% and 92.9%) ([Figure 2C](#)). The proportion of missing BUSCO genes was below 11% in all four groups ([Figure 2B](#)), and so genome completeness does not explain this pattern. Within Acoelomorpha, almost half (43.2%) of the genes were shared between Acoela and Nemertodermatida ([Figure 2A](#)), 41.5% were unique to Acoela, and 15.3% to Nemertodermatida.

3.4 Ultrafiltration excretory system

Nine genes, selected because of their participation in the development of the ultrafiltration excretory system, were investigated: three structural proteins (*Nephrin*, *Kirrel*, and *ZO1*), and six transcription factors (*Eya*, *Lhx1/5*, *Osr*, *POU3*, *Sall*, and *Six1*). All of them were annotated in both protostomes and deuterostomes. In Acoelomorpha, all genes but *Osr* were annotated, whereas only three out of the nine genes were found in the two cnidarian species (*ZO1*, *Six*, and *Lhx*; [Figure 3A](#)). According to GenBank, three more genes (*Nephrin*, *Eya*, and *POU3*) are also present in this phylum ([Figure 3A](#)).

The gene architecture (in terms of protein length, number of exons per gene, and average exon length) was compared for the nine genes among four clades: Cnidaria, Acoelomorpha, Deuterostomia, and Protostomia. Almost half of the 27 comparisons returned statistically significant differences among clades, most of them related to acoelomorphs ([Figure 3B](#)). Despite the evident variation in protein length, both within and among clades, only three out of the nine genes were considered to be statistically significant: *Kirrel*, which is significantly longer in acoelomorphs; *ZO1*, longer in deuterostomes; and *Lhx*, but in this case the

differences were only significant between acoelomorphs (longer) and protostomes (shorter). As for the number of exons per gene, *ZO1* and *Eya* presented fewer exons in acoelomorphs than in both deuterostomes and protostomes. Finally, the last gene with a significantly different number of exons is *POU3*. This is a relatively short protein, on average shorter than 500 amino acids in all clades, and with very few exons: only one exon in all deuterostomes but *Branchiostoma floridae* (three), between one and three in protostomes, and between one and four in acoelomorphs. Only the differences between deuterostomes and acoelomorphs were statistically significant. Two remarkable outliers were found when comparing the number of exons per gene. Three chordate *ZO1* sequences were divided into more than 80 exons (average 29.5) and one of the *POU3* sequences annotated in *P. naikaiensis* presented 15 exons (average in Acoelomorpha: 2.6). Nonetheless, these proteins were roughly of the same size as the others and their identity to the most similar protein was above 90%.

In an attempt to avoid the misleading effect of errors in the annotation (partial proteins will be generally shorter and with fewer exons), the average exon length was also considered. In this case, six out of the nine proteins were significantly different among clades. The average exon length was significantly longer in acoelomorphs in three genes (*Kirrel*, *Eya*, and *Lhx*), and two in deuterostomes (*Sall* and *Osr*, although the latter was only present in deuterostomes and protostomes). The only instance with significantly shorter exon lengths is the protostome's *ZO1* gene. Finally, among the nine comparisons including at least one cnidarian species (three genes, three metrics) no significant differences were found but in the average exon length of *Lhx*, which is significantly shorter than that of acoelomorphs, as also observed in deuterostomes and protostomes.

4 Discussion

4.1 Performance of the Ultra-Low DNA Input protocol for sequencing large genomes

The steady development of sequencing technologies is allowing the generation of genomes spanning the diversity of life, which now includes minute organisms. Indeed, thanks to the latest low and ultra-low DNA input protocols, sequencing high-quality genomes from millimetric animals is now possible ([Yoshida et al., 2018](#); [Schneider et al., 2021](#); [Lord et al., 2023](#)). In this study, we used the Pacbio Ultra-Low DNA Input protocol to sequence the genome of *N. westbladi*, reporting the first nemertodermatid genome, sequenced from a single microscopic worm. The estimated genome length is shorter than in any sequenced acoel, and considerably shorter than *S. roscoffensis* and *H. miamia* ([Gehrke et al., 2019](#)). Although the *P. naikaiensis* genome is slightly more contiguous than *N. westbladi*, all the metrics compared are similar between the two genomes. In contrast, both *S. roscoffensis* and *H. miamia* were scaffolded using proximity ligation data, and hence both show much higher contiguity. Beyond the differences in contiguity, annotation metrics are comparable among *N. westbladi*, *P. naikaiensis*, and *S. roscoffensis*. In this case, *N. westbladi* is more similar to *S. roscoffensis* than to *P. naikaiensis*,

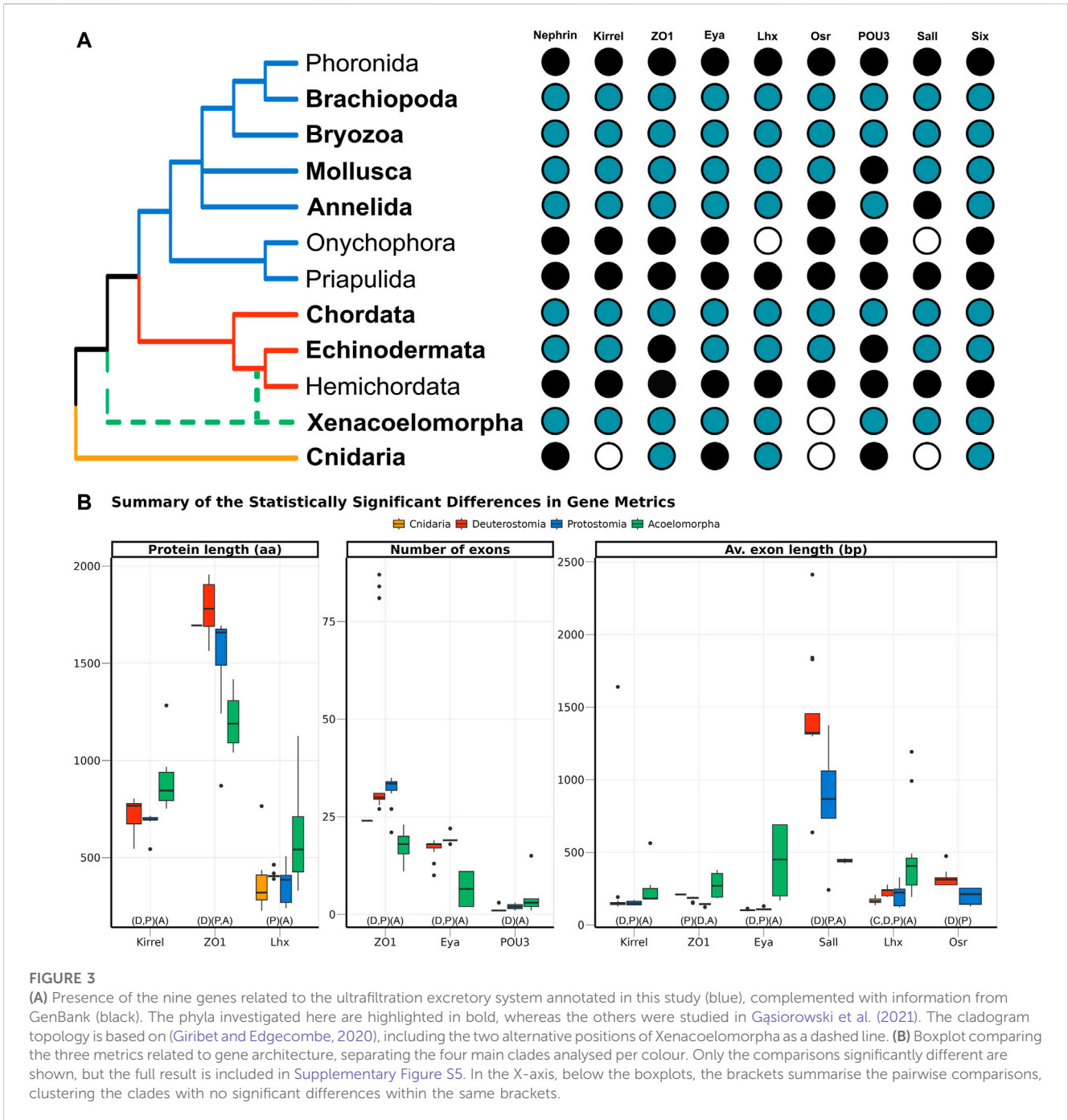


FIGURE 3

(A) Presence of the nine genes related to the ultrafiltration excretory system annotated in this study (blue), complemented with information from GenBank (black). The phyla investigated here are highlighted in bold, whereas the others were studied in Gąsiorowski et al. (2021). The cladogram topology is based on (Giribet and Edgecombe, 2020), including the two alternative positions of Xenacoelomorpha as a dashed line. (B) Boxplot comparing the three metrics related to gene architecture, separating the four main clades analysed per colour. Only the comparisons significantly different are shown, but the full result is included in Supplementary Figure S5. In the X-axis, below the boxplots, the brackets summarise the pairwise comparisons, clustering the clades with no significant differences within the same brackets.

which shows the lowest genome completeness and number of gene models. In particular, the analysis of gene architecture shows that the number of exons per gene and intron size is also comparable, likely meaning that the annotated proteins are complete or nearly complete, facilitating the study of gene properties, such as intron-exon structure. Likewise, all genomes are similarly repetitive: *N. westbladi* 59.66%; *P. naikaiensis* 69.8%; *S. roscoffensis* 61.14%; and *H. miamia* 53%, but this is where the difference between the short- and long-read genomes of *N. westbladi* strikes the most. Although they have similar completeness and number of gene models, the Illumina genome is only 62.2 Mb long and only 16.4% repeats,

which is probably explained by the difficulty to assemble repetitive areas of the genome (Tørresen et al., 2019).

It is obvious from the comparisons above that achieving a highly contiguous genome from single-millimetre worms is still challenging. One potential explanation for this is genome size. Although the *N. westbladi* genome is shorter than the maximum genome size advised by PacBio, due to the contaminants the raw assembly is still almost 700 Mb long. The Ultra-Low DNA Input protocol has insofar been tested in animals whose genome size ranges between 200 and 300 Mb, returning significantly more contiguous genomes than that of *N. westbladi* (Kingan et al.,

2019; Korklach, 2020; Schneider et al., 2021). The generally lower coverage of the raw assembly of the nemertodermatid genome, due to its larger size, could have resulted in a more fragmented assembly. Yet, sequencing a second HiFi SMRT cell was not feasible due to the low DNA yield. This protocol has been successfully tested on *C. elegans* and freshwater flatworms, demonstrating its good performance (Laumer, pers. comm.). However, it has never been used for sequencing >300 Mb genomes. Thus, it would be interesting to sequence the genome from a novel species with a similar genome size to confirm low coverage, and not unexpected artefacts, explain the fragmentation observed. Alternatively, one straightforward solution to improve genome contiguity is complementing this approach with ligation data, which has shown great results both in *S. roscoffensis* and *H. miamia* (Gehrke et al., 2019; Martinez et al., 2022). However, this approach would require pooling tens of individuals to obtain the required amount of DNA, which is not feasible for all animals. *N. westbladi* cannot be cultured in the lab and collecting worms in enough numbers is challenging. Interestingly, the *P. naikaiensis* genome (the most similar to *N. westbladi*) was sequenced from a pool of individuals in 52 SMRT Cells (Arimoto et al., 2019), whereas the *N. westbladi* genome comes from a single worm and one HiFi SMRT Cell. Altogether, these results highlight the potential of combining this protocol and HiFi reads to generate good-quality genomes from single, microscopic organisms, even for relatively large genomes.

The BlobTools analysis identified a high degree of contamination in the raw assembly of *N. westbladi*, which is to be expected from a microscopic organism caught in the wild. Although *N. westbladi* is known to not carry internal symbionts (based on hundreds of observations), a TEM analysis revealed the presence of gram-negative bacteria throughout the epidermal cilia (Lundin, 1998). All bacteria identified during the decontamination are gram-negative, matching this observation. Thus far, DNA extraction was performed from a whole specimen, thus sequencing the gut microbiome, and other contaminants might have been transferred from the DNA suspended in the seawater. A common practice to limit the presence of contaminants in the organism, and that was applied to both *P. naikaiensis* and *S. roscoffensis*, but not *N. westbladi*, is to starve the animals before DNA extraction. Besides, the two acoel genomes were sequenced from juveniles, before they incorporate the symbiotic algae, and rinsed with filtered seawater (Arimoto et al., 2019; Martinez et al., 2022). However, as seen here this is not enough to prevent the presence of contaminants. This was particularly problematic in the case of *P. naikaiensis*, as almost 4% of the contigs (75 Mb, over 10% of the genome) were identified as bacterial contigs. It is important to notice that a big fraction of the genomes did not have any hit against the Uniprot database (*N. westbladi* 18.1%; *P. naikaiensis* 8.4%; *S. roscoffensis* 1.9%; Supplementary Table S3), showing the importance of sequencing underrepresented groups to improve the reference databases.

4.2 Evolution of Acoelomorpha genomes

The increasing availability of animal genomes has unveiled a remarkable diversity in genome sizes, ranging from 15.3 Mb in the

orthonectid *Intoshia variabilis* to the 43 Gb of the lungfish genome (Slyusarev et al., 2020; Meyer et al., 2021). It has been observed that miniaturised animals tend to have smaller genomes, which has been noted both in vertebrates and invertebrates (Liu et al., 2012; Decena-Segarra et al., 2020; Xu et al., 2021), but with notable exceptions to this rule, as observed in nematodes and platyhelminths (Consortium, 2019). Genome length in the latter ranges between 700 and 1200 Mb, the same size range as birds, some gastropods, and many freshwater fish, among others (Zhang et al., 2014; Nam et al., 2017; Yuan et al., 2018). Similarly, acoelomorph genomes vary between 407 (the genome of *N. westbladi* is the shortest reported for Acoelomorpha) and 1059 Mb but contrast with the chromosome-level genome of *X. bocki*, estimated at 110 Mb (Schiffer et al., 2023). Comparisons of eukaryotic genomes proposed that variations in genome sizes and proportion of repeat elements are correlated (Elliott and Gregory, 2015; Shah et al., 2020), which might also apply within Xenacoelomorpha. Acoelomorph genomes show a much higher repeat content than the small genome of *Xenoturbella* (Schiffer et al., 2023).

In turn, acoelomorph genomes seem to be characterised by an important reduction of gene content. Indeed, almost 60% of the genes shared between protostomes and deuterostomes are missing in acoelomorphs, which could be explained by the morphological simplicity of these worms compared with other bilaterians, but the evolutionary interpretation depends on the phylogenetic hypothesis. Under the Xenambulacraria hypothesis, their absence must be explained by massive secondary losses. The Nephrozoa hypothesis, on the other hand, suggests that the evolution of the genes exclusively shared by deuterostomes and protostomes occurred in the stem line of Nephrozoa and no *ad hoc* hypotheses of gene loss are required.

4.3 Evolution of the genes related to the ultrafiltration excretory system

Despite the absence of a specialised excretory system in Xenacoelomorpha, Andrikou et al. (2019) detected active excretion in this group through the digestive tissue and annotated several genes known to participate in the excretory mechanisms of nephrozoan animals. Here, we annotated in the genomes of Acoela and Nemertodermatida seven of the nine genes involved in the development of the nephridia and one more (Sall) in Acoela. Regardless of their phylogenetic position, whether as a sister to Ambulacraria or Nephrozoa, the presence of these genes might be explained by their participation in other important functions. A spatial transcriptomics analysis in the acoel *Isodiametra pulchra* and the nemertodermatid *Meara stichopi* located the expression of *Nephrin* in the brain and the nerve cords (Andrikou et al., 2019), which resembles observations in mammals and *Drosophila*, the latter through the *Nephrin* homolog *Sns* (Putala et al., 2000; Putala et al., 2001; Bali et al., 2022). In contrast, no homologs to the *Osr* gene (named *Odd* in *Drosophila*) could be annotated in any of the acoelomorph genomes. A BLAST search over the two *Xenoturbella* transcriptomes failed to annotate this gene in these species, confirming its absence is a general trait of Xenacoelomorpha. This is noteworthy, as *Osr* is essential in the formation of the excretory organs: in vertebrates, it participates in the formation

of the pronephros, the first stage in kidney formation, and its knock-out results in the absence of kidneys (James et al., 2006); whereas in *Drosophila*, *Odd* participates in the embryogenesis of the tubules of Malpigi (Tena et al., 2007). Overall, it seems that the molecular machinery that participates in the functioning of a complex ultrafiltration excretory system is present in acoelomorphs, but they lack the one gene necessary to promote the formation of discrete excretory organs.

This pattern fits well within the Nephrozoa hypothesis. In this scenario, the origin of the molecular machinery associated with the excretory organs would be the result of gene co-option, a common phenomenon in the origin of key innovations, such as the development of the radula and shell evolution in molluscs (Hilgers et al., 2018) or the multiple origins of cnidarian eyes (Picciani et al., 2018). Interestingly, six of the nine genes investigated have been annotated in different cnidarian species, strengthening the idea of these genes pre-dating the appearance of this specialised excretory system (Gąsiorowski et al., 2021). Thus far, *Osr* has not been annotated in any phylum outside of Nephrozoa, supporting the origin of this gene in the ancestor of this clade. Nevertheless, given the ongoing debate around the phylogenetic position of xenacoelomorphs, the Xenambulacraria hypothesis also needs to be taken into consideration. If Xenacoelomorpha is the sister group of Ambulacraria, additional *ad hoc* hypotheses have to be invoked: either the *Osr* gene was independently gained in Protostomia, Ambulacraria, and Chordata or it was lost in Xenacoelomorpha. The *Drosophila Odd* gene has been shown to activate the formation of kidney tissue in vertebrates (Tena et al., 2007), which suggests a common origin of both genes in protostomes and deuterostomes. Likewise, the function of this gene is not limited to the development of the excretory organs, but it participates in the development of the foregut in vertebrates (Han et al., 2017) and it is known to be expressed in the digestive tract of spiralian and hemichordates (Gąsiorowski et al., 2021). Although its general anatomy is variable, the presence of a sack-like gut is considered a plesiomorphy within Xenacoelomorpha (Gavilán et al., 2019) and the involvement of *Osr* in its development could be expected. In this light, the reduction of the excretory organs alone would not explain the secondary loss of *Osr*, as it would need to be completely nonfunctionalized before that. Regardless of which of the two phylogenetic hypotheses is true, the acquisition of *Osr* and the development of a discrete excretory system seem to be interconnected. The two are a nephrozoan novelty and their origin likely dates back to the most recent common ancestor of protostomes and deuterostomes. More importantly, knock-out experiments demonstrate that there are no other current transcription factors capable of replacing *Osr* in the early formation of the pronephros (James et al., 2006), which suggests this gene must have assumed this function rather early.

We found statistically significant differences in the gene architecture of all genes but *Nephrin* and *Six*, six of them related to the average exon length. Acoelomorpha is responsible for two-thirds of the differences observed, which fits with the co-option of these genes into the development of the excretory system in the ancestor of Nephrozoa. Changes in gene structure are a strong generator of diversity, particularly after gene duplication, as part of the neofunctionalization of proteins (Xu et al., 2012). Alternatively, the

differences observed might simply be explained by changes in the selective pressures during the acquisition or the reduction of this system, something that might be supported by the observations in Bryozoa. Within protostomes, Bryozoa, which also lack an excretory system, is responsible for most of the variation observed. Notably, half of the gene metrics that are visibly different in this phylum are shared with acoelomorphs: *ZO1* and *Lhx* length, *ZO1* number of exons, and *Sall* average exon length. However, the variation does not always go in the same direction (e.g., the number of exons in *ZO1* increases in Acoelomorpha, but decreases in Bryozoa), likely because the absence of the excretory organs in the two animal groups represents two independent evolutionary events. Some authors have argued that the rapid evolutionary rates observed in Acoelomorpha might be associated with other traits observed in this group, such as chromosomal rearrangements or changes in gene content, misleading comparative analyses and making *Xenoturbella* a better model for studying the evolution of Xenacoelomorpha (Philippe et al., 2019; Schiffer et al., 2023). Unfortunately, the genomic data of *X. bocki* is yet not available so we have inferred a gene tree for each of the nine genes analysed and compared the differences in branch lengths among clades to explore this possibility (Supplementary Figure S4). Although branch lengths are indeed significantly longer in acoelomorphs than in any other clade (except in *Lhx* and *Six*), they are also longer in deuterostomes compared to protostomes despite the similarities between the two clades. In more detail, protostomes present the shortest branches in the gene trees, while Bryozoa is one of the phyla with the most changes in gene architecture. Hence, the accelerated evolutionary rates of Acoelomorpha do not seem to be the main factor underlying the differences observed in these genes, although it would be interesting to confirm this once all the data from the *Xenoturbella* genome is publicly available.

5 Conclusion

In this study, we have generated the first draft of a nemertodermatid genome, sequenced from a single, microscopic individual using the Ultra-Low DNA Input protocol and HiFi. We show that this approach is capable of producing genomes of relatively good quality even from small organisms with large genomes. The main drawback is genome contiguity, which remains the main challenge and one of the avenues in genome sequencing that need the most attention. Nevertheless, genome quality is good enough to annotate full proteins, allowing detailed analysis of gene architecture. We prove this by analysing the genes related to the ultrafiltration excretory system. We observe that the molecular machinery related to this system predates its origin, as most of the genes were present in Urbilateria or even in the cnidarian-bilateria ancestor. Interestingly, all genes but *Osr*, the one gene triggering the formation of these organs, were annotated in Xenacoelomorpha. Thus far, gene architecture is markedly different in Acoelomorpha, which cannot be explained either by the accelerated evolution of this clade or the lack of the excretory system alone. All these findings are more easily explained under the Nephrozoa hypothesis.

Data availability statement

The datasets presented in this study can be found in online repositories. The names of the repository/repositories and accession number(s) can be found below: <https://www.ncbi.nlm.nih.gov/genbank/>, GenBank Bioproject PRJNA981986 https://figshare.com/projects/2023_Nemertoderma_westbladi_genome/169818.

Author contributions

SA, OV, and UJ conceived the project; SA performed DNA extractions; JH was responsible for library preparations and sequencing; CT-R carried out the post-sequencing analyses, from quality filtering to genome assembly; SA decontaminated and annotated the genome and performed comparative analyses; SA and UJ led the writing of the manuscript. All authors contributed to the article and approved the submitted version.

Funding

This project was funded by the VR project 2018-05191, granted to UJ, and the “2021 Riksmusei Vänner” and “2020 Helge Ax:son Johnsons stiftelse” stipends to SA. Work performed at NGI/Uppsala Genome Center has been funded by RFI/VR and Science for Life Laboratory, Sweden. Analyses and data handling were enabled by resources in projects SNIC 2020/15-191 and SNIC 2021/22-562 provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS) at UPPMAX, funded by the Swedish Research Council through grant agreement no. 2018-05191.

References

- Albertin, C. B., Simakov, O., Mitros, T., Wang, Z. Y., Pungor, J. R., Edsinger-Gonzales, E., et al. (2015). The octopus genome and the evolution of cephalopod neural and morphological novelties. *Nature* 524, 220–224. doi:10.1038/nature14668
- Alonge, M., Lebeigle, L., Kirsche, M., Aganezov, S., Wang, X., Lippman, Z. B., et al. (2022). Automated assembly scaffolding using RagTag elevates a new tomato system for high-throughput genome editing. *Genome Biol.* 23, 258. doi:10.1186/s13059-022-02823-7
- Andrikou, C., Thiel, D., Ruiz-Santesteban, J. A., and Hejnol, A. (2019). Active mode of excretion across digestive tissues predates the origin of excretory organs. *PLoS Biol.* 17, e3000408–e3000422. doi:10.1371/journal.pbio.3000408
- Arimoto, A., Hikosaka-Katayama, T., Hikosaka, A., Tagawa, K., Inoue, T., Ueki, T., et al. (2019). A draft nuclear-genome assembly of the acoel flatworm *Praesagittifera naikaiensis*. *Gigascience* 8, 1–8. doi:10.1093/gigascience/giz023
- Bali, N., Lee, H. K., and Zinn, K. (2022). Sticks and Stones, a conserved cell surface ligand for the Type IIa RPTP Lar, regulates neural circuit wiring in *Drosophila*. *Elife* 11, 714699–e71530. doi:10.7554/eLife.71469
- Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., et al. (2012). SPAdes: A new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* 19, 455–477. doi:10.1089/cmb.2012.0021
- Boscaro, V., Holt, C. C., Van Steenkiste, N. W. L., Herranz, M., Irwin, N. A. T., Álvarez-Campos, P., et al. (2022). Microbiomes of microscopic marine invertebrates do not reveal signatures of phyllosymbiosis. *Nat. Microbiol.* 7, 810–819. doi:10.1038/s41564-022-01125-9
- Brown, C. T., Hug, L. A., Thomas, B. C., Sharon, I., Castelle, C. J., Singh, A., et al. (2015). Unusual biology across a group comprising more than 15% of domain Bacteria. *Nature* 523, 208–211. doi:10.1038/nature14486
- Brüna, T., Hoff, K. J., Lomsadze, A., Stanke, M., and Borodovsky, M. (2021). BRAKER2: automatic eukaryotic genome annotation with GeneMark-ep+ and AUGUSTUS supported by a protein database. *Nar. Genomics Bioinforma.* 3, 1–11. doi:10.1093/nargab/lqaa108
- Brüna, T., Lomsadze, A., and Borodovsky, M. (2020). GeneMark-EP+: eukaryotic gene prediction with self-training in the space of genes and proteins. *Nar. Genomics Bioinforma.* 2, lqaa026–14. doi:10.1093/nargab/lqaa026
- Buchfink, B., Reuter, K., and Drost, H. G. (2021). Sensitive protein alignments at tree-of-life scale using DIAMOND. *Nat. Methods* 18, 366–368. doi:10.1038/s41592-021-01101-x
- Cannon, J. T., Vellutini, B. C., Smith, J., Ronquist, F., Jondelius, U., and Hejnol, A. (2016). Xenacoelomorpha is the sister group to Nephrozoa. *Nature* 530, 89–93. doi:10.1038/nature16520
- Challis, R., Richards, E., Rajan, J., Cochrane, G., and Blaxter, M. (2020). BlobToolKit - interactive quality assessment of genome assemblies. *G3 Genes, Genomes, Genet.* 10, 1361–1374. doi:10.1534/g3.119.400908
- Cheng, H., Concepcion, G. T., Feng, X., Zhang, H., and Li, H. (2021). Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* 18, 170–175. doi:10.1038/s41592-020-01056-5
- Consortium, I. H. G. (2019). Comparative genomics of the major parasitic worms. *Nat. Genet.* 51, 163–174. doi:10.1038/s41588-018-0262-1
- Crisuolo, A., and Gribaldo, S. (2010). BMGE (block mapping and gathering with entropy): A new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* 10, 210. doi:10.1186/1471-2148-10-210
- Decena-Segarra, L. P., Bizjak-Mali, L., Kladnik, A., Sessions, S. K., and Rovito, S. M. (2020). Miniaturization, genome size, and biological size in a diverse clade of salamanders. *Am. Nat.* 196, 634–648. doi:10.1086/711019
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., et al. (2013). Star: ultrafast universal RNA-seq aligner. *Bioinformatics* 29, 15–21. doi:10.1093/bioinformatics/bts635
- Dos Reis, M., Thawornwattana, Y., Angelis, K., Telford, M. J., Donoghue, P. C. J., and Yang, Z. (2015). Uncertainty in the timing of origin of animals and the limits of precision in molecular timescales. *Curr. Biol.* 25, 2939–2950. doi:10.1016/j.cub.2015.09.066

Acknowledgments

We are thankful to C. Laumer (NHM) for his advice during the early stages of this project. The authors would like to acknowledge support of the National Genomics Infrastructure (NGI)/Uppsala Genome Center for providing assistance in massive parallel sequencing and computational infrastructure.

Conflict of interest

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Supplementary material

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fgene.2023.1244493/full#supplementary-material>

- Dunwell, T. L., Paps, J., and Holland, P. W. H. (2017). Novel and divergent genes in the evolution of placental mammals. *Proc. R. Soc. B Biol. Sci.* 284, 20171357. doi:10.1098/rspb.2017.1357
- Elliott, T. A., and Gregory, T. R. (2015). What's in a genome? The C-value enigma and the evolution of eukaryotic genome content. *Philos. Trans. R. Soc. B Biol. Sci.* 370, 20140331. doi:10.1098/rstb.2014.0331
- Emms, D. M., and Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biol.* 20, 238. doi:10.1186/s13059-019-1832-y
- Fu, L., Niu, B., Zhu, Z., Wu, S., and Li, W. (2012). CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28, 3150–3152. doi:10.1093/bioinformatics/bts565
- Gąsiorowski, L., Andrikou, C., Janssen, R., Bump, P., Budd, G. E., Lowe, C. J., et al. (2021). Molecular evidence for a single origin of ultrafiltration-based excretory organs. *Curr. Biol.* 31, 3629–3638.e2. doi:10.1016/j.cub.2021.05.057
- Gavilán, B., Sprecher, S. G., Hartenstein, V., and Martínez, P. (2019). The digestive system of xenacoelomorphs. *Cell Tissue Res.* 377, 369–382. doi:10.1007/s00441-019-03038-2
- Gehrke, A. R., Neverett, E., Luo, Y. J., Brandt, A., Ricci, L., Hulett, R. E., et al. (2019). Acoel genome reveals the regulatory landscape of whole-body regeneration. *Science* 80, 363. doi:10.1126/science.aau6173
- Giribet, G., and Edgecombe, G. D. (2020). *The invertebrate tree of life*. Princeton, New Jersey: Princeton University Press. doi:10.2307/j.ctvscxrhnm
- Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., et al. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat. Biotechnol.* 29, 644–652. doi:10.1038/nbt.1883
- Gross, V., Treffkorn, S., Reichelt, J., Eppl, L., Lüter, C., and Mayer, G. (2019). Miniaturization of tardigrades (water bears): morphological and genomic perspectives. *Arthropod Struct. Dev.* 48, 12–19. doi:10.1016/j.asd.2018.11.006
- Guimaraes, K. (2018). *DNA extraction (Salting out) V.4*. California: Protocols.io. doi:10.17504/protocols.io.vwfe7bn
- Gurevich, A., Saveliev, V., Vyahhi, N., and Tesler, G. (2013). Quast: quality assessment tool for genome assemblies. *Bioinformatics* 29, 1072–1075. doi:10.1093/bioinformatics/btt086
- Han, L., Xu, J., Grigg, E., Slack, M., Chaturvedi, P., Jiang, R., et al. (2017). Osr1 functions downstream of Hedgehog pathway to regulate foregut development. *Dev. Biol.* 427, 72–83. doi:10.1016/j.ydbio.2017.05.005
- Haszprunar, G. (2016). Review of data for a morphological look on Xenacoelomorpha (Bilateria incertae sedis). *Org. Divers. Evol.* 16, 363–389. doi:10.1007/s13127-015-0249-z
- Hilgers, L., Hartmann, S., Hofreiter, M., and von Rintelen, T. (2018). Novel genes, ancient genes, and gene Co-option contributed to the genetic basis of the radula, a Molluscan innovation. *Mol. Biol. Evol.* 35, 1638–1652. doi:10.1093/molbev/msy052
- James, R. G., Kamei, C. N., Wang, Q., Jiang, R., Schulthesis, T. M., and Schultheiss, T. M. (2006). Odd-skipped related 1 is required for development of the metanephric kidney and regulates formation and differentiation of kidney precursor cells. *Dev. Dis.* 133, 2995–3004. doi:10.1242/dev.02442
- Jondelius, U., Ruiz-Trillo, I., Bagaña, J., and Riutort, M. (2002). The Nemertodermatida are basal bilaterians and not members of the Platyhelminthes. *Zool. Scr.* 31, 201–215. doi:10.1046/j.1463-6409.2002.00090x
- Juravel, K., Porras, L., Höhna, S., Pisani, D., and Wörheide, G. (2023). Exploring genome gene content and morphological analysis to test recalcitrant nodes in the animal phylogeny. *PLoS One* 18, e0282444. doi:10.1371/journal.pone.0282444
- Kapli, P., and Telford, M. J. (2020). Topology-dependent asymmetry in systematic errors affects phylogenetic placement of Ctenophora and Xenacoelomorpha. *Sci. Adv.* 6, eabc5162–12. doi:10.1126/sciadv.abc5162
- Katoh, K., and Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* 30, 772–780. doi:10.1093/molbev/mst010
- Kim, D., Paggi, J. M., Park, C., Bennett, C., and Salzberg, S. L. (2019). Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* 37, 907–915. doi:10.1038/s41587-019-0201-4
- Kingan, S. B., Heaton, H., Cudini, J., Lambert, C. C., Baybayan, P., Galvin, B. D., et al. (2019). A high-quality de novo genome assembly from a single mosquito using pacbio sequencing. *Genes (Basel)* 10, 62. doi:10.3390/genes10010062
- Kolmogorov, M., Yuan, J., Lin, Y., and Pevzner, P. A. (2019). Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* 37, 540–546. doi:10.1038/s41587-019-0072-8
- Korlach, J. (2020). *A high-quality PacBio insect genome from 5 ng of input*. DNA.
- Kück, P., and Longo, G. C. (2014). FASconCAT-G: extensive functions for multiple sequence alignment preparations concerning phylogenetic studies. *Front. Zool.* 11, 81–88. doi:10.1186/s12983-014-0081-x
- Langmead, B., and Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat. Methods* 9, 357–359. doi:10.1038/nmeth.1923
- Lefort, V., Desper, R., and Gascuel, O. (2015). FastME 2.0: A comprehensive, accurate, and fast distance-based phylogeny inference program. *Mol. Biol. Evol.* 32, 2798–2800. doi:10.1093/molbev/msv150
- Li, H. (2021). New strategies to improve minimap2 alignment accuracy. *Bioinformatics* 37, 4572–4574. doi:10.1093/bioinformatics/btab705
- Liu, S., Hui, T. H., Tan, S. L., and Hong, Y. (2012). Chromosome evolution and genome miniaturization in minifish. *PLoS One* 7, e37305–e37307. doi:10.1371/journal.pone.0037305
- Lord, A., Cunha, T. J., de Medeiros, B. A. S., Sato, S., Khost, D. E., Sackton, T. B., et al. (2023). Expanding on our knowledge of ecdysozoan genomes, a contiguous assembly of the meiofaunal prapulan Tubiluchus corallicola. *Genome Biol. Evol.* 15 (6) evad103. doi:10.1093/gbe/evad103
- Lundin, K. (1998). Symbiotic bacteria on the epidermis of species of the Nemertodermatida (Platyhelminthes, Acoelomorpha). *Acta Zool.* 79, 187–191. doi:10.1111/j.1463-6395.1998.tb01157x
- Marçais, G., and Kingsford, C. (2011). A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 27, 764–770. doi:10.1093/bioinformatics/btr011
- Martín-Durán, J. M., Pang, K., Børve, A., Lè, H. S., Furu, A., Cannon, J. T., et al. (2018). Convergent evolution of bilaterian nerve cords. *Nature* 553, 45–50. doi:10.1038/nature25030
- Martinez, P., Ustyantsev, K., Biryukov, M., Mouton, S., Glasenburg, L., Sprecher, S. G., et al. (2022). Genome assembly of the acoel flatworm *Symsagittifera roscoffensis*, a model for research on body plan evolution and photosymbiosis. *G3 Genes/Genomes/Genetics* 13. doi:10.1093/g3journal/gkac336
- Meyer, A., Schloissnig, S., Franchini, P., Du, K., Woltering, J. M., Irisarri, I., et al. (2021). Giant lungfish genome elucidates the conquest of land by vertebrates. *Nature* 590, 284–289. doi:10.1038/s41586-021-03198-8
- Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., Von Haeseler, A., et al. (2020). IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol. Biol. Evol.* 37, 1530–1534. doi:10.1093/molbev/msaa015
- Mistry, J., Bateman, A., and Finn, R. D. (2007). Predicting active site residue annotations in the Pfam database. *BMC Bioinforma.* 8, 298–314. doi:10.1186/1471-2105-8-298
- Nam, B. H., Kwak, W., Kim, Y. O., Kim, D. G., Kong, H. J., Kim, W. J., et al. (2017). Genome sequence of pacific abalone (*Haliotis discus hannai*): the first draft genome in family haliotidae. *Gigascience* 6, 1–8. doi:10.1093/gigascience/gix014
- Nguyen, L. T., Schmidt, H. A., von Haeseler, A., and Minh, B. Q. (2015). IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* 32, 268–274. doi:10.1093/molbev/msu300
- Philippe, H., Poustka, A. J., Chiodin, M., Hoff, K. J., Dessimoz, C., Tomiczek, B., et al. (2019). Mitigating anticipated effects of systematic errors supports sister-group relationship between Xenacoelomorpha and Ambulacraria. *Curr. Biol.* 29, 1818–1826. doi:10.1016/j.cub.2019.04.009
- Picciani, N., Kerlin, J. R., Sierra, N., Swafford, A. J. M., Ramirez, M. D., Roberts, N. G., et al. (2018). Prolific origination of eyes in Cnidaria with Co-option of non-visual opsins. *Curr. Biol.* 28, 2413–2419. doi:10.1016/j.cub.2018.05.055
- Putala, H., Sainio, K., Sariola, H., and Tryggvason, K. (2000). Primary structure of mouse and rat nephrin cDNA and structure and expression of the mouse gene. *J. Am. Soc. Nephrol.* 11, 991–1001. doi:10.1681/asn.v116991
- Putala, H., Soininen, R., Kilpeläinen, P., Wartiovaara, J., and Tryggvason, K. (2001). The murine nephrin gene is specifically expressed in kidney, brain and pancreas: inactivation of the gene leads to massive proteinuria and neonatal death. *Hum. Mol. Genet.* 10, 1–8. doi:10.1093/hmg/10.1.1
- Ranallo-Benavidez, T. R., Jaron, K. S., and Schatz, M. C. (2020). GenomeScope 2.0 and Smudgeplot for reference-free profiling of polyploid genomes. *Nat. Commun.* 11, 1432. doi:10.1038/s41467-020-14998-3
- Rubin, B. E. R., Jones, B. M., Hunt, B. G., and Kocher, S. D. (2019). Rate variation in the evolution of non-coding DNA associated with social evolution in bees. *Philos. Trans. R. Soc. B Biol. Sci.* 374, 20180247. doi:10.1098/rstb.2018.0247
- Ruiz-Trillo, I., Riutort, M., Timothy, D., Littlewood, J., Herniou, E. A., and Bagaña, J. (1999). Acoel flatworms: earliest extant bilaterian metazoans, not members of platyhelminthes. *Sci. (80-)* 283, 1919–1923. doi:10.1126/science.283.5409.1919
- Sarmashghi, S., Bohmann, K., Thomas, P., Gilbert, M., Bafna, V., and Mirarab, S. (2017). Assembly-free and alignment-free sample identification using genome skims. *Genome Biol.* 20, 1–20. doi:10.1101/230409
- Schiffer, P. H., Natsidis, P., Leite, D. J., Robertson, H. E., Lapraz, F., Marlétaz, F., et al. (2023). *The slow evolving genomes of the xenacoelomorph worm Xenoturbella bocki*. bioRxiv. doi:10.1101/2022.06.24.497508
- Schmieder, R., and Edwards, R. (2011). Quality control and preprocessing of metagenomic datasets. *Bioinformatics* 27, 863–864. doi:10.1093/bioinformatics/btr026
- Schneider, C., Woehle, C., Greve, C., D'Haese, C. A., Wolf, M., Hiller, M., et al. (2021). Two high-quality de novo genomes from single ethanol-preserved specimens of tiny metazoans (Collembola). *Gigascience* 10, giab035–12. doi:10.1093/gigascience/giab035
- Seppy, M., Manni, M., and Zdobnov, E. M. (2019). “BUSCO: assessing genome assembly and annotation completeness,” in *Gene prediction. Methods in molecular biology* (New York: Humana Press).
- Shah, A., Hoffman, J. I., and Schielzeth, H. (2020). Comparative analysis of genomic repeat content in gomphocerine grasshoppers reveals expansion of satellite DNA and

- helitrons in species with unusually large genomes. *Genome Biol. Evol.* 12, 1180–1193. doi:10.1093/GBE/EVAA119
- Slyusarev, G. S., Starunov, V. V., Bondarenko, A. S., Zorina, N. A., and Bondarenko, N. I. (2020). Extreme genome and nervous system streamlining in the invertebrate parasite *Intoshia variabili*. *Curr. Biol.* 30, 1292–1298. doi:10.1016/j.cub.2020.01.061
- Smit, A. F. A., Hubley, R., and Green, P. (2015). *RepeatMasker open-4.0*.
- Smit, A. F. A., and Hubley, R. (2015). *RepeatModeler open-1.0*.
- Tena, J. J., Neto, A., de la Calle-Mustienes, E., Bras-Pereira, C., Casares, F., and Gómez-Skarmeta, J. L. (2007). Odd-skipped genes encode repressors that control kidney development. *Dev. Biol.* 301, 518–531. doi:10.1016/j.ydbio.2006.08.063
- Thalén, F., Kocot, K. M., and Haddock, S. (2021). *PhyloPyPruner: Tree-based orthology inference for phylogenomics*.
- Torresen, O. K., Star, B., Mier, P., Andrade-Navarro, M. A., Bateman, A., Jarnot, P., et al. (2019). Tandem repeats lead to sequence assembly errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res.* 47, 10994–11006. doi:10.1093/nar/gkz841
- Xu, G., Guo, C., Shan, H., and Kong, H. (2012). Divergence of duplicate genes in exon-intron structure. *Proc. Natl. Acad. Sci. U. S. A.* 109, 1187–1192. doi:10.1073/pnas.1109047109
- Xu, H., Ye, X., Yang, Y., Yang, Y., Sun, Y. H., Mei, Y., et al. (2021). Comparative genomics sheds light on the convergent evolution of miniaturized wasps. *Mol. Biol. Evol.* 38, 5539–5554. doi:10.1093/molbev/msab273
- Yoshida, Y., Konno, S., Nishino, R., Murai, Y., Tomita, M., and Arakawa, K. (2018). Ultralow input genome sequencing library preparation from a single tardigrade specimen. *J. Vis. Exp.*, 57615–57618. doi:10.3791/57615
- Yuan, Z., Liu, S., Zhou, T., Tian, C., Bao, L., Dunham, R., et al. (2018). Comparative genome analysis of 52 fish species suggests differential associations of repetitive elements with their living aquatic environments. *BMC Genomics* 19, 141–210. doi:10.1186/s12864-018-4516-1
- Zhang, C., Sayyari, E., and Mirarab, S. (2017). “ASTRAL-III: increased scalability and impacts of contracting low support branches,” in *Comparative genomics. RECOMB-CG 2017. Lecture notes in computer science* (Cham: Springer).
- Zhang, G., Li, C., Li, Q., Li, B., Larkin, D. M., Lee, C., et al. (2014). Comparative genomics reveals insights into avian genome evolution and adaptation. *Sci.* (80- 346, 1311–1320. doi:10.1126/science.1251385
- Zhu, B. H., Xiao, J., Xue, W., Xu, G. C., Sun, M. Y., and Li, J. T. (2018). P_RNA_scaffolder: A fast and accurate genome scaffolder using paired-end RNA-sequencing reads. *BMC Genomics* 19, 175–213. doi:10.1186/s12864-018-4567-3