

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25

Near chromosome-level and highly repetitive genome assembly of the snake pipefish *Entelurus aequoreus* (Syngnathiformes: Syngnathidae)

Magnus Wolf^{1,2,3}, Bruno Lopes da Silva Ferrette¹, Raphael T. F. Coimbra^{1,2}, Menno de Jong¹, Marcel Nebenfuehr^{1,2}, David Prochotta^{1,2}, Yannis Schöneberg^{1,2}, Konstantin Zapf^{1,2}, Jessica Rosenbaum², Hannah A. McIntyre², Julia Maier², Clara C.S. de Souza², Lucas M. Gehlhaar², Melina J. Werner², Henrik Oechler², Marie Wittekind², Moritz Sonnewald⁴, Maria A. Nilsson^{1,5}, Axel Janke^{1,2,5}, Sven Winter^{1,2,6}

¹ Senckenberg Biodiversity and Climate Research Centre (BiK-F), Frankfurt am Main, Germany

² Institute for Ecology, Evolution, and Diversity, Goethe University, Frankfurt am Main, Germany

³ Institute for Evolution and Biodiversity, University of Münster, Münster, Germany

⁴ Senckenberg Research Institute, Department of Marine Zoology, Section Ichthyology, Frankfurt am Main, Germany

⁵ LOEWE-Centre for Translational Biodiversity Genomics (TBG), Frankfurt am Main, Germany

⁶ Research Institute of Wildlife Ecology, University of Veterinary Medicine, Vienna, Austria

Authors for correspondence:

Sven Winter | Email: Sven.Winter@senckenberg.de

Magnus Wolf | Email: Magnus.Wolf@senckenberg.de

Name	e-mail	ORCID
Magnus Wolf (MW)	magnus.wolf@senckenberg.de	0000-0001-9212-9861
Bruno L. S. Ferrette (BF)	bruno.ferrette@senckenberg.de	0000-0002-3108-9867
Raphael T. F. Coimbra (RC)	raphael.t.f.coimbra@gmail.com	0000-0002-6075-7203
Menno de Jong (MDJ)	menno.de-jong@senckenberg.de	0000-0003-2131-9048
Marcel Nebenfuehr (MN)	marcel.nebenfuehr@senckenberg.de	0000-0001-8802-2105
David Prochotta (DP)	david.prochotta@senckenberg.de	0009-0000-6275-7752
Yannis Schöneberg (YS)	yannis.schoeneberg@gmx.de	0000-0003-1113-973X
Konstantin Zapf (KZ)	konstantin.zapf@senckenberg.de	
Jessica Rosenbaum (JR)	s5932510@stud.uni-frankfurt.de	0009-0008-2306-9015
Hannah A. Mc Intyre (HMI)	s3630184@stud.uni-frankfurt.de	0009-0002-3275-5048
Julia Maier (JM)	juliaomaier@googlemail.com	
Clara C.S. de Souza (CDS)	s6795932@stud.uni-frankfurt.de	0009-0000-9560-8905
Lucas M. Gehlhaar (LG)	s6244190@stud.uni-frankfurt.de	
Melina J. Werner (MJW)	s5033619@stud.uni-frankfurt.de	0009-0007-1081-009X

Henrik Oechler (HO)	henrik.oechler@uni-bayreuth.de	0009-0001-0413-8731
Marie Wittekind (MWI)	s0097351@stud.uni-frankfurt.de	0009-0001-2443-7552
Moritz Sonnewald (MS)	moritz.sonnewald@senckenberg.de	0000-0003-3042-8107
Maria A. Nilsson (MAN)	maria.nilsson-janke@senckenberg.de	0000-0002-8136-7263
Axel Janke (AJ)	axel.janke@senckenberg.de	0000-0002-9394-1904
Sven Winter (SW)	sven.winter@vetmeduni.ac.at	0000-0002-1890-0977

26

27

28

Abstract

29

30 The snake pipefish, *Entelurus aequoreus* (Linnaeus, 1758), is a slender, up to 60 cm long,
31 northern Atlantic fish that dwells in open seagrass habitats and has recently expanded its
32 distribution range. The snake pipefish is part of the family Syngnathidae (seahorses and
33 pipefish) that has undergone several characteristic morphological changes, such as loss of
34 pelvic fins and elongated snout. Here, we present a highly contiguous, near chromosome-scale
35 genome of the snake pipefish assembled as part of a university master's course. The final
36 assembly has a length of 1.6 Gbp in 7,391 scaffolds, a scaffold and contig N50 of 62.3 Mbp
37 and 45.0 Mbp and L50 of 12 and 14, respectively. The largest 28 scaffolds (>21 Mbp) span
38 89.7% of the assembly length. A BUSCO completeness score of 94.1% and a mapping rate
39 above 98% suggest a high assembly completeness. Repetitive elements cover 74.93% of the
40 genome, one of the highest proportions so far identified in vertebrate genomes. Demographic
41 modeling using the PSMC framework indicates a peak in effective population size (50 – 100
42 kya) during the last interglacial period and suggests that the species might largely benefit
43 from warmer water conditions, as seen today. Our updated snake pipefish assembly forms an
44 important foundation for further analysis of the morphological and molecular changes unique
45 to the family Syngnathidae.
46

47

48

49

Keywords

50 long reads, proximity-ligation scaffolding, genome annotation, demographic history,

51 repetitive elements

52

53

Introduction

54 The snake pipefish *Entelurus aequoreus* (Linnaeus 1758) is a member of the family
55 Syngnathidae, which currently includes over 300 species of seahorses and pipefishes [1]. The

56 species shares typical features with other pipefishes such as a unique, elongated body plan and
57 fused jaws [2]. However, unlike most pipefish, which are found in benthic habitats, the snake
58 pipefish inhabits more open and deeper seagrass environments and occurs even in pelagic
59 waters [2]. They are ambush predators on small crustaceans and other invertebrates, thereby
60 indirectly contributing to the overall biodiversity and stability of these fragile habitats [3].
61 Adult snake pipefish are poor swimmers with small fins and rely on their elongated, thin
62 bodies for crypsis in eelgrass habitats [4–6].

63 The snake pipefish historically ranged from the waters of Azores northwards to the
64 waters of Norway and Iceland, and eastward to the Baltic [7]. However, since 2003, the
65 species has expanded its distribution [8] into the arctic waters of Spitsbergen [9], the Barents
66 Sea, and the Greenland Sea [10]. Simultaneously, population sizes seem to increase within its
67 former range, as indicated by substantially increased catch rates [11, 12]. Several factors have
68 been proposed to cause this expansion and population growth, including rising sea
69 temperatures, an increased potential for long-distance dispersal of juveniles via ocean currents
70 [4, 7] and an increased reproductive success facilitated by the dispersal of invasive seaweeds
71 [6, 8–10, 13]. The latter explanation has been confirmed in local field experiments in the
72 northern Wadden Sea, suggesting a mutual co-occurrence of the invasive Japanese seaweed
73 (*Sargassum muticum*) and the snake pipefish [5]. Studies based on mtDNA marker regions
74 did not discern any population structure thus far and suggest a previous population expansion
75 in the Pleistocene ca. 50–100 kya [6]. Yet, a comprehensive analysis of demographic events is
76 better studied from genomic data, requiring a high-quality reference genome of ideally the
77 same species or at least a closely related one.

78 Previously, genomes of Syngnathidae have been used to study the evolution of highly
79 specialized morphologies and life-history traits unique to pipefishes and seahorses [14–16].
80 The transition to male pregnancy was associated with major genomic restructuring events and
81 parallel modifications of the adaptive immune system. There is a remarkable variability in

82 genome sizes within the family, with estimates ranging from 350 Mbp to 1.8 Gbp [14]. The
83 major shifts in body shape are assumed to be related to gene-family loss and expansion events
84 and higher rates of protein and nucleotide evolution [16]. Genomic data using a direct
85 sequencing approach of ultra-conserved elements (UCEs) improved the understanding of the
86 phylogeny of pipefishes [15] and identified a likely radiation of the group in the waters of the
87 modern Indo-Pacific Ocean. Nevertheless, high-quality genomes of Syngnathidae are only
88 available for a few species, and according to the NCBI genome database, only 7% of the
89 known species diversity has genome sequences available.

90 A draft genome of the snake pipefish was previously assembled using a combination
91 of paired-end and mate-pair sequencing techniques, yielding an assembly with low continuity
92 (N50 3.5 kbp, BUSCO C: 21%) and a large difference between the estimated and assembled
93 genome sizes (1.8 Gbp vs. 557 Mbp) [14]. To obtain a higher quality, near chromosome-scale
94 genome assembly for the snake pipefish for future population genomic, conservation, and
95 evolutionary studies of fish, we used long-read sequencing technologies. This allows us to
96 gain insight into the genetic properties of the species and to perform demographic analysis
97 based on the PSMC framework [17]. The data generation and analyses presented here were
98 conducted during a six-week master course in 2021 at the Goethe University, Frankfurt am
99 Main, Germany. The concept of high-quality genome sequencing in a course setting has so far
100 yielded three reference-quality genomes of fish and has proven to be a successful approach to
101 introduce the technology to a new generation of scientists [18–21].

102

103

104

Results and Discussion

105

Genome sequencing and assembly

106 PacBio's continuous long reads (CLR) technology generated 401 Gbp of long-read
107 data in ~60 million reads with an N50 of 7.9 kb (Table 1). Illumina sequencing yielded 38
108 Gbp of standard short-read data in ~257 million reads with a mean length of 148 bp after
109 filtering. Sequencing of the Omni-C library generated 54.7 Gbp of raw short-read data.

110 The snake pipefish's genome was assembled *de novo* to a total size of 1.7 Gbp. It
111 consisted of 2,204 scaffolds, with a scaffold N50 of 62 Mbp and an L50 of 11 (Table 1, Fig.
112 1A). The finalized assembly has 1.0 N's per 100 kbp and a GC content of 38.84%. A BUSCO
113 completeness assessment resulted in 94.1% complete core genes, given the
114 *actinopterygii_obd10* set, indicating high completeness of the assembly. Both long- and short-
115 read data mapped onto the assembly with high mapping rates of 98.6% and 99.5%,
116 respectively. HI-C mapping resulted in 28 larger scaffolds (Fig. 1B), indicating the near-
117 chromosome level of the *de novo* assembly as past karyotype estimations of other pipefish
118 and seahorses predicted 22 and 22-24 chromosomes, respectively [22–24]. The rest of the
119 genome comprises only smaller scaffolds and contigs, which may result from the high
120 amounts of repetitive regions described in the following section. Our Blobtools analysis of
121 both long- and short-read data (Fig. 1C+D) found no apparent signs of contamination,
122 although background noise of unknown origin was detected and removed in both datasets.

123 Variant calling resulted in ~301 million sites (including monomorphic sites), of which
124 ~1.3 million were found to be biallelic. Genome-wide heterozygosity was determined to be
125 0.387%, which is in line with other fish species [25, 26]. The GenomeScope results based on
126 short reads suggested a haploid genome size of 1.15 Gbp and an expected genome-wide
127 heterozygosity of 1%, around 362 Mbp shorter and 0.57% more heterozygous when compared
128 to the final assembly. This, again, might be explained by the high repeat content in the
129 genome.

130

Annotation

131 In total, 0.9 Gbp, or 74.93%, of the entire assembly, were identified as repetitive
132 during our *de novo* repeat-modeling and repeat-masking (Fig. 2). This high repeat content
133 contrasts that of other fish genomes [27], but is similar, although at a smaller scale, to the
134 closest relative *Nerophis ophidion* (65.7%) [14] and other genomes of syngnathid fish like
135 e.g., seadragons [28]. The first draft assembly of the snake pipefish had a repeat content of
136 57.2% [14] and our improved long-read assembly identified 17.7% additional repeats that
137 were missing from the previous assembly [14]. So far, among vertebrates, only the lungfish
138 *Neoceratodus forsteri* [29] has more transposable elements (TEs) than the snake pipefish.

139 The annotation of the genome featuring *de novo* and homology-based identification
140 approaches resulted in 33,202 genes with an average length of 13,828 bp. Each gene had on
141 average 7.32 exons and 6.25 introns with average lengths of 188 bp and 2,240 bp,
142 respectively. In total, we identified 243,038 exons and 207,467 introns within our annotation.
143 The total number of genes is ~30% larger compared to other annotated genomes in the order
144 of Syngnathiformes like, e.g., 23,458 for the tiger tail seahorse (*Hippocampus comes*) [16] or
145 24,927 for the greater pipefish (*Syngnathus acus*) [30] made by the NCBI Eukaryotic Genome
146 Annotation pipeline. Given that these two genomes are also considerably smaller, 492 Mbp
147 and 324 Mbp, respectively, it can be assumed that the large-scale genome increase in this
148 species also included many coding sequences. A high content of repetitive regions as well as a
149 lack of transcriptomic data might also have increased the number of false positive gene-calls;
150 however, a BUSCO completeness analysis of the predicted proteins resulted in 82.6%
151 complete sequences, of which only 6.8% were duplicated. 5.3% of the coding sequences
152 appeared fragmented, and 12.1% were missing from the *actinopterygii_obd10* OrthoDB set.
153 A functional annotation resulted in hits for 89% of the predicted proteins.

154

Demographic inference

178 Fair and Equitable Sharing of Benefits Arising from Their Utilization’. The sample was
179 initially frozen at -20 °C and later stored at -80 °C.

180 High molecular weight genomic DNA was extracted from muscle tissue, following the
181 protocol by Mayjonade et al. [33] with the addition of Proteinase K. We evaluated the
182 quantity and quality of the DNA with the Genomic DNA ScreenTape on the Agilent 2200
183 TapeStation system (Agilent Technologies), as well as with the Qubit® dsDNA BR Assay
184 Kit.

185 For long-read sequencing, a PacBio SMRT Bell continuous long read (CLR) library was
186 prepared using the SMRTbell Express Prep kit v3.0 kit (Pacific Biosciences – PacBio, Menlo
187 Park, CA, USA) and sequenced on the PacBio Sequel IIe platform. A proximity-ligation
188 library was compiled with muscle tissue following the Dovetail™ Omni-C protocol (Dovetail
189 Genomics, Santa Cruz, California, USA). In addition, a standard whole-genome 150 base pair
190 (bp) paired-end Illumina library was prepared using the NEBNext Ultra II library preparation
191 kit (New England Biolabs Inc., Ipswich, USA). Finally, the proximity ligation and the paired-
192 end library were shipped to Novogene (UK) for sequencing on the Illumina NovaSeq 6000
193 platform.

194 **Pre-processing & Genome size estimation**

195 The PacBio subreads were converted from BAM into FASTQ format using the PacBio
196 Secondary Analysis Tool BAM2fastx v.1.3.0 ([https://github.com/PacificBiosciences/
197 pbbioconda](https://github.com/PacificBiosciences/pbbioconda)). Quality control, trimming, and filtering of the Illumina reads were performed
198 using fastp v0.23.1 [34] with the settings “-g -3 -l 40 -y -Y 30 -q 15 -u 40 -c -p -j -h -R -w N”.
199 To estimate the genome size of the snake pipefish, we performed *k*-mer profiling using the
200 standard short-read Illumina data. We first ran Jellyfish v2.3.0 [35] to generate a histogram of
201 *k*-mers with a length of 21 bp. Subsequently, we used this data to obtain a genome profile

202 using GenomeScope v2.0 [36]. We further tested alternative k-mer lengths between 17- and
203 25-mers which resulted in no meaningful differences of the estimated genome size except for
204 the 17-mer, which resulted in a smaller genome size estimation of ~500 Mbp.

205

206 **Genome Assembly & polishing**

207 We assembled the genome from the PacBio long-read data using WTDBG v.2.5 [37]. The
208 resulting assembly was first polished using the PacBio data with Flye v.2.9 [38], using
209 Minimap v.2.17 [39] for mapping, followed by two rounds of short-read polishing by
210 mapping reads onto the assembly with BWA-MEM v.0.7.17 [40] and error correction with
211 Pilon v1.23 [41].

212 **Assembly QC & Scaffolding**

213 The polished assembly contigs were anchored into chromosome-scale scaffolds utilizing the
214 generated proximity-ligation Omni-C data. First, the data were mapped and filtered to the
215 assembly following the Arima Hi-C mapping pipeline used by the Vertebrate Genome Project
216 (https://github.com/VGP/vgpassembly/blob/master/pipeline/salsa/arima_mapping_pipeline.sh
217). In brief, reads were mapped using BWA-MEM v.0.7.17 [40], mapped reads were filtered
218 with samtools v.1.14 [42], and duplicated reads were removed with “MarkDuplicates” in
219 Picard v.2.26.10 (Broad Institute, 2019). The filtered mapped reads were then used for
220 proximity-ligation scaffolding in YaHs v.1.1 [43]. Gaps in the scaffolded assembly were
221 closed with TGS-GapCloser v.1.1.1 [44] using a subset (25%) of the PacBio subreads due to
222 computational constraints. To further improve the assembly's contiguity, scaffolding and gap-
223 closing were performed a second time using a different subset of PacBio reads for gap-
224 closing. The PacBio read subsets were generated with seqtk v.1.3
225 (<https://github.com/lh3/seqtk>) using the random number generator seeds 11 and 18. Gene set

226 completeness was analyzed with BUSCO v.5.4.7 [45] using the Actinopterygii set of core
227 genes (*actinopterygii_odb10*). Assembly continuity was evaluated using QUAST v5.0.2 [46],
228 and mapping rates were assessed by Qualimap v2.2.1 [47]. BlobToolsKit v.4.0.6 [48] was
229 used to perform contamination screening.

230 **Repeat landscape analysis & genome annotation**

231 The TE annotation was done in three steps. First, we used RepeatMasker v4.1.5 [49] to
232 annotate, and hard-mask known Actinopterygii repeats from RepBase, which comprises a
233 database of eukaryotic repetitive DNA element sequences [50]. Secondly, a *de novo* library of
234 transposable elements was created from the hard-masked genome assembly using
235 RepeatModeler v2.0.4 [51] which includes RECON v1.08 [52], RepeatScout v1.0.6 [53], and
236 LTRharvest/LTR_retriever [54, 55]. Finally, predicted repeats were annotated with a second
237 run of RepeatMasker on the hard-masked assembly obtained in the first run. The results from
238 both RepeatMasker runs were then combined. A summary of transposable elements and the
239 relative abundance of repeat classes in the genome are shown in Table 2 and Fig. 2.

240 The genome was annotated using the BRAKER3 pipeline [56–61], combining a *de novo* gene
241 calling and a homology-based gene annotation. For protein references, we combined the
242 vertebrate-specific protein collection from OrthoDB and the protein collection of the greater
243 pipefish (*Syngnathus acus*) genome [30] made by the NCBI (see: GCF_901709675.1, last
244 accessed 12th Oct. 2023). To further filter genes based on the support of introns by extrinsic
245 homology evidence, we used TSEBRA [62] with an “intron_support=0.1”. The resulting set
246 of proteins was tested for completeness using BUSCO v.5.4.7 [45] in “protein mode” and run
247 against the Actinopterygii-specific set of core genes. Functional annotation was done using
248 InterProScan v5 [63].

249 **Variant calling & demographic inference**

250 The preprocessed short reads were mapped to the final assembly using BWA-MEM v. 0.7.17
251 [40] followed by removal of duplicate reads with "MarkDuplicates" in Picard v.2.26.10
252 (Broad Institute, 2019) and evaluation of mapping quality using Qualimap v2.2.1 [47]. Indels
253 in the BAM files were first identified and then realigned with "RealignerTargetCreator" and
254 "IndelRealigner" as part of the Genome Analysis Toolkit (GATK) v3.8-1
255 (<https://gatk.broadinstitute.org/>). Subsequently, samtools v.1.14 [42] was used to check and
256 remove unmapped, secondary, QC failed, duplicated, and supplementary reads keeping only
257 reads mapped in proper pairs in non-repetitive regions of the 28 chromosome-scale scaffolds.

258 Sambamba v 1.0.0 [64] was used to estimate site depth statistics. Minimum and maximum
259 thresholds for the global site depth were set to $d \pm (5 \times \text{MAD})$, where d is the global site depth
260 distribution median and MAD is the median absolute deviation. Variant calling was
261 performed using the bcftools v1.17 [65] commands "mpileup" and "call" [-m]. Variants were
262 then filtered with bcftools "filter" [-e "DP < d - (5 × MAD) || DP > d + (5 × MAD) ||
263 QUAL < 30"] removing sites with low quality and out of range depth. Finally, bcftools was
264 used to estimate the genome-wide heterozygosity as the proportion of heterozygous sites
265 using the "stats" command.

266 Long-term changes in effective population size (N_e) over time were estimated with the
267 Pairwise Sequentially Markovian Coalescent (PSMC) model [17] based on the diploid
268 consensus genome sequences generated by bcftools v1.17 [65] with the script "*vcfutils.pl*"
269 from the processed BAM files, as described above. Sites with read-depth up to a third of the
270 average depth or above twice each sample's median depth and with a consensus base quality
271 < 30 were removed. PSMC was executed using 25 iterations with a maximum $2N_0$ -scaled
272 coalescent time of 15, an initial θ/ρ ratio of 5, and 64 atomic time intervals (4 + 25 × 2 + 4 +
273 6) to infer the scaled mutation rate, the recombination rate, and the free population size
274 parameters, respectively. We performed 100 bootstrap replicates by randomly sampling with

275 replacement of 1 Mb blocks from the consensus sequence for all individuals. A mutation rate
276 (μ) of 1.7×10^{-9} per site per generation [66] and a generation length of 2.5 years [67] were
277 employed for plotting.

278 **Availability of Supporting Data**

279 The *de novo* genome and all underlying raw data were uploaded to NCBI under the
280 BioProject PRJNA1005573, BioSample SAMN36988691, genome assembly
281 JAVRRV000000000. All other data, including the repeat and gene annotation, was uploaded
282 to the GigaDB repository: DOI:XXXXX. [Rawdata available for review at
283 <https://dataview.ncbi.nlm.nih.gov/object/PRJNA1005573?reviewer=2i5vm98fdb4r9j0asoff8m>
284 [sn3](#)]

285 **Author Contributions**

286 MW, BF, MS, AJ, and SW designed the study. SW, JR, HMI, JM, CDS, LG, MJW, HO, and
287 MWI performed laboratory procedures and sequencing. MW, BF, RC, MDJ, MN, DP, YS,
288 KZ, JR, HMI, JM, CDS, LG, MJW, HO, MWI, MAN, and SW conducted bioinformatic
289 processing and analyses. All authors contributed to writing this manuscript.

290 **List of Abbreviations**

291 bp: base pairs; BUSCO: Benchmarking Universal Single-Copy Orthologs; CLR: continuous
292 long reads; Gbp: Gigabase pairs; kbp: kilobase pairs; kya: thousand years ago; Mbp:
293 megabase pairs; Mya: million years ago; N_e : effective population size; PacBio: Pacific
294 Biosciences; PSMC: Pairwise Sequentially Markovian Coalescent; TEs: transposable
295 elements; UCEs: ultra-conserved elements.

296 **Conflict of Interest**

297 The authors declare that they have no competing interests.

298

Ethics Approval and Consent for Publication

299 Not Applicable.

300

Acknowledgments

301 The present study is a result of the Centre for Translational Biodiversity Genomics (LOEWE-
302 TBG) and was supported through the program ‘LOEWE-Landes-Offensive zur Entwicklung
303 Wissenschaftlich-ökonomischer Exzellenz’ of Hesse’s Ministry of Higher Education,
304 Research, and the Arts.

305

References

- 306 1. Froese R, Pauly D. FishBase. 2023. www.fishbase.org. Accessed 9 Aug 2023.
- 307 2. Dawson C. Syngnathidae. In: Smith M, Heemstra P, editors. *Smiths' sea fishes*. Berlin: Springer-Verlag;
308 1986. p. 445–458.
- 309 3. O’Gorman EJ. Multitrophic diversity sustains ecological complexity by dampening top-down control of
310 a shallow marine benthic food web. *Ecology*. 2021;102:e03274. doi:10.1002/ecy.3274.
- 311 4. Vincent ACJ, Berglund A, Ahnesj I. Reproductive ecology of five pipefish species in one eelgrass
312 meadow. *Environ Biol Fish*. 1995;44:347–61. doi:10.1007/BF00008250.
- 313 5. Polte P, Buschbaum C. Native pipefish *Entelurus aequoreus* are promoted by the introduced seaweed
314 *Sargassum muticum* in the northern Wadden Sea, North Sea. *Aquat. Biol.* 2008;3:11–8.
315 doi:10.3354/ab00071.
- 316 6. Braga Goncalves I, Cornetti L, Couperus AS, van Damme CJG, Mobley KB. Phylogeography of the
317 snake pipefish, *Entelurus aequoreus* (Family: Syngnathidae) in the northeastern Atlantic Ocean.
318 *Biological Journal of the Linnean Society*. 2017;122:787–800. doi:10.1093/biolinnean/blx112.
- 319 7. Wheeler A. *Key to the Fishes of Northern Europe: A guide to the identification of more than 350*
320 *species*. London: Frederick Warne & Co. Ltd; 1978.
- 321 8. Harris MP, Beare D, Toresen R, Nøttestad L, Kloppmann M, Dörner H, et al. A major increase in snake
322 pipefish (*Entelurus aequoreus*) in northern European seas since 2003: potential implications for seabird
323 breeding success. *Mar Biol*. 2007;151:973–83. doi:10.1007/s00227-006-0534-7.
- 324 9. Fleischer D, Schaber M, Piepenburg D. Atlantic snake pipefish (*Entelurus aequoreus*) extends its
325 northward distribution range to Svalbard (Arctic Ocean). *Polar Biol.* 2007;30:1359–62.
326 doi:10.1007/s00300-007-0322-y.
- 327 10. Rusyaev SM, Dolgov AV, Karamushko OV. Captures of snake pipefish *Entelurus aequoreus* in the
328 Barents and Greenland Seas. *J. Ichthyol.* 2007;47:544–6. doi:10.1134/S0032945207070090.
- 329 11. Kloppmann MHF, Ulleweit J. Off-shelf distribution of pelagic snake pipefish, *Entelurus aequoreus*
330 (*Linnaeus, 1758*), west of the British Isles. *Mar Biol*. 2007;151:271–5. doi:10.1007/s00227-006-0480-4.
- 331 12. van Damme CJ, Couperus AS. Mass occurrence of snake pipefish in the Northeast Atlantic: Result of a
332 change in climate? *Journal of Sea Research*. 2008;60:117–25. doi:10.1016/j.seares.2008.02.009.
- 333 13. Lindley J, Kirby R, Johns D, Reid C. Exceptional abundance of the snake pipefish (*Entelurus*
334 *aequoreus*) in the north-eastern North Atlantic Ocean. ICES Document. 2006.
- 335 14. Roth O, Solbakken MH, Tørresen OK, Bayer T, Matschiner M, Baalsrud HT, et al. Evolution of male
336 pregnancy associated with remodeling of canonical vertebrate immunity in seahorses and pipefishes.
337 *Proc Natl Acad Sci U S A*. 2020;117:9431–9. doi:10.1073/pnas.1916251117.
- 338 15. Stiller J, Short G, Hamilton H, Saarman N, Longo S, Wainwright P, et al. Phylogenomic analysis of
339 Syngnathidae reveals novel relationships, origins of endemic diversity and variable diversification rates.
340 *BMC Biol*. 2022;20:75. doi:10.1186/s12915-022-01271-w.
- 341 16. Lin Q, Fan S, Zhang Y, Xu M, Zhang H, Yang Y, et al. The seahorse genome and the evolution of its
342 specialized morphology. *Nature*. 2016;540:395–9. doi:10.1038/nature20595.

- 343 17. Li H, Durbin R. Inference of human population history from individual whole-genome sequences.
344 Nature. 2011;475:493–6. doi:10.1038/nature10231.
- 345 18. Prost S, Winter S, Raad J de, Coimbra RTF, Wolf M, Nilsson MA, et al. Education in the genomics era:
346 Generating high-quality genome assemblies in university courses. *Gigascience* 2020.
347 doi:10.1093/gigascience/giaa058.
- 348 19. Prost S, Petersen M, Grethlein M, Hahn SJ, Kuschik-Maczollek N, Olesiuk ME, et al. Improving the
349 Chromosome-Level Genome Assembly of the Siamese Fighting Fish (*Betta splendens*) in a University
350 Master's Course. *G3 (Bethesda)*. 2020;10:2179–83. doi:10.1534/g3.120.401205.
- 351 20. Winter S, Prost S, Raad J de, Coimbra RTF, Wolf M, Nebenführ M, et al. Chromosome-level genome
352 assembly of a benthic associated Syngnathiformes species: the common dragonet, *Callionymus lyra*.
353 *GigaByte*. 2020;2020:gigabyte6. doi:10.46471/gigabyte.6.
- 354 21. Winter S, Raad J de, Wolf M, Coimbra RTF, Jong MJ de, Schöneberg Y, et al. A chromosome-scale
355 reference genome assembly of the great sand eel, *Hyperoplus lanceolatus*. *J Hered*. 2023;114:189–94.
356 doi:10.1093/jhered/esad003.
- 357 22. Vitturi R, Catalano E. Karyotypes in two species of the genus *Hippocampus* (Pisces: Syngnathiformes).
358 *Mar Biol*. 1988;99:119–21. doi:10.1007/BF00644985.
- 359 23. Vitturi R, Libertini A, Campolmi M, Calderazzo F, Mazzola A. Conventional karyotype, nucleolar
360 organizer regions and genome size in five Mediterranean species of Syngnathidae (Pisces,
361 Syngnathiformes). *Journal of Fish Biology*. 1998;52:677–87. doi:10.1111/j.1095-8649.1998.tb00812.x.
- 362 24. Small CM, Bassham S, Catchen J, Amores A, Fuiten AM, Brown RS, et al. The genome of the Gulf
363 pipefish enables understanding of evolutionary innovations. *Genome Biol*. 2016;17:258.
364 doi:10.1186/s13059-016-1126-6.
- 365 25. Tigano A, Jacobs A, Wilder AP, Nand A, Zhan Y, Dekker J, Therikildsen NO. Chromosome-Level
366 Assembly of the Atlantic Silverside Genome Reveals Extreme Levels of Sequence Diversity and
367 Structural Genetic Variation. *Genome Biol Evol* 2021. doi:10.1093/gbe/evab098.
- 368 26. Barry P, Broquet T, Gagnaire P-A. Age-specific survivorship and fecundity shape genetic diversity in
369 marine fishes. *Evol Lett*. 2022;6:46–62. doi:10.1002/evl3.265.
- 370 27. Shao F, Han M, Peng Z. Evolution and diversity of transposable elements in fish genomes. *Sci Rep*.
371 2019;9:15399. doi:10.1038/s41598-019-51888-1.
- 372 28. Small CM, Healey HM, Currey MC, Beck EA, Catchen J, Lin ASP, et al. Leafy and weedy seadragon
373 genomes connect genic and repetitive DNA features to the extravagant biology of syngnathid fishes.
374 *Proc Natl Acad Sci U S A*. 2022;119:e2119602119. doi:10.1073/pnas.2119602119.
- 375 29. Meyer A, Schloissnig S, Franchini P, Du K, Woltering JM, Irisarri I, et al. Giant lungfish genome
376 elucidates the conquest of land by vertebrates. *Nature*. 2021;590:284–9. doi:10.1038/s41586-021-
377 03198-8.
- 378 30. Scott-Somme K, McTierney S, Brittain R, Perry F, Brenen M. The genome sequence of the greater
379 pipefish, *Syngnathus acus* (Linnaeus, 1758). *Wellcome Open Res*. 2023;8:274.
380 doi:10.12688/wellcomeopenres.19528.1.
- 381 31. Obrochta SP, Crowley TJ, Channell JE, Hodell DA, Baker PA, Seki A, Yokoyama Y. Climate
382 variability and ice-sheet dynamics during the last three glaciations. *Earth and Planetary Science Letters*.
383 2014;406:198–212. doi:10.1016/j.epsl.2014.09.004.
- 384 32. Armstrong E, Hopcroft PO, Valdes PJ. A simulated Northern Hemisphere terrestrial climate dataset for
385 the past 60,000 years. *Sci Data*. 2019;6:265. doi:10.1038/s41597-019-0277-1.
- 386 33. Mayjonade B, Gouzy J, Donnadieu C, Pouilly N, Marande W, Callot C, et al. Extraction of high-
387 molecular-weight genomic DNA for long-read sequencing of single molecules. *Biotechniques*.
388 2016;61:203–5. doi:10.2144/000114460.
- 389 34. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*.
390 2018;34:i884-i890. doi:10.1093/bioinformatics/bty560.
- 391 35. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-
392 mers. *Bioinformatics*. 2011;27:764–70. doi:10.1093/bioinformatics/btr011.
- 393 36. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC.
394 *GenomeScope*: fast reference-free genome profiling from short reads. *Bioinformatics*. 2017;33:2202–4.
395 doi:10.1093/bioinformatics/btx153.
- 396 37. Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *Nat Methods*. 2020;17:155–8.
397 doi:10.1038/s41592-019-0669-3.
- 398 38. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs.
399 *Nat Biotechnol*. 2019;37:540–6. doi:10.1038/s41587-019-0072-8.
- 400 39. Li H. *Minimap2*: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34:3094–100.
401 doi:10.1093/bioinformatics/bty191.
- 402 40. Li H. *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*. arXiv. 2013.

- 403 41. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated tool for
404 comprehensive microbial variant detection and genome assembly improvement. *PLoS One*.
405 2014;9:e112963. doi:10.1371/journal.pone.0112963.
- 406 42. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
407 format and SAMtools. *Bioinformatics*. 2009;25:2078–9. doi:10.1093/bioinformatics/btp352.
- 408 43. Zhou C, McCarthy SA, Durbin R. YaHS: yet another Hi-C scaffolding tool. *Bioinformatics* 2023.
409 doi:10.1093/bioinformatics/btac808.
- 410 44. Xu M, Guo L, Gu S, Wang O, Zhang R, Peters BA, et al. TGS-GapCloser: A fast and accurate gap
411 closer for large genomes with low coverage of error-prone long reads. *Gigascience* 2020.
412 doi:10.1093/gigascience/giaa094.
- 413 45. Manni M, Berkeley MR, Seppey M, Zdobnov EM. BUSCO: Assessing Genomic Data Quality and
414 Beyond. *Curr Protoc*. 2021;1:e323. doi:10.1002/cpz1.323.
- 415 46. Mikhchenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A. Versatile genome assembly evaluation
416 with QUAST-LG. *Bioinformatics*. 2018;34:i142-i150. doi:10.1093/bioinformatics/bty266.
- 417 47. Okonechnikov K, Conesa A, García-Alcalde F. Qualimap 2: advanced multi-sample quality control for
418 high-throughput sequencing data. *Bioinformatics*. 2016;32:292–4. doi:10.1093/bioinformatics/btv566.
- 419 48. Challis R, Richards E, Rajan J, Cochrane G, Blaxter M. BlobToolKit – Interactive quality assessment of
420 genome assemblies; 2019.
- 421 49. Smit A, Hubley R, Green P. RepeatMasker Open-4.0. 2013. <http://www.repeatmasker.org>.
- 422 50. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic
423 genomes. *Mob DNA*. 2015;6:11. doi:10.1186/s13100-015-0041-9.
- 424 51. Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. RepeatModeler2 for
425 automated genomic discovery of transposable element families. *Proc Natl Acad Sci U S A*.
426 2020;117:9451–7. doi:10.1073/pnas.1921046117.
- 427 52. Bao Z, Eddy SR. Automated de novo identification of repeat sequence families in sequenced genomes.
428 *Genome Res*. 2002;12:1269–76. doi:10.1101/gr.88502.
- 429 53. Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large genomes.
430 *Bioinformatics*. 2005;21 Suppl 1:i351-8. doi:10.1093/bioinformatics/bti1018.
- 431 54. Ou S, Jiang N. LTR_retriever: A Highly Accurate and Sensitive Program for Identification of Long
432 Terminal Repeat Retrotransposons. *Plant Physiol*. 2018;176:1410–22. doi:10.1104/pp.17.01310.
- 433 55. Ellinghaus D, Kurtz S, Willhoeft U. LTRharvest, an efficient and flexible software for de novo
434 detection of LTR retrotransposons. *BMC Bioinformatics*. 2008;9:18. doi:10.1186/1471-2105-9-18.
- 435 56. Bruna T, Lomsadze A, Borodovsky M. GeneMark-ETP: Automatic Gene Finding in Eukaryotic
436 Genomes in Consistency with Extrinsic Data; 2023.
- 437 57. Bruna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. BRAKER2: automatic eukaryotic genome
438 annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genom*
439 *Bioinform*. 2021;3:lqaa108. doi:10.1093/nargab/lqaa108.
- 440 58. Kovaka S, Zimin AV, Pertea GM, Razaghi R, Salzberg SL, Pertea M. Transcriptome assembly from
441 long-read RNA-seq alignments with StringTie2. *Genome Biol*. 2019;20:278. doi:10.1186/s13059-019-
442 1910-1.
- 443 59. Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: Unsupervised RNA-Seq-Based
444 Genome Annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics*. 2016;32:767–9.
445 doi:10.1093/bioinformatics/btv661.
- 446 60. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat Methods*.
447 2015;12:59–60. doi:10.1038/nmeth.3176.
- 448 61. Gabriel L, Bruna T, Hoff KJ, Ebel M, Lomsadze A, Borodovsky M, Stanke M. BRAKER3: Fully
449 Automated Genome Annotation Using RNA-Seq and Protein Evidence with GeneMark-ETP,
450 AUGUSTUS and TSEBRA. *bioRxiv* 2023. doi:10.1101/2023.06.10.544449.
- 451 62. Gabriel L, Hoff KJ, Bruna T, Borodovsky M, Stanke M. TSEBRA: transcript selector for BRAKER.
452 *BMC Bioinformatics*. 2021;22:566. doi:10.1186/s12859-021-04482-0.
- 453 63. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale
454 protein function classification. *Bioinformatics*. 2014;30:1236–40. doi:10.1093/bioinformatics/btu031.
- 455 64. Tarasov A, Vilella AJ, Cuppen E, Nijman IJ, Prins P. Sambamba: fast processing of NGS alignment
456 formats. *Bioinformatics*. 2015;31:2032–4. doi:10.1093/bioinformatics/btv098.
- 457 65. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, et al. Twelve years of SAMtools
458 and BCFtools. *Gigascience* 2021. doi:10.1093/gigascience/giab008.
- 459 66. He L, Long X, Qi J, Wang Z, Huang Z, Wu S, et al. Genome and gene evolution of seahorse species
460 revealed by the chromosome-level genome of *Hippocampus abdominalis*. *Mol Ecol Resour*.
461 2022;22:1465–77. doi:10.1111/1755-0998.13541.
- 462 67. Schultz J. *Entelurus aequoreus*: IUCN Red List of Threatened Species, e.T18258072A44775951; 2014.

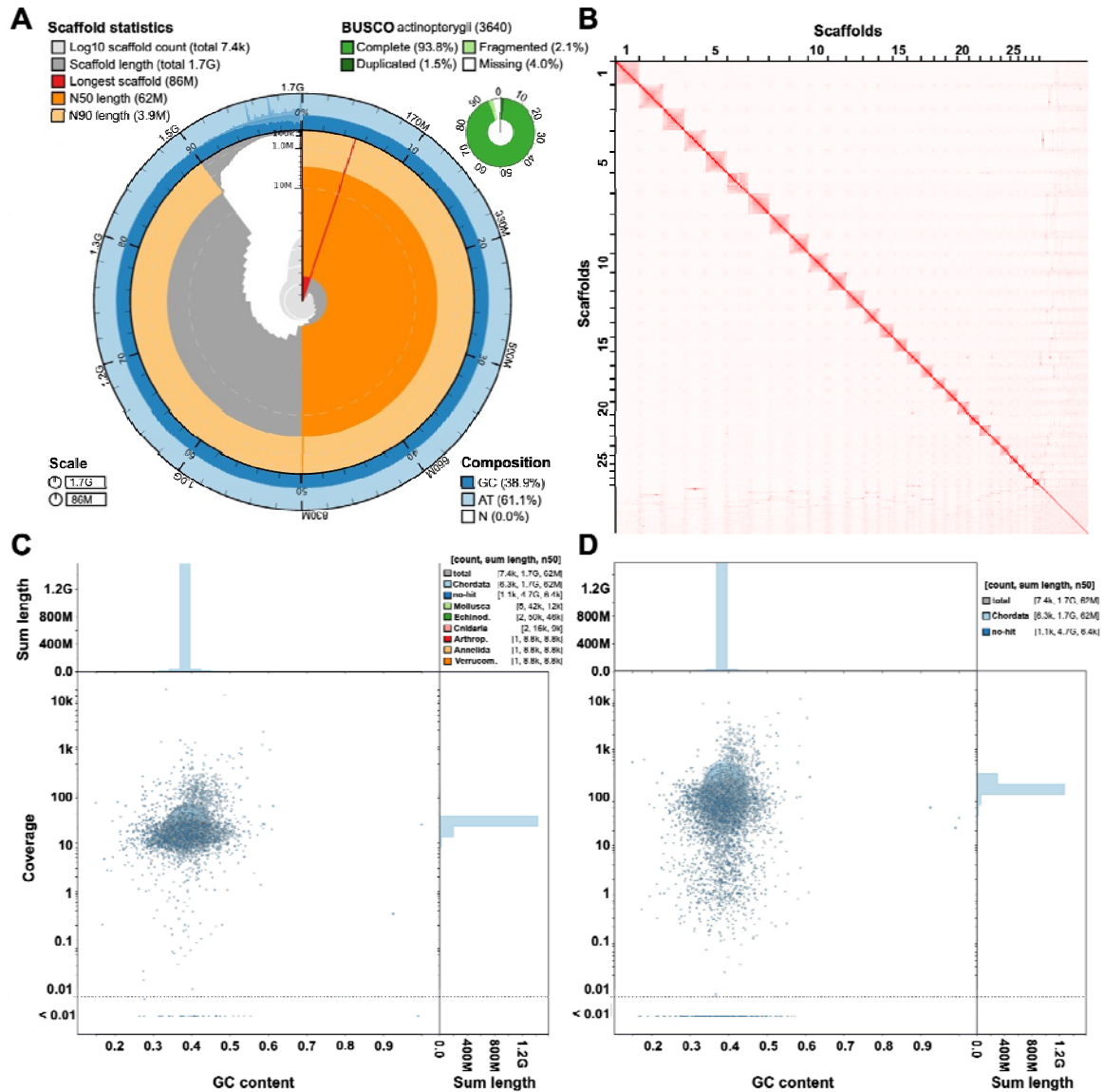
463

464

465

466

Figures

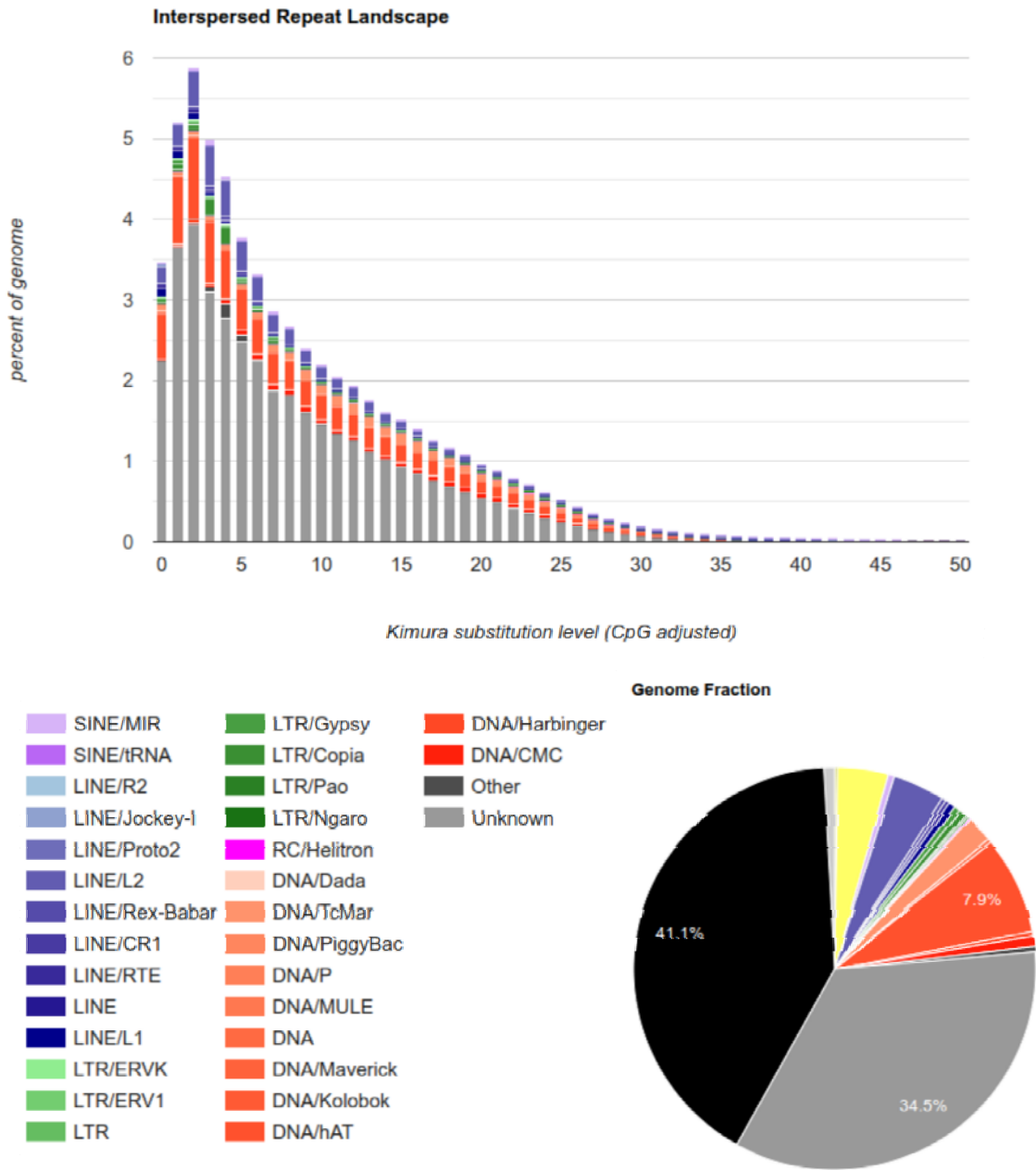


467

468 **Figure 1** Assembly characteristics and quality assessments of the *de novo* *Entelurus aequoreus* genome. **A** The
 469 snail plot summarizes different assembly properties. Scaffold statistics are depicted in the innermost circle and
 470 the colors red to orange represent the longest scaffold N50 and N90, respectively. GC composition is shown in
 471 the outer blue circle. BUSCO completeness statistics are depicted in the small green circle. **B** Omni-C contact
 472 density map indicating 28 larger scaffolds and the near-chromosome level of the assembly. **C-D** The BlobPlot
 473 analysis compares GC content (x-axis), assembly coverage (y-axis) and taxonomic BLAST assignments of
 474 contigs (color) for both the Omni-C short reads (C) and PacBio long reads (D).

475

476



477

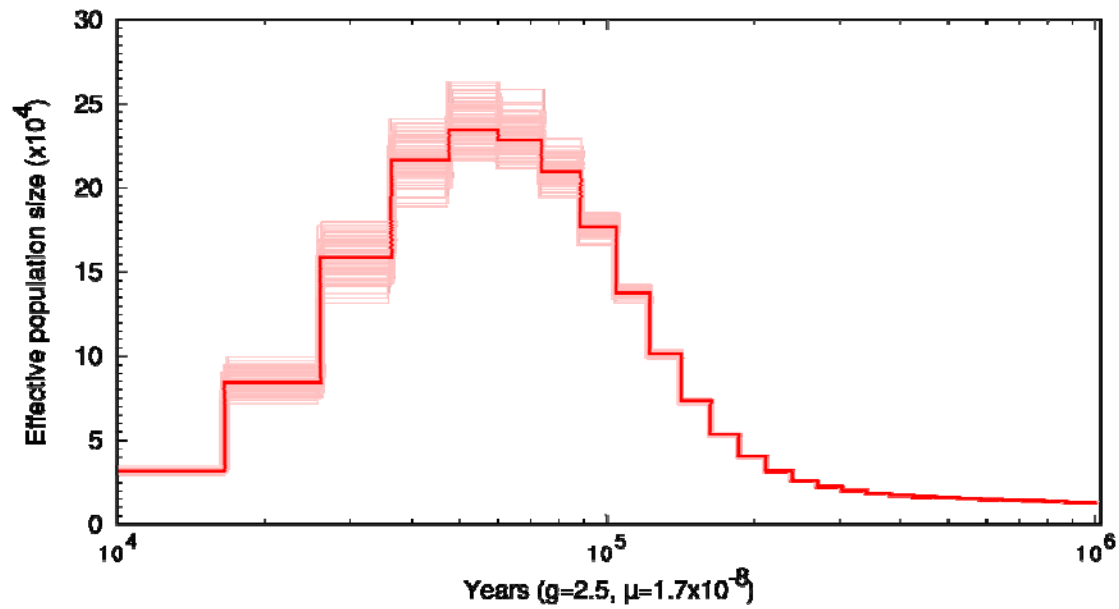
478 **Figure 2** Repeat landscape of the *de novo* *Entelurus aequoreus* genome. Colors represent repetitive element
 479 types, gray areas indicate unclassified types of repetitive regions.

480

481

482

483



484

485 **Figure 3** Demographic history of the snake pipefish estimated using the PSMC framework. Using a generation
486 time of 2.5 years [67] and a substitution rate of 1.7×10^{-8} per site per generation [66] a model was created
487 covering the last 10 kya to 1 Mya. The x-axis represents time in number of years ago and the y-axis shows the
488 effective population (N_e) size in tens of thousands of individuals. The model indicates a peak in N_e of 250
489 thousand individuals during the Pleistocene at around 100 thousand years ago.

490

491

492

493

494

495

496

497

498

499

500

501

502

503

Tables

504 **Table 1** Summary statistics of the snake pipefish reference genome. The table includes information for **A** the
 505 raw read sequencing, and **B** the scaffold- and contig-level *de novo* assembly and **C** the BUSCO completeness
 506 statistics.

(A)	Raw read statistics	
No. short reads	264,111,731	
Mapped short reads (%)	99.53	
Mean short read coverage (x)	23	
No. long reads	130,590,372	
Mapped long reads (%)	98.61	
Mean long read coverage (x)	205.2	
(B)	Assembly statistics (scaffold/contig)	
No. scaffolds/contigs	7,387	7,473
No. scaffolds/contigs (>50 kbp)	466	526
scaffold/contig L50	12	14
scaffold/contig N50 (bp)	62,341,166	45,010,074
Total length (bp)	1,662,053,046	1,662,035,846
GC (%)	38.87	38.87
No. of N's per 100 kb	1.03	0.0
heterozygosity (%)	0.387	
Total interspersed repeats (bp)	1,237,929,559 (74.93 %)	
(C)	BUSCO completeness	
Clade: <i>Actinopterygii</i>	C:94.1%[S:92.6%, D:1.5%] F:2.0%, M:3.9% n:3640	

507
 508
 509
 510
 511
 512
 513
 514
 515
 516
 517
 518
 519

BUSCO: Benchmarking Universal Single Copy Orthologs (65); C, complete;
 S, single copy; D, duplicated; F, fragmented; M, missing.

520

521 **Table 2** Repeat content of the genome assembly. Class, class of the repetitive regions. Count, number of
522 occurrences of the repetitive region. bpMasked, number of base pairs masked; %Masked, percentage of base pairs
523 masked. LINE, Long Interspersed Nuclear Elements (include retroposons); LTR, Long Terminal Repeat
524 elements (including retroposons); SINE, Short Interspersed Nuclear Elements; RC, Rolling Circle.
525

Class	Count	bpMasked	%masked
ARTEFACT	4	84	0.00%
DNA	2765297	372407739	22.40%
LINE	850222	167337419	10.06%
LTR	177214	55439687	3.33%
PLE	1	0	0.00%
RC	32348	3385084	0.20%
SINE	435464	32709572	1.95%
Unknown	3628328	534216084	32.14%
Low complexity	127733	3095322	0.19%
Satellite	21221	7145469	0.43%
Simple repeat	1437090	61077339	3.67%
rRNA	4394	534599	0.03%
scRNA	5	504	0.00%
snRNA	695	46845	0.00%
tRNA	6029	533812	0.03%
Total	9486045	1237929559	74.93%

526