

The complete chromosome-level genome and mitogenome assembly of *Seriola lalandi lalandi* (Yellowtail Kingfish) for aquaculture and fisheries management.

Carla Finn¹, Maren Wellenreuther², David Chagne², Tom Oosting¹, Yvan Papa³, Alvin Setiawan⁴, Vinko Besic⁵, and Peter Ritchie¹

¹Victoria University of Wellington

²The New Zealand Institute for Bioeconomy Science Limited (Bioeconomy Science Institute – formerly Plant & Food Research)

³Victoria University of Wellington School of Biological Science

⁴Earth Sciences New Zealand

⁵PHF Science (New Zealand Institute for Public Health and Forensic Science)

April 03, 2026

Abstract

++slugcomment: Draft version

Seriola lalandi comprises three genetically distinct clades, with *S. lalandi lalandi* inhabiting the Southern Hemisphere. While a chromosome-level genome exists for the Northwestern Pacific *S. lalandi aureovittata*, no such resource is available for *S. lalandi lalandi*, limiting genetic research and aquaculture advancements. This study presents a high-quality chromosome-level genome assembly of *S. lalandi lalandi* using Illumina short-read, long-read (ONT) and Hi-C. A total of 289.12 Gb of raw genomic data was generated, and a final assembly of 644 Mb was produced with 99.64% of the assembled bases anchored into 24 pseudo-chromosomes. Genome annotation identified 24,600 protein-coding genes, with 96.1% completeness based on 3,640 BUSCO orthologs. Additionally, the complete mitogenome was assembled and annotated for *S. lalandi lalandi*. This genomic resource will help to inform comparative genomic studies, improve aquaculture breeding programs, and support appropriate fishery strategies, particularly in Aotearoa New Zealand, where *S. lalandi lalandi* is a culturally significant taonga species for Māori.

1 The complete chromosome-level genome and
2 mitogenome assembly of *Seriola lalandi lalandi*
3 (Yellowtail Kingfish) for aquaculture and fisheries
4 management.

5 Carla H. Finn^{*1}, Vinko Besic², Tom Oosting¹, Yvan Papa¹, Alvin
6 Setiawan³, Maren Wellenreuther^{4,5}, David Chagné⁴, and Peter Ritchie¹

7 ¹Te Herenga Waka – Victoria University of Wellington, Aotearoa New
8 Zealand

9 ²New Zealand Institute for Public Health and Forensic Science, Aotearoa
10 New Zealand

11 ³Earth Sciences New Zealand, Aotearoa New Zealand

12 ⁴Bioeconomy Science Institute, Aotearoa New Zealand

13 ⁵The University of Auckland, School of Biological Sciences, Aotearoa New
14 Zealand

15 **Abstract**

16 *Seriola lalandi* comprises three genetically distinct clades, with *S. lalandi lalandi*
17 inhabiting the Southern Hemisphere. While a chromosome-level genome exists for
18 the Northwestern Pacific *S. lalandi aureovittata*, no such resource is available for
19 *S. lalandi lalandi*, limiting genetic research and aquaculture advancements. This
20 study presents a high-quality chromosome-level genome assembly of *S. lalandi lalandi*
21 using Illumina short-read, long-read (ONT) and Hi-C. A total of 289.12 Gb of raw
22 genomic data were generated, and a final assembly of 644 Mb was produced with
23 99.64% of the assembled bases anchored into 24 pseudo-chromosomes. Genome
24 annotation identified 24,600 protein-coding genes, with 96.1% completeness based
25 on 3,640 Benchmarking Universal Single-Copy Orthologues (BUSCO). Additionally,

*Corresponding author: carla.finn@vuw.ac.nz, Victoria University of Wellington Level 2, Te Toki a Rata Building, TTR 206, Gate 7, Kelburn Parade, Wellington, NZ 6012

26 the complete mitogenome was assembled and annotated for *S. lalandi lalandi*. This
27 genomic resource will help to inform comparative genomic studies, improve aquacul-
28 ture breeding programmes, and support appropriate fishery strategies, particularly
29 in Aotearoa New Zealand, where *S. lalandi lalandi* is a culturally significant species
30 for Māori.

31 **Keywords:** Genome; assembly; *Seriola lalandi*; fisheries; aquaculture

32 1 Introduction

33 *Seriola lalandi* (Figure 1) is a circumglobally distributed marine pelagic piscivore. Genetic
34 analyses have identified three distinct clades, each corresponding to a geographic region:
35 the Northwestern Pacific (*S. lalandi aureovittata*), the Northeastern Pacific (*S. lalandi*
36 *dorsalis*), and the Southern Hemisphere (*S. lalandi lalandi*) (Premachandra et al., 2017).
37 This population structure is not only relevant to understand the evolutionary history of
38 the species, but also has practical implications for aquaculture, since farmed stocks are
39 typically derived from wild populations specific to the region. Collectively, these clades
40 form an important component of global aquaculture production, owing to their rapid
41 growth and high market value (Symonds et al., 2012).

42 In aquaculture, high-quality reference genomes are a critical resource for selective
43 breeding and genetic improvement programmes. They enable the discovery of genetic
44 markers associated with economically important traits such as growth rate, feed efficiency,
45 flesh quality, and disease resistance (Houston et al., 2020; Yáñez et al., 2023). While a
46 chromosome-level reference genome assembly exists for the Northwestern Pacific *S. lalandi*
47 *aureovittata* population (Li et al., 2022), no equivalent assembly is available for Southern
48 Hemisphere *S. lalandi lalandi*. This gap is significant because the use of a reference genome
49 from a genetically divergent population can introduce reference bias—systematic errors
50 in read mapping and variant calling that arise when sample genomes differ substantially
51 from the reference sequence (Günther and Nettelblad, 2019). Such biases can obscure
52 true genetic variation and reduce the accuracy of downstream analyses. Recent studies
53 underscore the value of contemporary and population-specific reference genome assembly
54 in reducing reference bias, provided genome quality is not compromised (Thorburn et al.,
55 2023).

56 Given these considerations, a high-quality de novo reference genome assembly using a
57 contemporary *S. lalandi lalandi* individual is preferable for accurate and confident studies
58 of the Southern Hemisphere population. Such an assembly will capture the unique genetic
59 characteristics, diversity, and adaptations of the local population, providing a critical
60 foundation for selective breeding, disease resistance research, and stock management
61 in aquaculture sourcing the local population. While this is particularly valuable for

62 advancing sustainable aquaculture practices in the Southern Hemisphere, it also has
63 broader conservation implications. Notably, as climate change drives rapid distributional
64 shifts, hybridization across *Seriola* species is becoming more common (Takahashi et al.,
65 2021), further underscoring the need for population-specific genomic resources to investigate
66 and monitor population-level changes.



Figure 1: A Yellowtail kingfish from Northland, Aotearoa New Zealand. Source: Photograph © Lukas Phan-huy, via iNaturalist NZ, licensed under CC BY-NC 4.0. Observation date: 28 September 2024, Northern Arch, Poor Knights Islands, New Zealand. No changes made.

67 1.1 Cultural acknowledgement and context

68 1.1.1 Naming of the genome — Haku Raukura

69 The individual used to generate this genome was sourced from Te Akau (Bream Bay
70 — Ruakākā, Aotearoa New Zealand), which is within the traditional territorial region
71 (‘rohe’) of Patuharakeke Te Iwi Trust hapū (kinship group). Under New Zealand’s Treaty
72 framework and established Indigenous data governance principles, Patuharakeke retain
73 custodial authority (‘kaitiaki’) over biological resources originating from their region and
74 over genomic data derived from those resources. We acknowledge the guidance of Te Pou

75 Taiao o Patuharakeke Te Iwi Trust throughout the research process, and are grateful
76 that the final *S. lalandi lalandi* genome has been named Haku Raukura by Patuharakeke
77 Pou Ahurea Roopu, a collective of kaumatua (elders) and tohunga (experts). The name
78 reflects a narrative of chiefly strength and spiritual significance: haku, the Māori name
79 for *S. lalandi lalandi*, and raukura refers to a feather, plume, or treasure, traditionally
80 worn by individuals of high rank. The raukura, as a symbol of mana (power/status) and
81 prestige, is likened to the haku – an ika (fish) renowned for its strength, agility, and
82 fighting spirit. The name Haku Raukura thus embodies a fusion of physical prowess
83 and cultural reverence, affirming the genome’s status as a taonga tuku iho (treasured
84 inheritance) and symbolises the indelible link between haku and tangata (people).

85 **1.1.2 Data availability and sovereignty**

86 The *S. lalandi lalandi* individual used in this study was sourced from within the rohe
87 (region) of Patuharakeke iwi, hapū and whānau (tribes, kinship groups and families) .
88 Therefore, Patuharakeke rightfully holds responsibilities of kaitiakitanga (guardianship)
89 over their rohe and, by extension, over the species and individuals that originate from it
90 — including the individual used in this study, and including the raw and analysed data
91 derived from it. Patuharakeke’s authority over their *S. lalandi lalandi* and resulting *Haku*
92 *Raukura* (genome assembly) gives effect to the Treaty of Waitangi, a treaty signed in 1840
93 between Māori chiefs and the British Crown. Patuharakeke Te Iwi Trust have expressed
94 that the use of such taonga (treasure) must align with the principles of Wai 262 / Ko
95 Aotearoa Tēnei, which affirms the rights of Māori to protect and control their cultural
96 heritage and biological resources. Therefore, the data is made available *by request* on the
97 Aotearoa Genomics Data Repository (AGDR) (Te Aika et al., 2023), with permissions
98 granted by the Patuharakeke Te Iwi Trust. The AGDR has been developed to follow the
99 principles of Māori Data Sovereignty, and to enable kaitiakitanga (guardianship), so that
100 can effectively exercise their responsibilities as guardians over biological entities that are
101 taonga (treasured). Internationally, Māori ownership and control of natural resources are
102 supported by the UN Declaration on the Rights of Indigenous Peoples 2007.

2 Materials and methods

2.1 Sampling and tissue extraction

Whole blood and liver tissue was collected from a captive bred *S. lalandi lalandi* (haku/yellowtail kingfish) (sourced by NIWA at Te Akau/Bream Bay, Ruakākā, Aotearoa New Zealand) individual (885g, 38cm, 344 days old). Tissue samples were obtained 0.5–1.5 hours post mortem following anaesthetic overdose. High-molecular-weight DNA was extracted from whole blood using a high-salt extraction protocol adapted from (Aljanabi and Martinez, 1997; Oosting et al., 2020) employing the Nanobind CBB Big DNA Kit (Circulomics). Prior to long-read sequencing, the integrity of DNA fragments was assessed by gel electrophoresis in 1% agarose. DNA quantity was also assessed using the Quant-iT™ PicoGreen™ dsDNA assay kit, while purity and quality was assessed using the NanoDrop (ThermoFisher) and a Genomic DNA ScreenTape (Agilent), respectively. DNA samples had concentrations > 200 ng/ μ l, $A_{260/280} \approx 1.8$, $A_{260/230} \approx 2$, and total yields > 20 μ g.

2.2 Sequencing

Long-read sequencing was performed at PHF Science (the New Zealand Institute for Public Health and Forensic Science) using Oxford Nanopore (ONT) sequencing on the MinION (SQK-LSK112) and PromethION (SQK-LSK114) instruments with R10 chemistry, generating 94.18 Gb of raw ONT sequence data. Additionally, genomic DNA was also sent to AGRF (Australian Genome Research Facility, Melbourne, Australia) for Illumina paired-end sequencing on the NovaSeq, with 150bp read lengths, generating 79.42Gb of raw Illumina data. Liver tissue was sent to BGI Tech Solutions Co., Ltd. (Hong Kong, China) for genomic DNA extraction, Hi-C library preparation, and sequencing on the DNBseq platform with paired-end (PE) reads of 2×150 bp. This resulted in 115.32Gb of data with a Q20(%) of 97.8.

2.3 Quality, contamination and mitochondrial filtering

2.3.1 ONT reads

DNA was sequenced using both the PromethION and MinION platforms, and base-calling was performed using different algorithms based on ONT’s software compatibility: The MinION sequencing data were base-called using Guppy (version 6.3.4+cfaa134) (Oxford Nanopore Technologies) and the PromethION sequencing data were base-called using Dorado (version 0.3.1) (Oxford Nanopore Technologies). After base calling, duplex reads were extracted from the ONT dataset yielding a total of 1,191,939 duplex pairs. Adapters were removed using Guppy (version 6.3.4+cfaa134, Oxford Nanopore Technologies) for MinION data, and Porechop (version 0.2.4) (Wick and Volkening, 2018) for PromethION

137 data. Further methodological details are available in the Appendix for methods. The 3kb
138 lambda DNA control sequence was removed using `clean` (v0.2.0 via nextflow 23.04.4.5881)
139 (Lataretu et al., 2023). Since all duplex reads had a Phred quality score ≥ 7 , no additional
140 quality control filtering was applied to this dataset; however, simplex data were filtered
141 and trimmed using `Chopper` (version 0.6.0) (De Coster and Rademakers, 2023) specifying
142 options `-headcrop 50` and `-qc 7`.

143 After filtering, the PromethION simplex and duplex reads, as well as the MinION
144 simplex and duplex reads, were concatenated for downstream assembly pipelines. The
145 final ONT dataset had a mean read length of 8,382bp and mean read quality of Q18.8
146 (Table 1)

147 **2.3.2 Illumina reads**

148 Adapters were removed from the Illumina reads using `Trimmomatic` (version 0.39) (Bolger
149 et al., 2014). `Kraken2` (version `kraken/2.0.7-beta`) (Wood and Salzberg, 2014) was used
150 to detect, and remove, sequences from archaea, bacteria, viruses, and humans, using the
151 `MiniKraken2 v2 8GB` database. This step resulted in 233,736,388 read pairs remaining
152 with an average mean read quality of Q34.83 (Table 1).

153 **2.4 Genome size, coverage, and heterozygosity estimation**

154 Genome size and sequencing coverage based on the filtered Illumina sequence reads was
155 performed using `Jellyfish` (v2.2.10) (Marcais and Kingsford, 2012) and analysed and
156 visualised using `GenomeScope` (Vurture et al., 2017).

157 **2.5 Assemblies and annotation**

158 For a general illustrative overview of the genome assembly pipeline, please refer to Figure
159 1, below.

160 The Maryland Super Read Cabog Assembler, `MaSuRCA` (version 3.2.9) (Zimin et al.,
161 2013) has been widely used for assembling eukaryotic genomes using hybrid short- and
162 long-read approaches, delivering consistently high-quality results across studies (Tan
163 et al., 2018); indeed, this method was found to be most suitable for our data. Before
164 assembly, the Illumina sequences were not trimmed or edited as per `MaSuRCA` author
165 recommendation (<https://github.com/alekseyzimin/masurca>). Unfiltered Illumina
166 short-reads and filtered ONT reads were used with recommended parameters. Methods,
167 including the configuration file, are provided in the Appendix. The assembly generated
168 was called the “First draft assembly”. The first draft assembly was subsequently put
169 through three iterations of `Pilon` (version 1.23) (Walker et al., 2014) using the `-frags`
170 and `-unpaired` options (i.e., a hybrid polishing step). The option `-fix-all` was also

Table 1: Overview of filtering steps applied to Illumina and Oxford Nanopore data. CR = control region. Illumina statistics assessed via MultiQC (version 1.15) (Ewels et al., 2016); ONT statistics via NanoPlot (version 1.32.0) (De Coster and Rademakers, 2023)

Platform	Instrument	Method	Condition	No. of reads	Bases (Gbp)	Mean length
Illumina	NovaSeq	-	Primary QC	262,988,602 pairs	79	150.0
	-	-	Adapter trimmed, contamination removed	233,736,388 pairs	70	149.97
Oxford Nanopore	PromethION	Simplex	Raw	9,632,490	74.69	7,744.9
		Simplex	CR removed, adapter trimmed, QC filtered	1,707,688	12.24	7,169.3
Oxford Nanopore	MinION	Duplex	Raw	544,717	4.30	7,888.1
		Duplex	CR removed, adapter trimmed	26,783	0.180	6,540.4
		Simplex	Raw + adapter trimmed	7,381,142	61.43	8,332.2
		Simplex	CR removed, adapter trimmed, QC filtered	1,207,050	10.46	8,663.8
		Duplex	Raw + adapter trimmed	480,403	4.07	8,465.0
		Duplex	CR removed	21,598	0.164	7,600.0
Final ONT reads (PromethION + MinION)				90,090,211	76.12	8,373.4

171 used to ensure a comprehensive correction of assembly errors. To produce a haploid-
172 collapsed reference genome suitable for downstream read mapping and variant-based
173 analyses, `PurgeHaplotigs` (version 1.1.2) (Roach et al., 2018) was used to identify and
174 remove alternative haplotigs. Parameters (`-l 3 -m 50 -h 200`) were selected based on
175 the observed coverage distribution in the graphical histogram (Appendix Figure S1).
176 Read depth for haplotig identification was estimated using Illumina short-read mappings,
177 as these provide more uniform and accurate coverage profiles than long-read data for
178 distinguishing haploid and diploid contigs.

179 The assembly generated is referred to as the “Second draft assembly”.

180 For Hi-C scaffolding, the Second draft assembly was indexed and aligned to Hi-C short
181 reads using `BWA`, `Burrows-Wheeler alignment` (version 0.7.17-r1188) (Li and Durbin,
182 2009). PCR duplicates were flagged with `SAMBLASTER` (version 0.1.26) (Faust and Hall,
183 2014) using default settings, and non-primary and unmapped aligned reads were removed
184 using `Samtools` (version 1.17) (Danecek et al., 2021). Finally, `YaHS`, `Yet another Hi-C`
185 `scaffolding tool` (version YaHS-1.1) (Zhou et al., 2023) was used to order and orientate
186 the Second draft assembly. The minimum read mapping quality was set to 20 (`-q 20`)
187 and `-e GATC` (*DpnII* restriction enzyme) was specified.

188 The resulting scaffolded genome is referred to as the “Scaffolded assembly”.

189 The statistics for each iteration of the assembly culminating to the Haku Raukura
190 (final) assembly are found in Appendix Table S1.

191 2.5.1 Contamination screening

192 Contamination was evaluated using `NCBI Foreign Contamination Screen (FCS)` mod-
193 `ule FCS-GX` (version 0.4.1) (Astashyn et al., 2023). `FCS-GX` classifies sequences as contam-
194 inant when their taxonomic assignment is different from the user provided taxonomic
195 identifier; for this assessment, the NCBI tax ID `-tax-id 302047` was used, which identifies
196 *Seriola lalandi*. No contamination was identified.

197 2.6 Repeat masking and protein annotation

198 Repetitive elements in the Haku Raukura assembly were identified using both de novo mod-
199 elling and homologue-based approaches. For de novo repeat identification, `RepeatModeler`
200 (version 2.0.1) (Flynn et al., 2020) was used with the `-LTRStruct` parameter enabled.
201 For homology-based repeat identification, known *Actinopterygii* repeats were retrieved
202 using the `famdb.py` utility of `RepeatMasker` (version 4.1.1) (Smit et al., 2015) from
203 a combined database comprising `Dfam` (version 3.8) (Storer et al., 2021) and `RepBase`
204 `RepeatMasker Edition (v20181026)` (Bao et al., 2015) specifying the parameters fam-
205 ilies `-include-class-in-name -ad -add-reverse-complement Actinopterygii`. The
206 repeat libraries generated by de novo modelling and homology-based identification were

207 concatenated to create a custom repeat library specific to *S. lalandi lalandi*. Genome assem-
208 bly sequences were then mapped to this library using RepeatMasker (v4.1.1) (`-xsmall`) to
209 classify repeat regions and produce a repeat annotation file. This process also generated a
210 soft-masked genome assembly. A hard-masked version was created by replacing lowercase
211 bases with 'N' characters.

212 **2.7 Protein sequence prediction and functional annotation**

213 Protein sequences were predicted using BRAKER3 (version v3.0.8 via singularity (version
214 3.7.3) (Gabriel et al., 2021; Stanke et al., 2006) on the soft-masked Haku Raukura
215 assembly. Homology-based gene predictions were then conducted with the OrthoDB v11
216 Vertebrates database (Kuznetsov et al., 2023), improving the accuracy of the predicted
217 models through conserved domain and orthology information. Functional annotation was
218 carried out in two sequential steps. First, protein predictions were aligned against the
219 UniProtKB/Swiss-Prot database (Apweiler et al., 2004) using BLAST+ (version 2.6.0)
220 (Camacho et al., 2009). An E-value cut-off of 1×10^{-6} was used (`-evalue 1e-6`), and only
221 the top-scoring alignment was retained by specifying `-max_hsps 1 -max_target_seqs`
222 `1`. Second, InterProScan (version 5.68-100.0) (Blum et al., 2021; Jones et al., 2014)
223 was employed to identify functional domains and additional annotations, specifying the
224 Pfam database (`-appl pfam`), Gene Ontology (GO) terms (`-goterms`), and pathway
225 annotations (`-dp`).

226 The resulting functional annotations were integrated into gene models and refined
227 using the AGAT toolkit (version 1.4.0) (Dainat, 2022). The detailed methodology for
228 this step is available in the Appendix for methods.

229 **2.8 Mitogenome sequence**

230 The mitogenome was identified by aligning a publicly available *Seriola lalandi* mitogenome
231 (*Seriola lalandi* isolate Serio_lalandi_NSW mitochondrion GenBank accession number
232 MW309824.1) to a ONT-only version of the *S. lalandi lalandi* assembly via blastn (version
233 BLAST+/2.14.1) (Camacho et al., 2009). A contig with >98% match was found. It was
234 annotated and visualised using MitoAnnotator (Zhu et al., 2023).

235 **2.9 Synteny**

236 To investigate the genome synteny between *S. lalandi lalandi* and *S. lalandi aureovittata*
237 (RefSeq assembly GCF_021018895.1 (ASM2101889v1), we utilised MUMMER4 (version 4.0.0.)
238 (Marçais et al., 2018). Nucmer was used for whole-genome alignment, specifying *S. lalandi*
239 *aureovittata* (RefSeq: GCF_021018895.1) as the reference and *S. lalandi lalandi* as the
240 query. To refine the alignment, delta-filter was applied with a minimum identity

241 threshold of 95% (`delta-filter -i 95`) and minimum alignment length of 10,000bp
242 (`delta-filter -l 10000`). Then, `mummerplot` was used to generate the visual results
243 for the synteny alignment, using the options `-filter -fat -layout`. Finally, `dnadiff`
244 was run on the filtered dataset to generate a report of the differences between the genomes
245 of *S. lalandi laladi* and *S. lalandi aureovittata*.

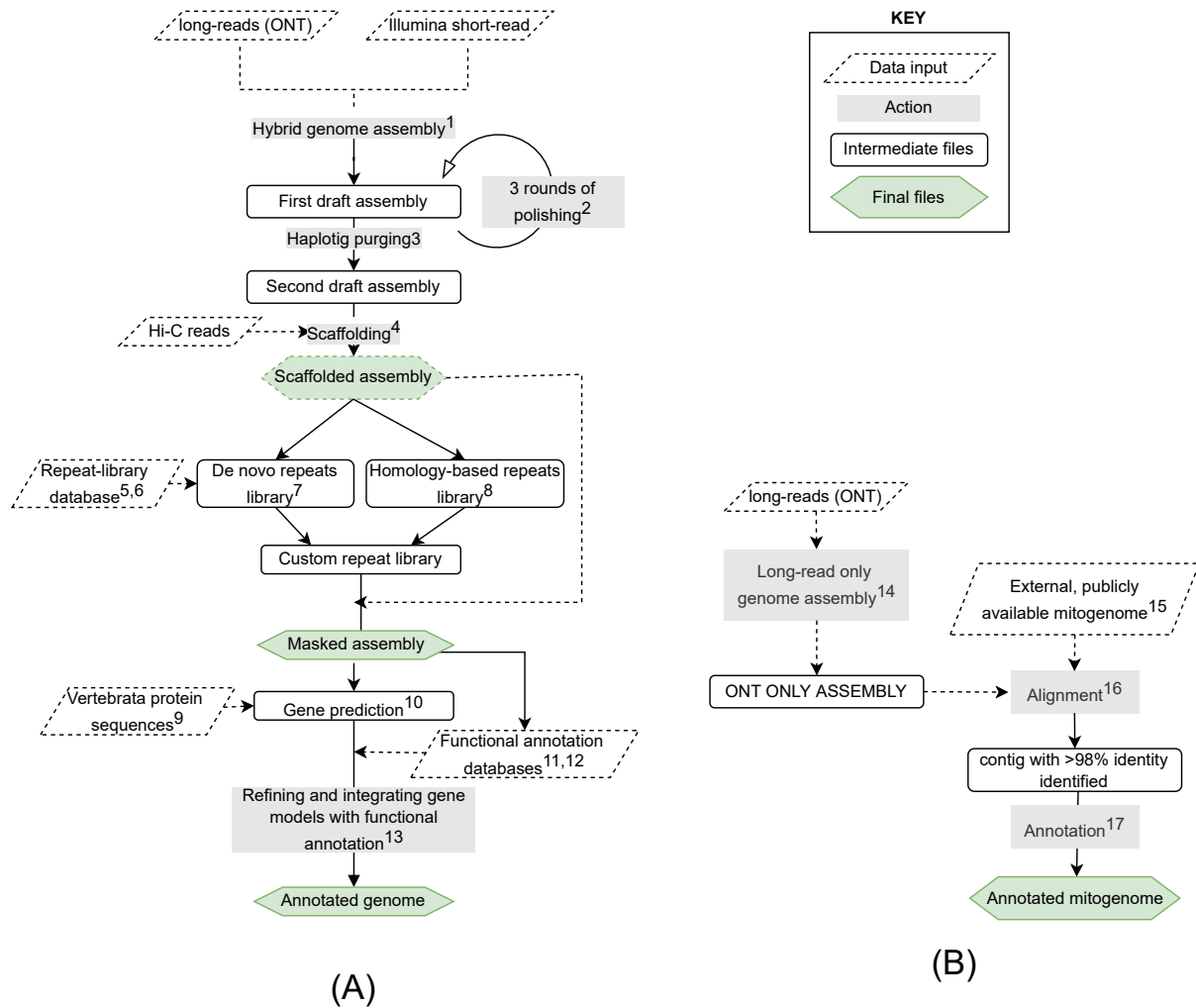


Figure 2: An outline of the assembly methods used to assemble the *Seriola lalandi lalandi* genome. Each step is sequentially represented; with input data in dotted outline, and resulting final data in green hexagons, illustrating the progression from raw data to the final assemblies. (a) The pipeline for developing Haku Raukura, i.e., the chromosome-level assembly for *Seriola lalandi lalandi*. The statistical results for each iteration of the genome assembly are available in Appendix Table S1. (b) The pipeline for extracting and annotating the mitogenome for *Seriola lalandi lalandi*.

¹ MaSuRCA v3.2.9 (Zimin et al., 2013); ² Pilon v1.23 (Walker et al., 2014); ³ Purge Haplotigs v1.1.2 (Roach et al., 2018); ⁴ YaHS v1.1 (Zhou et al., 2023); ⁵ Dfam v3.8 (Storer et al., 2021); ⁶ RepBase v20181026 (Bao et al., 2015); ⁷ RepeatModeler v2.0.1 (Flynn et al., 2020); ⁸ RepeatMasker v4.1.1 (Smit et al., 2015); ⁹ OrthoDB v11 (Kuznetsov et al., 2023); ¹⁰ Gene prediction (Gabriel et al., 2021); ¹¹ SwissProt database (Apweiler et al., 2004); ¹² InterProScan 5.68-100.0 (Jones et al., 2014); ¹³ AGAT Toolkit v1.4.0 (Dainat, 2022); (B) ¹⁴ FLYE v2.9.2-b1795 (Kolmogorov et al., 2019); ¹⁵ *Seriola lalandi* isolate *Seriola_lalandi_NSW* mitochondrion, complete genome (GenBank accession MW309824.1); ¹⁶ BLAST+ v2.14.1 (Camacho et al., 2009); ¹⁷ MitoAnnotator (Zhu et al., 2023).

246 3 Results

247 3.1 Genome size, coverage, and heterozygosity estimate post 248 sequencing

249 Genome size, coverage, and heterozygosity were estimated from Illumina short-read
250 data using k-mer-based analysis. Among the k-mer sizes evaluated, the 27-mer model
251 demonstrated the best fit, with model accuracy ranging from 92.48% to 97.01% (Table
252 2). Based on this model, the genome size was estimated to be approximately 588Mbp.
253 Heterozygosity was calculated at 0.43%–0.44%, indicating a relatively low level of genetic
254 variation within the genome. The heterozygous coverage of x42.9 was considered sufficient
255 for performing genome assembly.

Table 2: 3 K-mer Analysis Statistics. The model based on 27-mers demonstrated the highest model fit, suggesting it as the most reliable for genome size estimation.

K	Total number of kmers	Erroneous kmers	Estimated genome size	Heterozygosity		Model fit	
				Min	Max	Min	Max
17	61,092,433,257	2,210,987,433	282,645,767	5.97%	17.74%	83.86%	91.45%
21	59,707,868,384	5,057,318,169	587,342,308	0.48%	0.49%	90.75%	96.62%
27	57,255,261,346	5,939,877,545	588,852,038	0.43%	0.44%	92.48%	97.01%

256 3.2 Final genome assembly

257 The final scaffolded assembly consisted of a total length of 644,985,034 bp in 146 contigs.
258 Visualisation of the Hi-C contact map confirmed 24 anchored scaffolds (Figure 2). These
259 24 scaffolds corresponded to 99.64% of the scaffolded assembly by base count, with an
260 average length of 26,779.82 kb. The unanchored 122 scaffolds were shorter with an average
261 length of 18.8kb, covering only 0.36% of the assembly. The 24 anchored scaffolds totalled
262 642 Mb, with an N50 scaffold length of 34 Mbp, N50 scaffold count of 11 and a complete
263 BUSCO (Benchmarking Universal Single-Copy Orthologs) score of 98.6% (BUSCO (version
264 5.5.0) using the *actinopterygii_odb10* database (Table 3)).

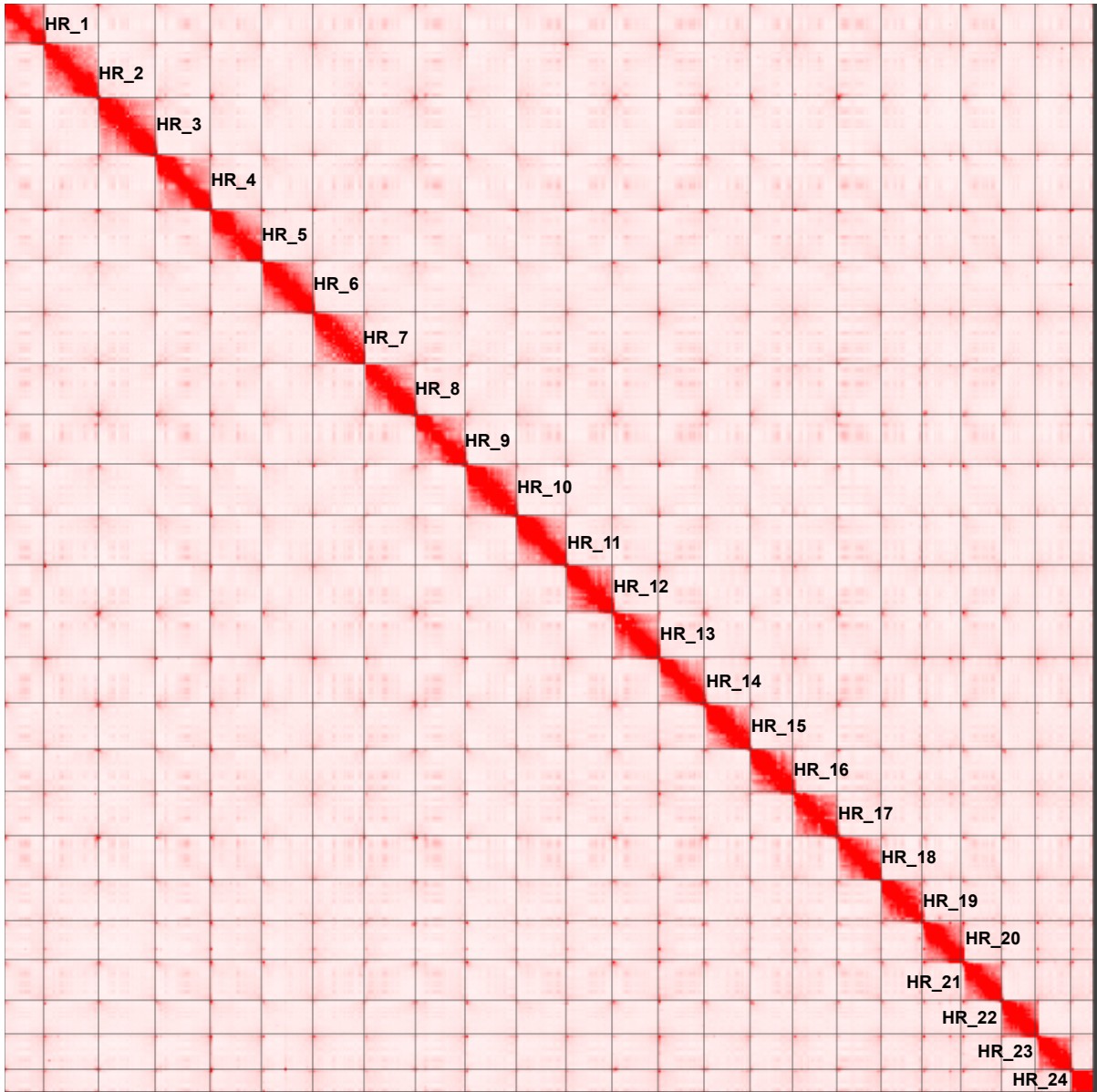


Figure 3: Hi-C contact map of Haku Raukura (*S. lalandi lalandi*). This heatmap illustrates the interaction frequencies between different regions within the genome. Interactions are depicted by the intensity of the red colouring, with darker shades indicating higher frequencies of contact between loci. The matrix is organised such that both the x-axis and y-axis represent the linear progression of genomic coordinates.

Table 3: Statistics for the Haku Raukura (*S. lalandi lalandi*) assembly. BUSCO - Benchmarking Universal Single-Copy Orthologs

Metric	Value
Span (bp)	642,715,660
N (%)	-
GC (%)	40.81
AT (%)	59.19
Scaffold count	24
Longest scaffold (bp)	34,161,391
Scaffold N50 length (bp)	28,040,511
Scaffold N50 count	11
Scaffold N90 length (bp)	22,307,305
Scaffold N90 count	21
Complete BUSCOs (C)	98.60%
Complete and single copy BUSCOs (S)	98%
Complete and duplicated BUSCOs (D)	0.70%
Fragmented BUSCOS (F)	0.20%
Missing BUSCOS (M)	1.20%

265 **3.3 Repeat masking**

266 The results identified 20Mb of repeats and included a total of 16.6Mb interspersed repeats
 267 (18.56% of the genome), 3.04Mb satellite repeats (3%) and 1Mb of small RNAs (0.16%).

268 In total, 22% of the genome was masked for repeats (Table 4).

Table 4: Summary of repeat elements in the Haku Raukura (*S. lalandi lalandi*) assembly.

Repeat elements	Copies	Bases	Coverage of genome (%)
Interspersed repeats			
SINES	15,039	1,393,562	0.22%
Penelope	80	6,580	0.00%
LINE	40,018	8,394,254	1.31%
LTR	8,394,254	2,763,167	0.43%
DNA transposons	2,763,167	2,763,167	3.92%
Rolling circles	2,469	814,453	0.13%
Unclassified	549,317	549,317	12.55%
Total interspersed repeats	11,764,344	16,684,500	18.56%
Satellite repeats			
Simple repeats	411,847	411,847	3.02%
Satellites	1,346	134,950	0.02%
Low complexity	45,848	2,495,677	0.39%
Total satellite repeats	459,041	3,042,474	3%
Small RNA	11,972	1,056,110	0.16%
Total	12,235,357	20,783,084	22%

269 3.4 Protein sequence prediction and functional annotation

270 Protein sequence prediction identified a total of 38,487 gene models and 42,264 mRNA
271 models for Haku Raukura. The mean gene length was 6,610bp and the mean mRNA
272 length was 7,144bp; the mean exon length was 174bp and the mean intron length in in
273 CDS (coding sequence) was 853bp. Functional annotation identifies 24,600 catalogued
274 genes. Of the predicted gene models, 24,600 (64%) were named. For the total predicted
275 gene models, 96.1% of 3640 actinopterygii conserved genes are complete (BUSCO 5.5.0
276 *actinopterygii_odb10* database). Similarly, 27,861 (66%) mRNA models were named,
277 while 14,403 (34%) were unnamed. A total of 21,143 genes were assigned to at least one
278 of the 5 queried databases (Table 5). Those remaining unnamed or ‘hypothetical’ likely
279 represent hypothetical or poorly characterised genes.

Table 5: Annotation of Haku Raukura *S. lalandi lalandi* genes to different databases.

Type	Database	Assigned gene number
Homologue	Gene ontology	14,135
	InterPro	22,614
	MetaCyc	12,473
	Pfam	23,118
	Reactome	20,247

280 3.5 Synteny

281 The 24 pseudo-chromosomes of *S. lalandi lalandi* had clear one-to-one relationship to *S.*
 282 *lalandi aureovittata* (RefSeq: GCF_021018895.1) pseudo-chromosomes, with 82% of bases
 283 aligning (Figure 3). A total of 16,613 one-to-one alignments were identified, exhibiting
 284 an average alignment length of approximately 32kb and a high average sequence identity
 285 of 98.44%. The analysis detected 33,310 breakpoints, along with relocations (36 in the
 286 reference vs. 27 in the query), translocations (13 in both assemblies), and inversions (3 in
 287 the reference vs. 4 in the query). In total, 3 million SNPs and 5.29 million indels were
 288 observed in each assembly, highlighting both the extensive synteny and the presence of
 289 structural rearrangements between these two closely related *Seriola* genomes.

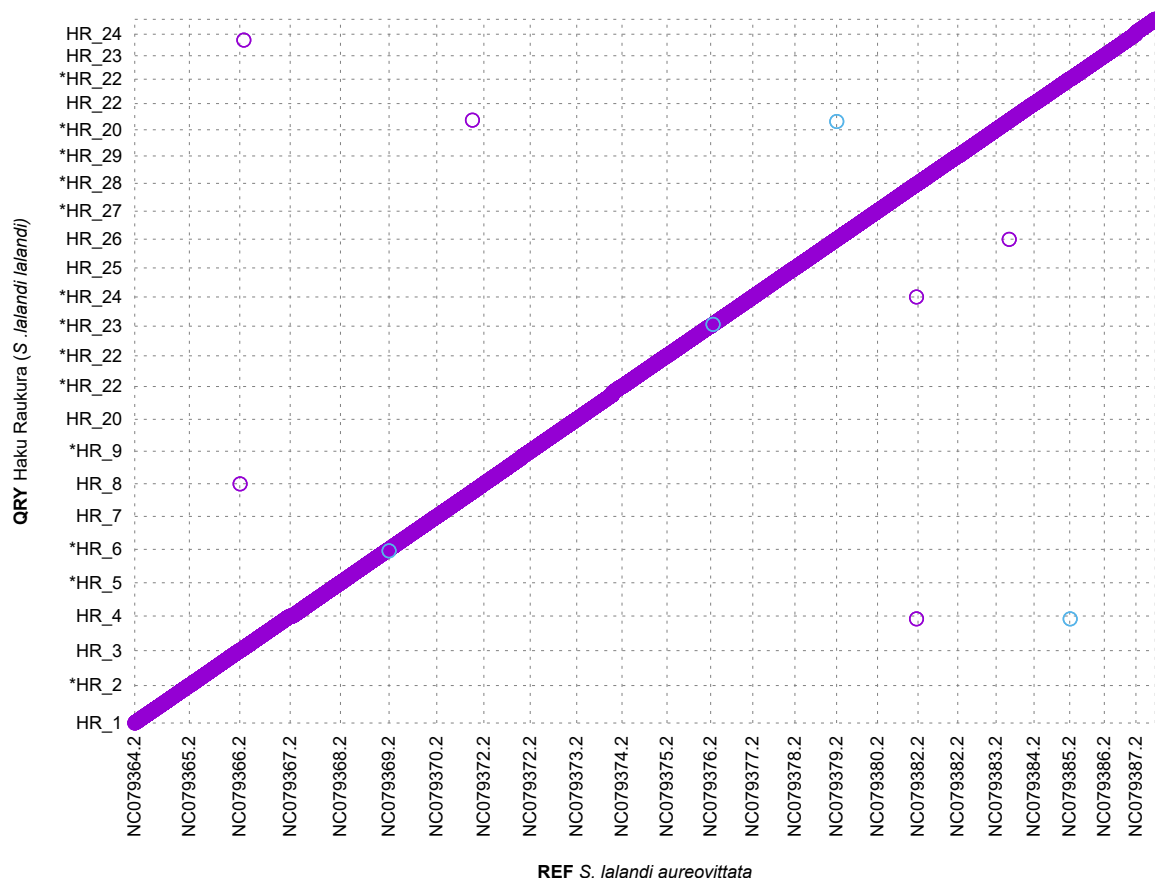


Figure 4: Whole genome plot of *S. lalandi aureovittata* (RefSeq: GCF_021018895.1) (reference, x axis) compared with the Haku Raukura assembly (query, y axis). Each colour dot/line indicates a match between the reference and query. Forward mappings are purple; reverse mappings are blue. Asterisked chromosomes indicate reverse-complemented sequences (i.e., sequences that have been flipped in orientation).

290 **3.6 Mitogenome results**

291 This contig identified was 16,535bp in length. It was annotated and visualised using
 292 MitoAnnotator (Zhu et al., 2023; Sato et al., 2018; Iwasaki et al., 2013). There were 13
 293 genes: 2 rRNAs and 22 tRNAs. The overall base composition was as follows: A, 26.6%;
 294 C, 30.31%; G, 17.86%; T, 25.23%. The mitogenome is shown in Figure 4.

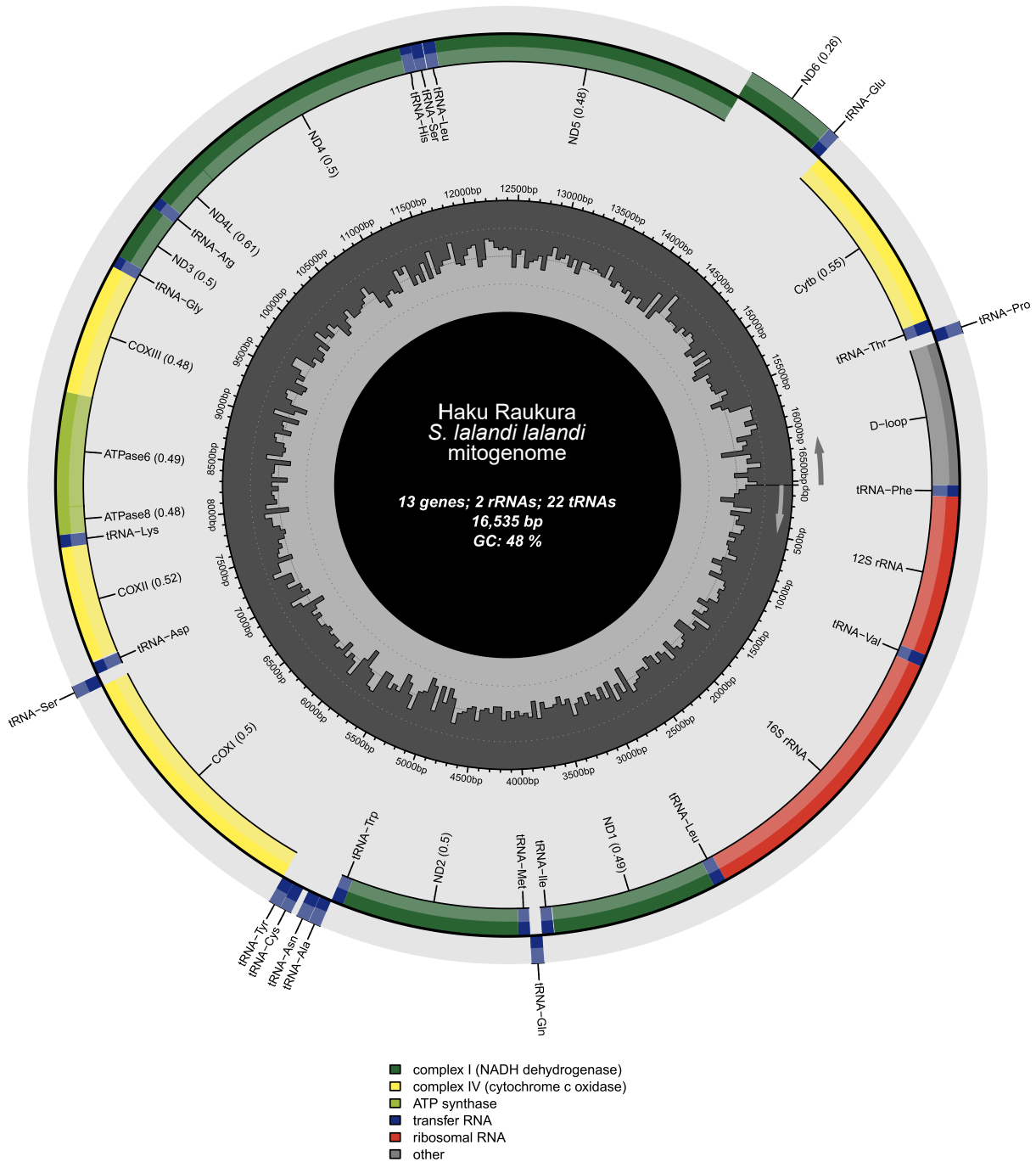


Figure 5: The mitogenome for Haku Raukura *S. lalandi lalandi*.

4 Discussion

4.1 Genome quality and relevance

In the context of other publicly available *Seriola* genomes, the final Haku Raukura (*S. lalandi lalandi*) assembly presented here demonstrates high contiguity and completeness, placing it among the more robust assemblies for the genus. Our scaffold N50 of 28 Mb is comparable to the other high-quality *Seriola* assemblies released thus far (e.g., 28.5 Mb for the Northwestern *S. lalandi aureovittata* (RefSeq: GCF_021018895.1) and substantially higher than assemblies of other *Seriola* species (e.g., *S. quinqueradiata* with 5.6 Mb). Furthermore, BUSCO analysis (`actinopterygii_odb10`) indicates a completeness of 98.6%, which is on par with or only marginally below the highest reported values (e.g., 99.4% for *S. lalandi aureovittata* and 99.8% for *S. dumerili*). This indicates that the assembled pseudochromosomes for *S. lalandi lalandi* cover nearly all expected core orthologues. The total repeat content (22%) aligns with that reported for other Carangidae genomes (20–22%), such as 20.24% for *Trachinotus ovatus* (Zhang et al., 2019) and 22.46% for *S. lalandi aureovittata* (Li et al., 2022), suggesting broadly consistent repetitive architecture across related species. Whilst the 27-mer model estimated the genome size to 588Mbp — which is notably lower than the c-value (pg) of 0.7 according to the Animal Genome Size Database (Animal Genome Size Database: Species Record id=2705)— it is common for k-mer estimated size and genome assembly size to be smaller than the size estimated with C-value (Pflug et al., 2020), so this was not considered a problem. Heterozygosity was calculated at 0.43%–0.44%, indicating a relatively low level of genetic variation within the genome. This is not an uncommon range; other members of the Carangidae demonstrate comparable levels of diversity, such as 0.31% in *Trachinotus ovatus* (Zhang et al., 2019) and 0.71% in *Pseudocaranx georgianus* (Catanach et al., 2021). These results provide a baseline estimate of standing genetic variation in *S. lalandi lalandi* and reinforce the value of a population-specific reference genome for downstream analyses.

Our syntenic analyses reveal that *S. lalandi lalandi* shares 98% sequence identity with *S. lalandi aureovittata*, yet exhibits structural rearrangements, numerous breakpoints, and unique SNPs/indels. These findings underscore the genomic distinctiveness of the Southern Hemisphere population and highlight the importance of a local reference for capturing population-specific variation accurately, and indeed, recent work highlights the benefits of using population-specific reference genomes to reduce reference bias, provided assembly quality remains high (Thorburn et al., 2023). Consistent with this low heterozygosity (0.43–0.44%) and the presence of a single dominant coverage peak, the final Haku Raukura assembly was generated as a haploid-collapsed reference, in which alternative haplotypes were intentionally removed. This avoided artificial duplication of allelic contigs and yielded a non-redundant genome representation. Such an approach is appropriate for

332 the primary applications of this resource, including read mapping, variant discovery,
333 comparative genomics, and aquaculture breeding programmes — but notably, future
334 phased or pangenome assemblies could complement this resource by explicitly representing
335 haplotypic diversity.

336 **4.2 Implications for aquaculture**

337 This reference genome greatly facilitates the identification of SNPs and indels in eco-
338 nomically important genes across breeding populations of *S. lalandi lalandi*. In practice,
339 aquaculture programmes can integrate such markers into breeding value prediction models
340 (e.g., genomic selection), enabling earlier and more precise selection for traits that are
341 otherwise challenging to measure (e.g., feed conversion or fillet quality). Additionally,
342 the Haku Raukura genome has already proven useful for sex determination applications:
343 a sex-identification SNP in Hsd17b1 on chromosome 17 (position 15,001,654) has been
344 pinpointed (Earth Sciences New Zealand, personal communication), allowing faster and
345 more reliable sexing of juvenile fish.

346 **4.3 Implications for fisheries management**

347 The *S. lalandi lalandi* genome is also valuable for managing wild fisheries. Genomic
348 data can complement traditional stock assessments by revealing fine-scale population
349 structure, genetic diversity, and cryptic stock boundaries that may not be evident from
350 morphological or tagging data. For instance, genome-wide analyses of the related *S.*
351 *dumerili* identified three genetic groups (one Mediterranean, two Atlantic), where only two
352 had been presumed previously (Katirtzoglou et al., 2024), suggesting that genome-wide
353 polymorphic markers can enable novel stock delimitations. A similar high-resolution
354 survey across the *S. lalandi lalandi* range (e.g., New Zealand, Australia, other South
355 Pacific populations) could identify sub-populations requiring distinct management. By
356 providing a template for resequencing data and marker discovery, this genome facilitates
357 the development of conservation strategies tailored to specific genetic stocks.

358 **4.4 Conclusion**

359 The high-quality *S. lalandi lalandi* reference genome presented here provides a com-
360 prehensive resource for aquaculture, fisheries management, and conservation. Its high
361 contiguity and completeness make it comparable to the best *Seriola* assemblies to date
362 and reveal genomic features—such as unique structural variants and population-specific
363 SNPs/indels—that underscore the distinctiveness of Southern Hemisphere *S. lalandi*
364 *lalandi*. This resolution enables more precise genomic selection strategies and breeding
365 programmes, supports fine-scale assessments of wild fish populations, and offers valu-

366 able insights into adaptive variation across different environments. As genomic research
367 continues to expand across Carangidae species, this assembly will serve as a critical
368 reference point for comparative studies, informing both scientific and cultural approaches
369 to preserving the genetic health of *S. lalandi lalandi* populations.

370 4.5 Data availability

371 Genomic data generated in this study are available via the Aotearoa Genomic Data
372 Repository under controlled access. Because the sampled individual originated from
373 within the traditional territory of Patuharakeke Te Iwi Trust (Aotearoa New Zealand),
374 access to the resulting genomic data requires approval by that authority (see Section
375 1.1.2 Data availability and sovereignty, for explanation). Data may be requested through:
376 <https://doi.org/10.57748/t37x-tz81> **Reviewer access** is supported through apply-
377 ing for review-based data access using the link above and this will be approved by
378 the kaitiaki for the dataset (further details available at [https://docs.agdr.org.nz/
379 general_information/information_for_publishers/](https://docs.agdr.org.nz/general_information/information_for_publishers/)).

380 The scripts used to generate these data are openly available at [https://github.com/
381 carla-hazelf/hakuRaukura_genomeAssembly_pipeline](https://github.com/carla-hazelf/hakuRaukura_genomeAssembly_pipeline)

382 5 Acknowledgments

383 We thank and acknowledge the guidance of Te Pou Taiao o Patuharakeke Te Iwi Trust
384 throughout the research process, and for the guidance and koha (gift) of knowledge
385 from Patuharakeke Pou Ahurea Roopu, a collective of kaumatua (elders) and tohunga
386 (experts). We thank Juliane Chetham (Resource Management and Customary Fisheries,
387 Patuharakeke Te Iwi Trust) and Tracey Godfrey (Vision Mātauranga Manager, Genomics
388 Aotearoa) for providing assistance with liaison between the research team and Te Pou
389 Taiao o Patuharakeke Te Iwi Trust. We are grateful for the guidance received from the
390 Genomics Aotearoa team. Computational analyses were performed with the support of
391 Raapoi HPC (Te Herenga Waka — Victoria University of Wellington).

392 We thank the Indigenous Genome Project of Genomics Aotearoa, and Te Herenga Waka –
393 Victoria University of Wellington, for providing funding for this project.

394 6 Author contributions

395 Carla H. Finn: Formal analysis (lead), investigation (lead), methodology (lead), writing –
396 original draft (lead), writing – review and editing (lead), Data curation (equal), visualisa-
397 tion (lead)

398 Vinko Besic: Data curation (supporting).

399 Tom Oosting: Data curation (supporting), Writing — review and editing (supporting).
400 Yvan Papa: Data curation (supporting)
401 Alvin Setiawan: Resources (equal)
402 Maren Wellenreuther: investigation (supporting), writing — review and editing (support-
403 ing), funding acquisition(equal), supervision (equal).
404 David Chagné: investigation (supporting), writing — review and editing (supporting),
405 funding acquisition(equal), supervision (equal)
406 Peter Ritchie: investigation (supporting), writing — review and editing (supporting),
407 supervision (lead), funding aquisition(equal).

408 Appendices

409 Appendix - Methods

410 6.0.1 Basecalling, filtering and duplex-calling ONT reads

411 The MinION sequencing data was base-called using Guppy (version 6.3.4+cfaa134, Oxford
412 Nanopore Technologies), with the `dna_r10.4_e8.1` super-accurate model. Adapter trim-
413 ming was enabled during basecalling, specifying `-c dna_r10.4_e8.1_sup.cfg --num_callers 28 --tr`
414 After this, Duplex Tools (version 0.3.3, Oxford Nanopore Technologies) was used to iden-
415 tify and prepare duplex read pairs for high-accuracy base-calling. Duplex read pairs were
416 extracted using `pairs_from_summary`, filtered with `filter_pairs`, and re-basecalled us-
417 ing the options `-c dna_r10.4_e8.1_sup.cfg -num_callers 28 -trim_adapters`
418 `-duplex_pairing_mode from_pair_list`. Raw PromethION FAST5 data was first
419 converted to `.pod5` format using the `Pod5 Python tools` (version 0.2.4, Oxford Nanopore
420 Technologies). This conversion was required for compatibility with Dorado (version 0.3.1).
421 To improve duplex base-calling efficiency, the `.pod5` files were sub-set into per-channel
422 unique `.pod5` files specifying `-subset -columns channel` as recommended by Dorado
423 (<https://github.com/nanoporetech/dorado>). Duplex and simplex reads were then
424 base-called simultaneously using the super-accurate model (`dna_r10.4.1_e8.2_400bps_sup@v4.1.0`).
425 The resulting BAM files were validated for integrity using `samtools` (version 1.17) (Danecek
426 et al., 2021), then merged, converted to FASTQ, and compressed using `samtools bgzip`
427 (`HTSLib` version 1.17) (Bonfield et al., 2021). The `dx` tag in the BAM records was used to
428 differentiate between duplex and simplex reads, defined as follows:

- 429 • `dx:i:1` for duplex reads.
- 430 • `dx:i:0` for simplex reads without duplex offspring.
- 431 • `dx:i:-1` for simplex reads with duplex offspring.

432 In this way, `samtools view` (version 1.17) (Danecek et al., 2021) was used to extract the
433 `'dx:i:1'` tag for duplex reads. Simplex reads were extracted as both `'dx:i:0'` and `'dx:i:-1'`
434 to avoid potential loss of information. Since Dorado version 0.3.1 lacked integrated adapter
435 removal functionality, `Porechop` (version 0.2.4) (Wick and Volkening, 2018) was used to
436 trim adapters from PromethION data with default settings.

437 6.0.2 FLYE Assembler

438 The following options were used FLYE (version 2.9.2-b1795) (Kolmogorov et al., 2019):
439 `-nano-hq -read-error 0.03 -genome-size 0.7g -scaffold`. The `-read-error` pa-
440 rameter was chosen as recommended for Q20-quality reads.

441 **6.1 MaSuRCA assembler**

442 Unfiltered Illumina sequences and filtered ONT reads were used for the assembly. This
443 was run on MaSuRCA (version 3.2.9) (Zimin et al., 2013) with recommended parameters,
444 automatic k-mer size computation, and a config input of;

```
445 PE = pe 150 22
446 NANOPORE =
447 EXTEND\_JUMP\_READS=0
448 GRAPH\_KMER\_SIZE = auto
449 USE\_LINKING\_MATES = 0
450 GRID\_BATCH\_SIZE=500000000
451 LHE\_COVERAGE=25
452 LIMIT\_JUMP\_COVERAGE = 300
453 CA\_PARAMETERS = cgwErrorRate=0.15 ovlRefBlockSize=40000000
454 CLOSE\_GAPS=1
455 NUM\_THREADS = 32
456 JF\_SIZE = 13000000000
```

457 Note that for the paired end parameter, the standard deviation was set to approxi-
458 mately 15% of the mean following MaSuRCA recommendations ([https://github.com/](https://github.com/alekseyzimin/masurca#configuration)
459 [alekseyzimin/masurca#configuration](https://github.com/alekseyzimin/masurca#configuration)).

460 **6.1.1 Protein sequence prediction and functional annotation**

461 Functional annotations were integrated into gene models using the AGAT toolkit (version
462 1.4.0) (Dainat, 2022). Specifically, `agat_sp_manage_functional_annotation.pl` merged
463 annotation results with the BRAKER3-generated gene models, refining gene product
464 names and descriptions. Both BLAST+ and InterProScan annotations were filtered using
465 an E-value threshold of 1×10^{-6} . Predicted gene models were refined using AGAT scripts
466 in the following order:

- 467 1. Premature stop codon removal (`agat_sp_flag_premature_stop_codons.pl`)
- 468 2. ORF Size filtering (`agat_sp_filter_by_ORF_size.pl -test ">=" -size 50`)
- 469 3. Filtering Incomplete Coding Models (`agat_sp_filter_incomplete_gene_coding_models.pl`)
- 470 4. Gene length filtering (`agat_sp_filter_gene_by_length.pl -test ">=" -size`
471 `50`)

472 **6.2 Appendix - Figures**

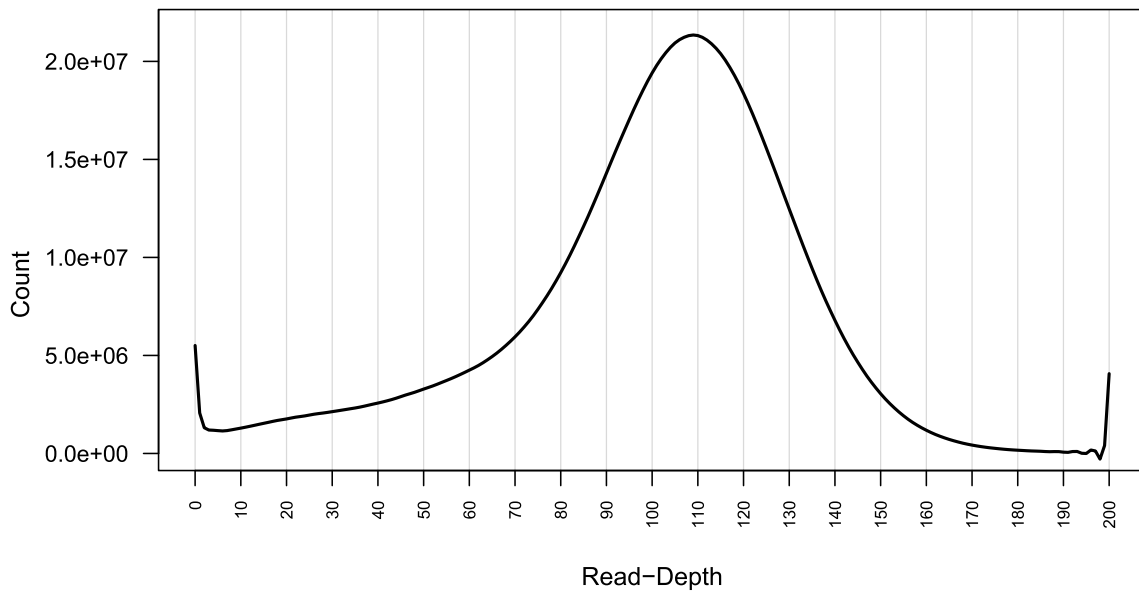


Figure S1: Graphical output of purge haplotigs, demonstrating 110 at the $2n$ peak. Read-depth histogram generated by Purge Haplotigs using Illumina paired-end reads mapped to the polished draft assembly. Illumina reads were used for coverage estimation due to their higher mapping accuracy relative to long reads. The histogram shows a single dominant peak at approximately $110\times$, corresponding to diploid coverage. The lack of a secondary haploid peak suggests minimal haplotig redundancy and supports the use of a haploid-collapsed assembly.

473 **6.3 Appendix - Tables**

Table S1: Assembly statistics and BUSCO completeness across major genome assembly stages, from raw draft (MaSuRCA) to final chromosome-level scaffolds (Haku Raukura - *i.e., the longest 24 scaffolds).

Metric	First draft	Pilon 1	Pilon 2	Pilon 3	Purged haplotigs	Mapped to Hi-C	Haku Raukura*
span (bp)	648,299,410	648,080,790	648,016,141	647,977,682	644,971,134	644,985,034	642,715,660
N (%)	-	-	-	-	-	-	-
GC (%)	40.81	40.81	40.81	40.81	40.81	40.81	40.81
AT (%)	59.19	59.19	59.19	59.19	59.19	59.19	59.19
scaffold count	390	390	390	390	273	146	24
longest scaffold (bp)	28,950,014	28,944,476	28,943,996	28,941,487	28,941,487	34,161,391	34,161,391
scaffold N50 length (bp)	14,193,346	14,192,664	14,192,744	14,192,743	14,192,743	28,040,511	28,040,511
scaffold N50 count	17	17	17	17	17	11	11
scaffold N90 length (bp)	2,230,765	2,228,891	2,228,875	2,228,872	2,272,965	22,307,305	22,307,305
scaffold N90 count	60	60	60	60	59	2	21
Complete BUSCOs (C)	98.70	98.70	98.70	98.70	98.60	98.70	98.60
Complete and single-copy BUSCOs (S)	97.90	97.90	97.90	97.90	97.90	98.00	98.00
Complete and duplicated BUSCOs (D)	0.80	0.80	0.80	0.80	0.70	0.70	0.70
Fragmented BUSCOs (F)	0.20	0.20	0.20	0.20	0.20	0.20	0.20
Missing BUSCOs (M)	1.10	1.10	1.10	1.10	1.20	1.10	1.20

Statistics

BUSCO (%)

References

- 474
- 475 Salah M. Aljanabi and Iciar Martinez. Universal and rapid salt-extraction of high quality
476 genomic DNA for PCR-based techniques. *Nucleic Acids Research*, 25(22), 1997. ISSN
477 03051048. doi: 10.1093/nar/25.22.4692.
- 478 Rolf Apweiler, Amos Bairoch, Cathy H. Wu, Winona C. Barker, Brigitte Boeckmann,
479 Serenella Ferro, Elisabeth Gasteiger, Hongzhan Huang, Rodrigo Lopez, Michelle Ma-
480 grane, Maria J. Martin, Darren A. Natale, Claire O'Donovan, Nicole Redaschi, and Lai
481 Su L. Yeh. UniProt: The universal protein knowledgebase. *Nucleic Acids Research*, 32,
482 2004. ISSN 03051048. doi: 10.1093/nar/gky092.
- 483 Alexander Astashyn, Eric S. Tvedte, Deacon Sweeney, Victor Sapojnikov, Nathan Bouk,
484 Victor Joukov, Eyal Mozes, Pooja K. Strobe, Pape M. Sylla, Lukas Wagner, Shelby L.
485 Bidwell, Karen Clark, Emily W. Davis, Brian Smith-White, Wratko Hlavina, Kim D.
486 Pruitt, Valerie A. Schneider, and Terence D. Murphy. Rapid and sensitive detection of
487 genome contamination at scale with FCS-GX. *bioRxiv*, 2023.
- 488 Weidong Bao, Kenji K. Kojima, and Oleksiy Kohany. Repbase Update, a database of
489 repetitive elements in eukaryotic genomes. *Mobile DNA*, 6(1), 2015. ISSN 17598753.
490 doi: 10.1186/s13100-015-0041-9.
- 491 Matthias Blum, Hsin Yu Chang, Sara Chuguransky, Tiago Grego, Swaathi Kandasamy,
492 Alex Mitchell, Gift Nuka, Typhaine Paysan-Lafosse, Matloob Qureshi, Shriya Raj,
493 Lorna Richardson, Gustavo A. Salazar, Lowri Williams, Peer Bork, Alan Bridge, Julian
494 Gough, Daniel H. Haft, Ivica Letunic, Aron Marchler-Bauer, Huaiyu Mi, Darren A.
495 Natale, Marco Necci, Christine A. Orengo, Arun P. Pandurangan, Catherine Rivoire,
496 Christian J.A. Sigrist, Ian Sillitoe, Narmada Thanki, Paul D. Thomas, Silvio C.E.
497 Tosatto, Cathy H. Wu, Alex Bateman, and Robert D. Finn. The InterPro protein
498 families and domains database: 20 years on. *Nucleic Acids Research*, 49(D1), 2021.
499 ISSN 13624962. doi: 10.1093/nar/gkaa977.
- 500 Anthony M. Bolger, Marc Lohse, and Bjoern Usadel. Trimmomatic: A flexible trimmer
501 for Illumina sequence data. *Bioinformatics*, 30(15), 2014. ISSN 14602059. doi: 10.1093/
502 bioinformatics/btu170.
- 503 James K. Bonfield, John Marshall, Petr Danecek, Heng Li, Valeriu Ohan, Andrew
504 Whitwham, and Thomas Keane. HTSlib: C library for reading/writing high-Throughput
505 sequencing data. *GigaScience*, 10(2), 2021. ISSN 2047217X. doi: 10.1093/gigascience/
506 giab007.
- 507 Christiam Camacho, George Coulouris, Vahram Avagyan, Ning Ma, Jason Papadopoulos,

508 Kevin Bealer, and Thomas L. Madden. BLAST+: Architecture and applications. *BMC*
509 *Bioinformatics*, 10, 2009. ISSN 14712105. doi: 10.1186/1471-2105-10-421.

510 Andrew Catanach, Mike Ruigrok, Deepa Bowatte, Marcus Davy, Roy Storey, Noémie
511 Valenza-Troubat, Elena López-Girona, Elena Hilario, Matthew J. Wylie, David Chagné,
512 and Maren Wellenreuther. The genome of New Zealand trevally (Carangidae: Pseu-
513 docaranx georgianus) uncovers a XY sex determination locus. *BMC Genomics*, 22(1),
514 2021. ISSN 14712164. doi: 10.1186/s12864-021-08102-2.

515 Jacques Dainat. AGAT: Another GFF Analysis Toolkit to handle annotations in any
516 GTF/GFF format. *Zenodo*, 4, 2022.

517 Petr Danecek, James K. Bonfield, Jennifer Liddle, John Marshall, Valeriu Ohan, Martin O.
518 Pollard, Andrew Whitwham, Thomas Keane, Shane A. McCarthy, and Robert M.
519 Davies. Twelve years of SAMtools and BCFtools. *GigaScience*, 10(2), 2021. ISSN
520 2047217X. doi: 10.1093/gigascience/giab008.

521 Wouter De Coster and Rosa Rademakers. NanoPack2: population-scale evaluation
522 of long-read sequencing data. *Bioinformatics*, 39(5), 2023. ISSN 13674811. doi:
523 10.1093/bioinformatics/btad311.

524 Philip Ewels, Måns Magnusson, Sverker Lundin, and Max Käller. MultiQC: Summarize
525 analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32
526 (19), 2016. ISSN 14602059. doi: 10.1093/bioinformatics/btw354.

527 Gregory G. Faust and Ira M. Hall. SAMBLASTER: Fast duplicate marking and structural
528 variant read extraction. In *Bioinformatics*, volume 30, 2014. doi: 10.1093/bioinformatics/
529 btu314.

530 Jullien M. Flynn, Robert Hubley, Clément Goubert, Jeb Rosen, Andrew G. Clark, Cédric
531 Feschotte, and Arian F. Smit. RepeatModeler2 for automated genomic discovery of
532 transposable element families. *Proceedings of the National Academy of Sciences of the*
533 *United States of America*, 117(17), 2020. ISSN 10916490. doi: 10.1073/pnas.1921046117.

534 Lars Gabriel, Katharina J. Hoff, Tomáš Brůna, Mark Borodovsky, and Mario Stanke.
535 TSEBRA: transcript selector for BRAKER. *BMC Bioinformatics*, 22(1), 2021. ISSN
536 14712105. doi: 10.1186/s12859-021-04482-0.

537 Torsten Günther and Carl Nettelblad. The presence and impact of reference bias on
538 population genomic studies of prehistoric human populations. *PLoS Genetics*, 15(7),
539 2019. ISSN 15537404. doi: 10.1371/journal.pgen.1008302.

540 Ross D. Houston, Tim P. Bean, Daniel J. Macqueen, Manu Kumar Gundappa, Ye Hwa
541 Jin, Tom L. Jenkins, Sarah Louise C. Selly, Samuel A.M. Martin, Jamie R. Stevens,

- 542 Eduarda M. Santos, Andrew Davie, and Diego Robledo. Harnessing genomics to fast-
543 track genetic improvement in aquaculture. *Nature Reviews Genetics*, 21(7), 2020. ISSN
544 14710064. doi: 10.1038/s41576-020-0227-y.
- 545 Wataru Iwasaki, Tsukasa Fukunaga, Ryota Isagozawa, Koichiro Yamada, Yasunobu Maeda,
546 Takashi P. Satoh, Tetsuya Sado, Kohji Mabuchi, Hirohiko Takeshima, Masaki Miya,
547 and Mutsumi Nishida. Mitofish and mitoannotator: A mitochondrial genome database
548 of fish with an accurate and automatic annotation pipeline. *Molecular Biology and*
549 *Evolution*, 30(11), 2013. ISSN 07374038. doi: 10.1093/molbev/mst141.
- 550 Philip Jones, David Binns, Hsin Yu Chang, Matthew Fraser, Weizhong Li, Craig McAnulla,
551 Hamish McWilliam, John Maslen, Alex Mitchell, Gift Nuka, Sebastien Pesseat, Antony F.
552 Quinn, Amaia Sangrador-Vegas, Maxim Scheremetjew, Siew Yit Yong, Rodrigo Lopez,
553 and Sarah Hunter. InterProScan 5: Genome-scale protein function classification.
554 *Bioinformatics*, 30(9), 2014. ISSN 14602059. doi: 10.1093/bioinformatics/btu031.
- 555 Aikaterini Katirtzoglou, Dimitris Tsaparis, Evangelos Kolios, Antonios Magoulas, Con-
556 stantinos C. Mylonas, Ioannis Fakriadis, Tereza Manousaki, and Costas S. Tsigenopou-
557 los. Population genomic analysis of the greater amberjack (*Seriola dumerili*) in the
558 Mediterranean and the Northeast Atlantic, based on SNPs, microsatellites, and mi-
559 tochondrial DNA sequences. *Frontiers in Fish Science*, 2, 2024. ISSN 28139097. doi:
560 10.3389/frish.2024.1356313.
- 561 Mikhail Kolmogorov, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner. Assembly of long, error-
562 prone reads using repeat graphs. *Nature Biotechnology*, 37(5), 2019. ISSN 15461696.
563 doi: 10.1038/s41587-019-0072-8.
- 564 Dmitry Kuznetsov, Fredrik Tegenfeldt, Mosè Manni, Mathieu Seppey, Matthew Berkeley,
565 Evgenia V. Kriventseva, and Evgeny M. Zdobnov. OrthoDB v11: annotation of orthologs
566 in the widest sampling of organismal diversity. *Nucleic Acids Research*, 51(1 D), 2023.
567 ISSN 13624962. doi: 10.1093/nar/gkac998.
- 568 Marie Lataretu, Sebastian Krautwurst, Adrian Viehweger, Christian Brandt, and Martin
569 Hölzer. Targeted decontamination of sequencing data with CLEAN. *bioRxiv*, 2023.
- 570 Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows-
571 Wheeler transform. *Bioinformatics*, 25(14), 2009. ISSN 13674803. doi: 10.1093/
572 bioinformatics/btp324.
- 573 Shuo Li, Kaiqiang Liu, Aijun Cui, Xiancai Hao, Bin Wang, Hong Yan Wang, Yan Jiang,
574 Qian Wang, Bo Feng, Yongjiang Xu, Changwei Shao, and Xuezhou Liu. A Chromosome-
575 Level Genome Assembly of Yellowtail Kingfish (*Seriola lalandi*). *Frontiers in Genetics*,
576 12, 2022. ISSN 16648021. doi: 10.3389/fgene.2021.825742.

- 577 G Marçais and C Kingsford. Jellyfish : A fast k-mer counter. *Tutorialis e Manuais*, (1),
578 2012.
- 579 Guillaume Marçais, Arthur L. Delcher, Adam M. Phillippy, Rachel Coston, Steven L.
580 Salzberg, and Aleksey Zimin. MUMmer4: A fast and versatile genome alignment system.
581 *PLoS Computational Biology*, 14(1), 2018. ISSN 15537358. doi: 10.1371/journal.pcbi.
582 1005944.
- 583 Tom Oosting, Elena Hilario, Maren Wellenreuther, and Peter A. Ritchie. DNA degradation
584 in fish: Practical solutions and guidelines to improve DNA preservation for genomic
585 research. *Ecology and Evolution*, 10(16), 2020. ISSN 20457758. doi: 10.1002/ece3.6558.
- 586 James M. Pflug, Valerie Renee Holmes, Crystal Burrus, J. Spencer Johnston, and David R.
587 Maddison. Measuring genome sizes using read-depth, k-mers, and flow cytometry:
588 Methodological comparisons in beetles (Coleoptera). *G3: Genes, Genomes, Genetics*,
589 10(9), 2020. ISSN 21601836. doi: 10.1534/g3.120.401028.
- 590 H. K.A. Premachandra, Fabiola Lafarga De La Cruz, Yutaka Takeuchi, Adam Miller,
591 Stewart Fielder, Wayne O'Connor, Celine H. Frère, Nguyen Hong Nguyen, Ido Bar,
592 and Wayne Knibb. Genomic DNA variation confirmed *Seriola lalandi* comprises three
593 different populations in the Pacific, but with recent divergence. *Scientific Reports*, 7(1),
594 2017. ISSN 20452322. doi: 10.1038/s41598-017-07419-x.
- 595 Michael J. Roach, Simon A. Schmidt, and Anthony R. Borneman. Purge Haplotigs: Allelic
596 contig reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics*, 19
597 (1), 2018. ISSN 14712105. doi: 10.1186/s12859-018-2485-7.
- 598 Yukuto Sato, Masaki Miya, Tsukasa Fukunaga, Tetsuya Sado, and Wataru Iwasaki.
599 MitoFish and mifish pipeline: A mitochondrial genome database of fish with an analysis
600 pipeline for environmental DNA metabarcoding. *Molecular Biology and Evolution*, 35
601 (6), 2018. ISSN 15371719. doi: 10.1093/molbev/msy074.
- 602 A F A Smit, R Hubley, and P Grenn. RepeatMasker Open-4.0, 2015.
- 603 Mario Stanke, Oliver Schöffmann, Burkhard Morgenstern, and Stephan Waack. Gene
604 prediction in eukaryotes with a generalized hidden Markov model that uses hints
605 from external sources. *BMC Bioinformatics*, 7, 2006. ISSN 14712105. doi: 10.1186/
606 1471-2105-7-62.
- 607 Jessica Storer, Robert Hubley, Jeb Rosen, Travis J. Wheeler, and Arian F. Smit. The Dfam
608 community resource of transposable element families, sequence models, and genome an-
609 notations. *Mobile DNA*, 12(1), 2021. ISSN 17598753. doi: 10.1186/s13100-020-00230-y.

610 J E Symonds, S P Walker, I van de Ven, A Marchant, G Irvine, S Pether, and Y Gublin.
611 Developing broodstock resources for farmed marine fish. In *Proceedings of the New*
612 *Zealand Society of Animal Production*, volume 72, 2012.

613 Hiroshi Takahashi, Taiki Kurogoushi, Ryo Shimoyama, and Hiroyuki Yoshikawa. First
614 report of natural hybridization between two yellowtails, *Seriola quinqueradiata* and
615 *S. lalandi*. *Ichthyological Research*, 68(1):139–144, 1 2021. ISSN 16163915. doi:
616 10.1007/s10228-020-00752-8.

617 Mun Hua Tan, Christopher M. Austin, Michael P. Hammer, Yin Peng Lee, Laurence J.
618 Croft, and Han Ming Gan. Finding Nemo: Hybrid assembly with Oxford Nanopore and
619 Illumina reads greatly improves the clownfish (*Amphiprion ocellaris*) genome assembly.
620 *GigaScience*, 7(3), 2018. ISSN 2047217X. doi: 10.1093/gigascience/gix137.

621 Ben Te Aika, Libby Liggins, Claire Rye, E. Owen Perkins, Jun Huh, Rudiger Brauning,
622 Tracey Godfery, and Michael A. Black. Aotearoa genomic data repository: An āhuru
623 mōwai for taonga species sequencing data. *Molecular Ecology Resources*, 2023. ISSN
624 17550998. doi: 10.1111/1755-0998.13866.

625 Doko Miles J. Thorburn, Kostas Sagonas, Mahesh Binzer-Panchal, Frederic J.J. Chain,
626 Philine G.D. Feulner, Erich Bornberg-Bauer, Thorsten B.H. Reusch, Irene E. Samonte-
627 Padilla, Manfred Milinski, Tobias L. Lenz, and Christophe Eizaguirre. Origin matters:
628 Using a local reference genome improves measures in population genomics. *Molecular*
629 *Ecology Resources*, 23(7), 2023. ISSN 17550998. doi: 10.1111/1755-0998.13838.

630 Gregory W. Vulture, Fritz J. Sedlazeck, Maria Nattestad, Charles J. Underwood, Han
631 Fang, James Gurtowski, and Michael C. Schatz. GenomeScope: Fast reference-free
632 genome profiling from short reads. In *Bioinformatics*, volume 33, 2017. doi: 10.1093/
633 bioinformatics/btx153.

634 Bruce J. Walker, Thomas Abeel, Terrance Shea, Margaret Priest, Amr Abouelliel, Sharadha
635 Sakthikumar, Christina A. Cuomo, Qiandong Zeng, Jennifer Wortman, Sarah K. Young,
636 and Ashlee M. Earl. Pilon: An integrated tool for comprehensive microbial variant
637 detection and genome assembly improvement. *PLoS ONE*, 9(11), 2014. ISSN 19326203.
638 doi: 10.1371/journal.pone.0112963.

639 Ryan Wick and Jeremy Volkening. Porechop: adapter trimmer for Oxford Nanopore reads.
640 *Github*, 2018.

641 Derrick E. Wood and Steven L. Salzberg. Kraken: Ultrafast metagenomic sequence
642 classification using exact alignments. *Genome Biology*, 15(3), 2014. ISSN 1474760X.
643 doi: 10.1186/gb-2014-15-3-r46.

- 644 José M. Yáñez, Agustín Barría, María E. López, Thomas Moen, Baltasar F. Garcia,
645 Grazyella M. Yoshida, and Peng Xu. Genome-wide association and genomic selection
646 in aquaculture, 2023. ISSN 17535131.
- 647 Dian Chang Zhang, Liang Guo, Hua Yang Guo, Ke Cheng Zhu, Shang Qi Li, Yan Zhang,
648 Nan Zhang, Bao Suo Liu, Shi Gui Jiang, and Jiong Tang Li. Chromosome-level genome
649 assembly of golden pompano (*Trachinotus ovatus*) in the family Carangidae. *Scientific*
650 *Data*, 6(1), 2019. ISSN 20524463. doi: 10.1038/s41597-019-0238-8.
- 651 Chenxi Zhou, Shane A. McCarthy, and Richard Durbin. YaHS: yet another Hi-C scaffolding
652 tool. *Bioinformatics*, 39(1), 2023. ISSN 13674811. doi: 10.1093/bioinformatics/btac808.
- 653 Tao Zhu, Yukuto Sato, Tetsuya Sado, Masaki Miya, and Wataru Iwasaki. MitoFish,
654 MitoAnnotator, and MiFish Pipeline: Updates in 10 Years. *Molecular Biology and*
655 *Evolution*, 40(3), 2023. ISSN 15371719. doi: 10.1093/molbev/msad035.
- 656 Aleksey V. Zimin, Guillaume Marçais, Daniela Puiu, Michael Roberts, Steven L. Salzberg,
657 and James A. Yorke. The MaSuRCA genome assembler. *Bioinformatics*, 29(21), 2013.
658 ISSN 13674803. doi: 10.1093/bioinformatics/btt476.