

# ERGA Assembly Report

v24.10.15

Tags: ATLASEa[INVALID TAG]

TxID	303921
ToLID	<b>fAthPrel</b>
Species	Atherina presbyter
Class	Actinopteri
Order	Atheriniformes

Genome Traits	Expected	Observed
Haploid size (bp)	986,362,630	1,040,595,279
Haploid Number	18 (source: ancestor)	24
Ploidy	2 (source: ancestor)	2
Sample Sex	Unknown	Unknown

## EBP metrics summary and curation notes

Obtained EBP quality metric for collapsed: 7.7.Q44

The following metrics were automatically flagged as below EBP recommended standards or different from expected:

- . Observed Haploid Number is different from Expected
- . Kmer completeness value is less than 90 for collapsed

### Curator notes

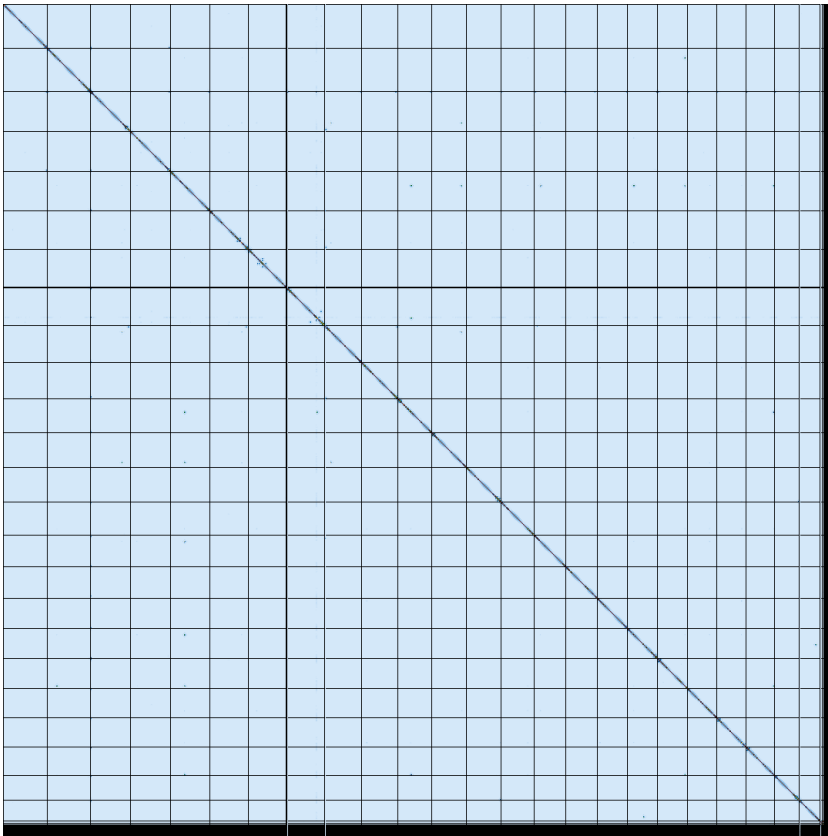
. Interventions/Gb: 19  
. Contamination notes: ""  
. Other observations: "The assembly of Atherina presbyter (fAthPrel.1) is based on 80X of PacBio data and 146X of Arima Hi-C data generated as part of the ATLASEa programme (<https://www.atlasea.fr>). The assembly process included the following steps: initial PacBio assembly generation with Hifiasm, removal of contaminant sequences using Context, removal of haplotypic duplications using purge\_dups, and Hi-C-based scaffolding with YaHS. In total, 1 contig of 0.013 Mb was identified as contaminant (bacterial). Additionally, 164 regions totaling 9 Mb (with the largest being 0.32 Mb) were identified as haplotypic duplications and removed. The mitochondrial genome was assembled using OATK. Finally, the primary assembly was analyzed and manually improved using Pretext. During manual curation, 25 haplotypic regions and 231 contaminant sequences were removed, totaling 15.2 Mb and 23.24 Mb respectively (with the largest being 1.59 Mb and 0.1 Mb. Chromosome-scale scaffolds confirmed by Hi-C data were named in order of size. "

# Quality metrics table

Metrics	Pre-curation collapsed	Curated collapsed
Total bp	1,069,086,011	1,040,595,279
GC %	39.2	39.34
Gaps/Gbp	261.91	203.73
Total gap bp	28,000	21,700
Scaffolds	394	127
Scaffold N50	43,223,242	43,223,242
Scaffold L50	11	11
Scaffold L90	22	22
Contigs	674	339
Contig N50	23,738,000	25,751,300
Contig L50	17	17
Contig L90	84	69
QV	35.2034	44.3653
Kmer compl.	80.4506	80.1427
BUSCO sing.	96.0%	97.0%
BUSCO dupl.	1.6%	0.6%
BUSCO frag.	0.6%	0.6%
BUSCO miss.	1.7%	1.8%

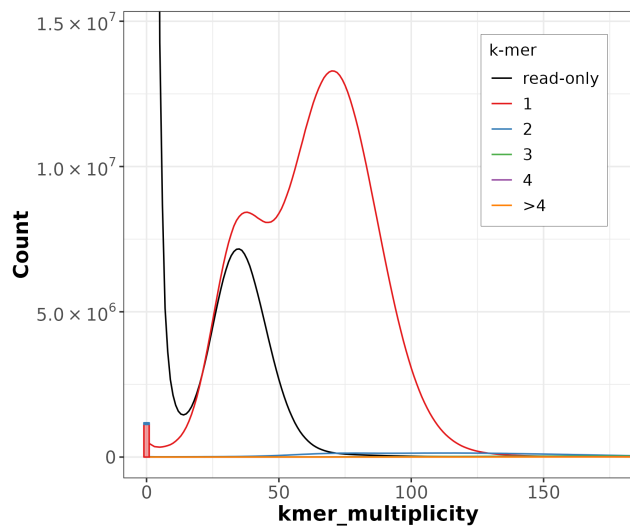
BUSCO: 5.8.2 (euk\_genome\_met, metaeuk) / Lineage: actinopterygii\_odb12 (genomes:75, BUSCOs:7207)

# HiC contact map of curated assembly

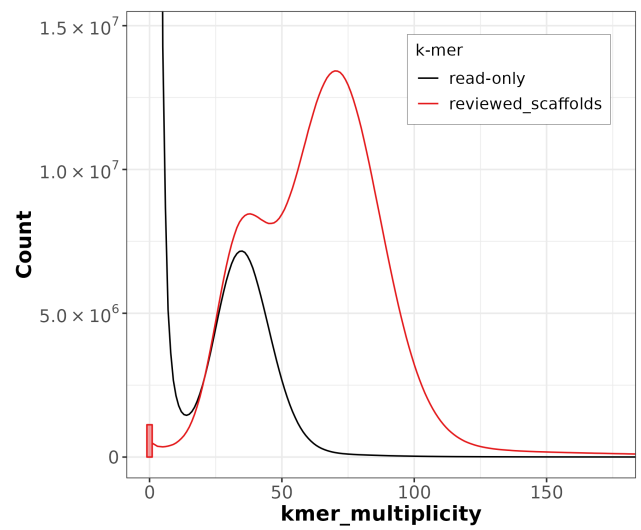


collapsed [\[LINK\]](#)

# K-mer spectra of curated assembly

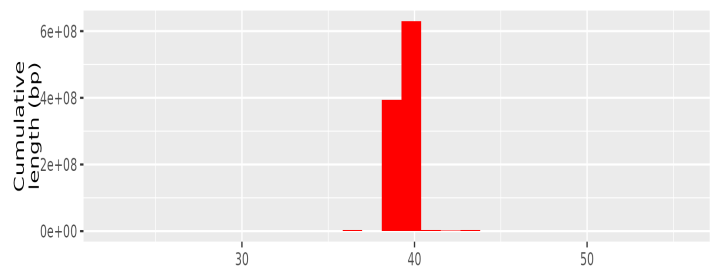


Distribution of k-mer counts per copy numbers found in asm

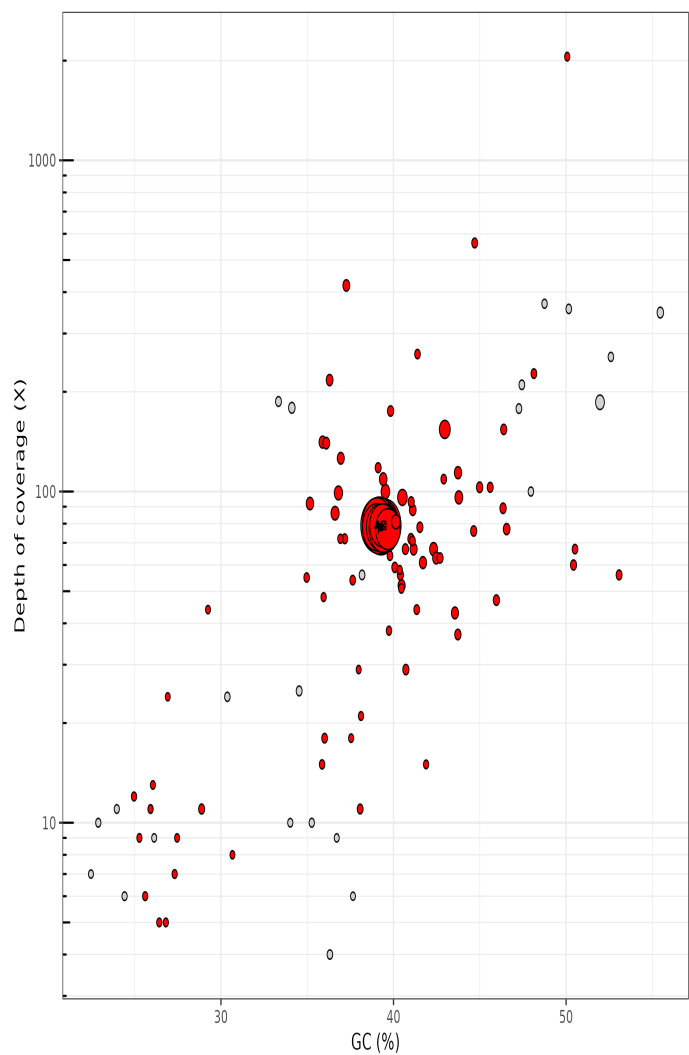


Distribution of k-mer counts coloured by their presence in reads/assemblies

# Post-curation contamination screening



TAPAs summary Graph



**collapsed.** Bubble plot circles are scaled by sequence length, positioned by coverage and GC proportion, and coloured by taxonomy. Histograms show total assembly length distribution on each axis.

## Data profile

Data	PACBIO Hifi	Arima
Coverage	80	146

## Assembly pipeline

- **Hifiasm**
  - |\_ *ver*: 0.19.5-r593
  - |\_ *key param*: NA
- **purge\_dups**
  - |\_ *ver*: 1.2.5
  - |\_ *key param*: NA
- **YaHS**
  - |\_ *ver*: 1.2
  - |\_ *key param*: NA

## Curation pipeline

- **PretextMap**
  - |\_ *ver*: 0.1.9
  - |\_ *key param*: NA
- **PretextView**
  - |\_ *ver*: 0.2.5
  - |\_ *key param*: NA

Submitter: Adama Ndar  
Affiliation: Genoscope

Date and time: 2025-05-28 14:53:34 CEST