

ERGA Assembly Report

v24.10.15

Tags: ATLASEa[INVALID TAG]

TxID	1365585
ToLID	fThoEph1
Species	Thorogobius ehippiatus
Class	Actinopteri
Order	Gobiiformes

Genome Traits	Expected	Observed
Haploid size (bp)	929,764,645	1,018,347,539
Haploid Number	24 (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.Q47

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

- . Kmer completeness value is less than 90 for collapsed

Curator notes

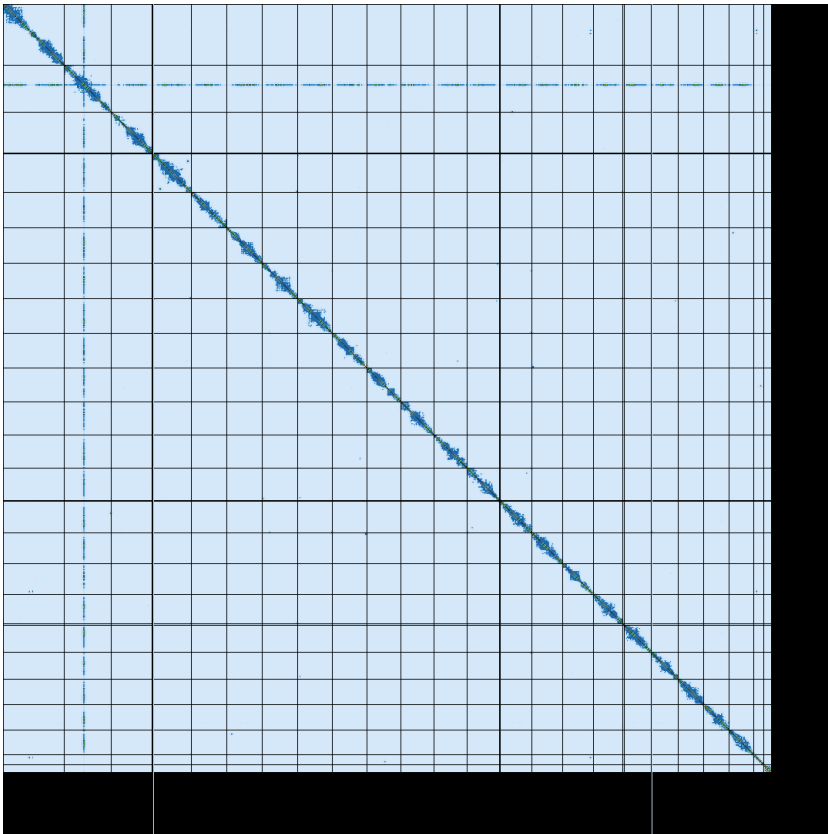
- . Interventions/Gb: 28
- . Contamination notes: ""
- . Other observations: "The assembly of Thorogobius ehippiat (fThoEph1.1) is based on 55X of PacBio data and Arima Hi-C data generated as part of the ATLASEa programme (<https://www.atlasea.fr>). The assembly process included the following steps: initial PacBio assembly generation with Hifiiasm, removal of contaminant sequences using Context, removal of haplotypic duplications using purge_dups, and Hi-C-based scaffolding with YaHS. In total, 6 contigs were identified as contaminants (bacterial, archaeal, or viral), totaling 0.40 Mb (with the largest being 0.20 Mb). Additionally, 270 regions totaling 19 Mb were identified as haplotypic duplications and removed (with the largest being 0.47 Mb). The mitochondrial genome was assembled using OATK. Finally, the primary assembly was analyzed and manually improved using Pretext. During manual curation, 2 haplotypic regions and 0 contaminant sequences were removed, totaling 1 Mb (with the largest being 0.78 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	1,019,431,561	1,018,347,539
GC %	41.22	41.22
Gaps/Gbp	167.74	182.65
Total gap bp	17,100	20,700
Scaffolds	1,098	1,051
Scaffold N50	40,862,349	40,862,349
Scaffold L50	11	11
Scaffold L90	23	23
Contigs	1,269	1,237
Contig N50	22,007,000	22,779,384
Contig L50	18	17
Contig L90	133	131
QV	47.2675	47.3055
Kmer compl.	88.5266	88.5115
BUSCO sing.	94.0%	94.0%
BUSCO dupl.	0.5%	0.5%
BUSCO frag.	1.9%	1.9%
BUSCO miss.	3.6%	3.5%

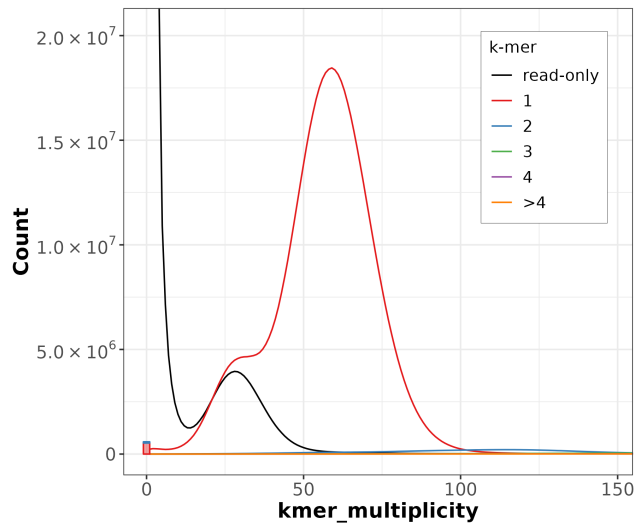
BUSCO: 5.8.2 (euk_genome_met, metaeuk) / Lineage: actinopterygii_odb12 (genomes:75, BUSCOs:7207)

HiC contact map of curated assembly

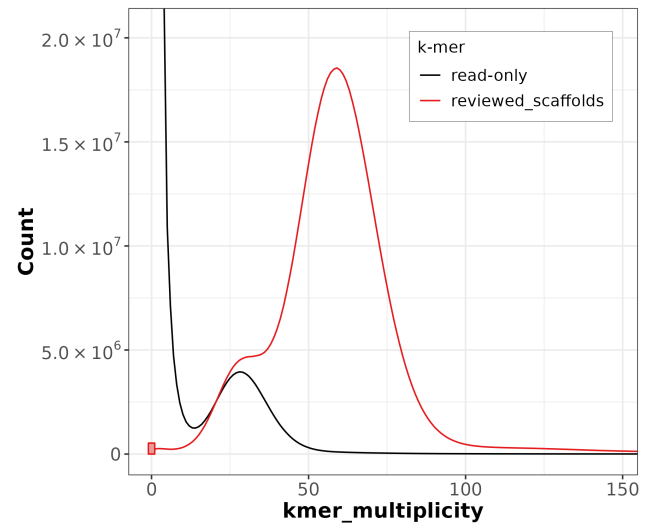


collapsed [\[LINK\]](#)

K-mer spectra of curated assembly

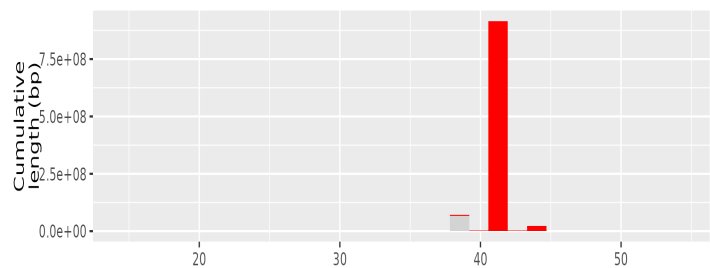


Distribution of k-mer counts per copy numbers found in asm

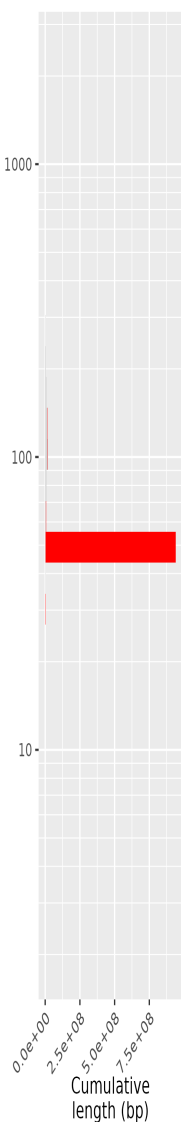
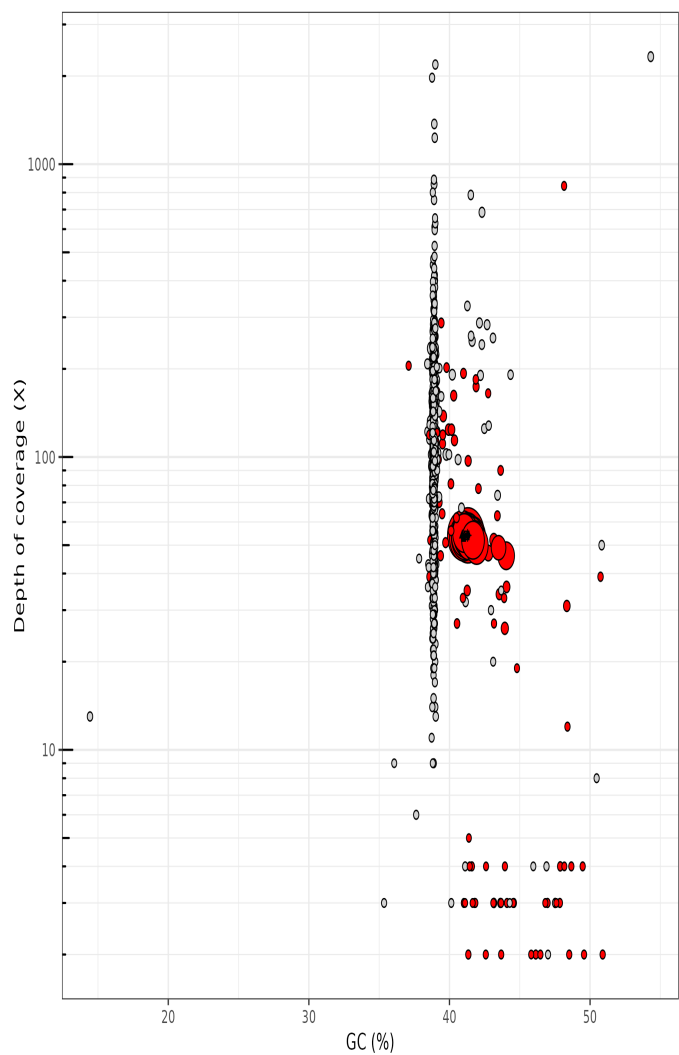


Distribution of k-mer counts coloured by their presence in reads/assemblies

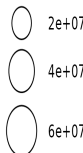
Post-curation contamination screening



TAPAs summary Graph



Length (bp)



Longest sequences (bp)

- fThoEph1_1 - 74877362 (Eukaryota)
- ▲ fThoEph1_2 - 57402124 (Eukaryota)
- fThoEph1_3 - 49882766 (Eukaryota)
- + fThoEph1_4 - 46514158 (Eukaryota)
- ▣ fThoEph1_5 - 43608956 (Eukaryota)

superkingdom

- Eukaryota
- N/A

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	55	141

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: Arnaud Couloux

Affiliation: Genoscope

Date and time: 2025-05-15 12:36:16 CEST