

# ERGA Assembly Report

v24.10.15

Tags: ATLASEa[INVALID TAG]

TxID	2951299
ToLID	<b>wpSigMath1</b>
Species	Sigalion mathildae
Class	Polychaeta
Order	Phyllodocida

Genome Traits	Expected	Observed
Haploid size (bp)	1,014,393,677	1,076,242,907
Haploid Number	10 (source: ancestor)	12
Ploidy	2 (source: ancestor)	2
Sample Sex	Unknown	Unknown

## EBP metrics summary and curation notes

Obtained EBP quality metric for collapsed: 6.7.Q45

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

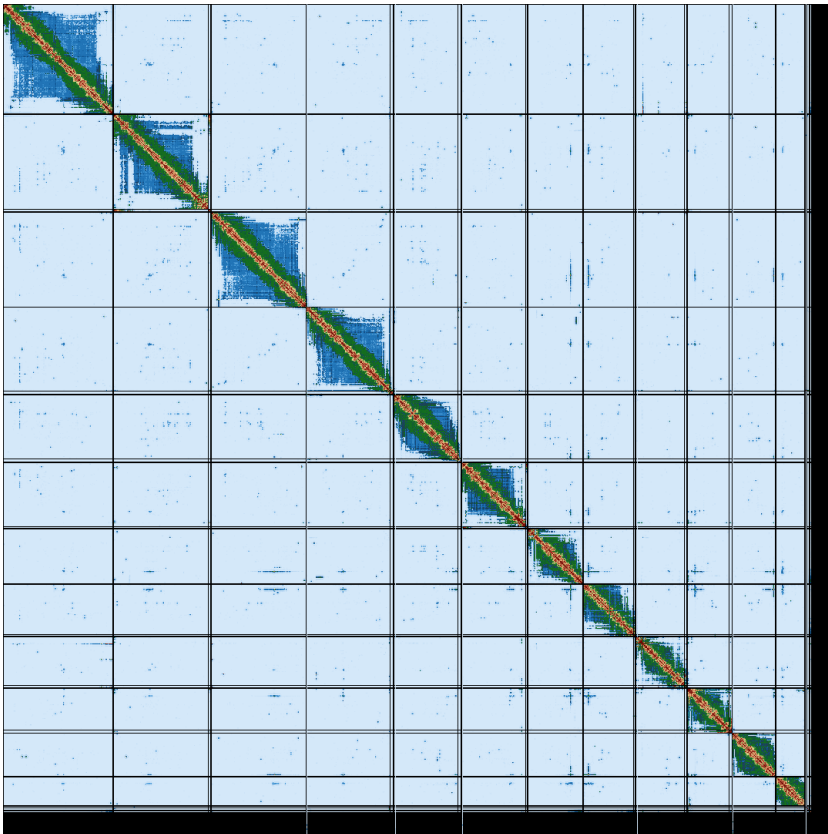
. Interventions/Gb: 111  
. Contamination notes: ""  
. Other observations: "The assembly of Sigalion mathildae (wpSigMath1) is based on 45X PacBio data and 144X Arima Hi-C data generated as part of the ATLASEa programme (<https://www.atlasea.fr>). The assembly process included the following steps: initial 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, 301 contigs were identified as contaminants (bacterial, archaeal, or viral), totaling 9.97 Mb (with the largest being 0.271 Mb). Additionally, 1003 regions totaling 95.292 Mb (with the largest being 2.801 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, 19 haplotypic regions and 3 contaminant sequences were removed, totaling 18.4Mb and 0.112Mb, respectively (with the largest being 1.9Mb and 0.055Mb). 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	1,094,840,021	1,076,242,907
GC %	39.49	39.49
Gaps/Gbp	825.69	806.51
Total gap bp	100,400	102,900
Scaffolds	478	416
Scaffold N50	82,238,604	83,475,863
Scaffold L50	5	5
Scaffold L90	12	11
Contigs	1,333	1,284
Contig N50	7,076,407	7,124,853
Contig L50	43	42
Contig L90	298	285
QV	45.0825	45.0899
Kmer compl.	72.0244	71.1385
BUSCO sing.	97.1%	97.3%
BUSCO dupl.	1.7%	1.6%
BUSCO frag.	0.3%	0.3%
BUSCO miss.	0.9%	0.8%

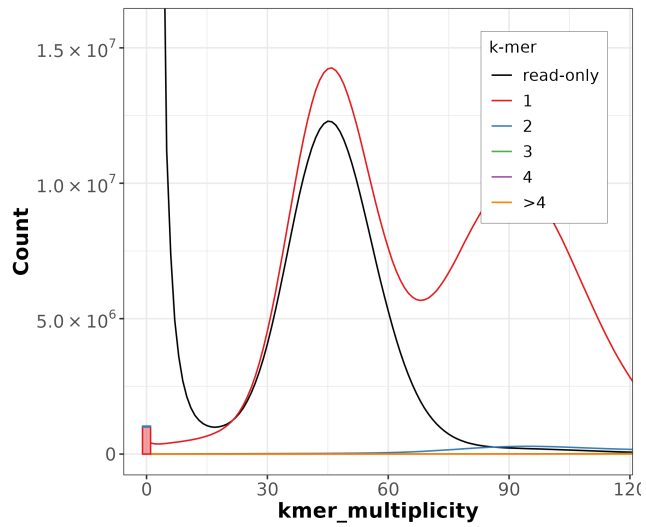
BUSCO: 6.0.0 (euk\_genome\_min, miniprot) / Lineage: lophotrochozoa\_odb12 (genomes:75, BUSCOs:1252)

# HiC contact map of curated assembly

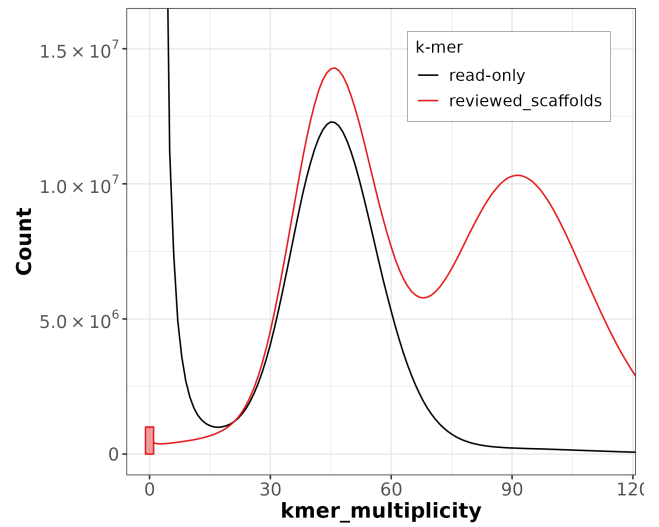


collapsed [\[LINK\]](#)

# K-mer spectra of curated assembly

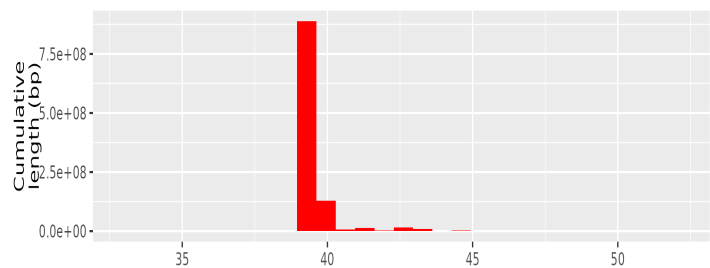


Distribution of k-mer counts per copy numbers found in asm

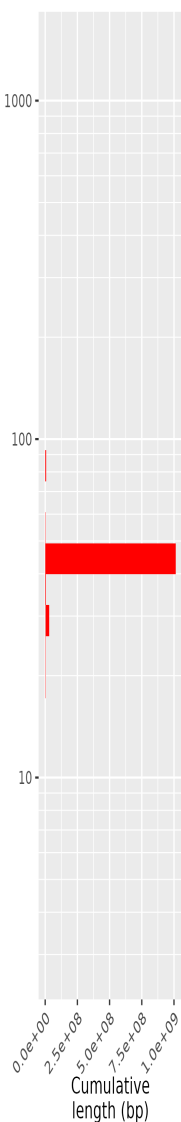
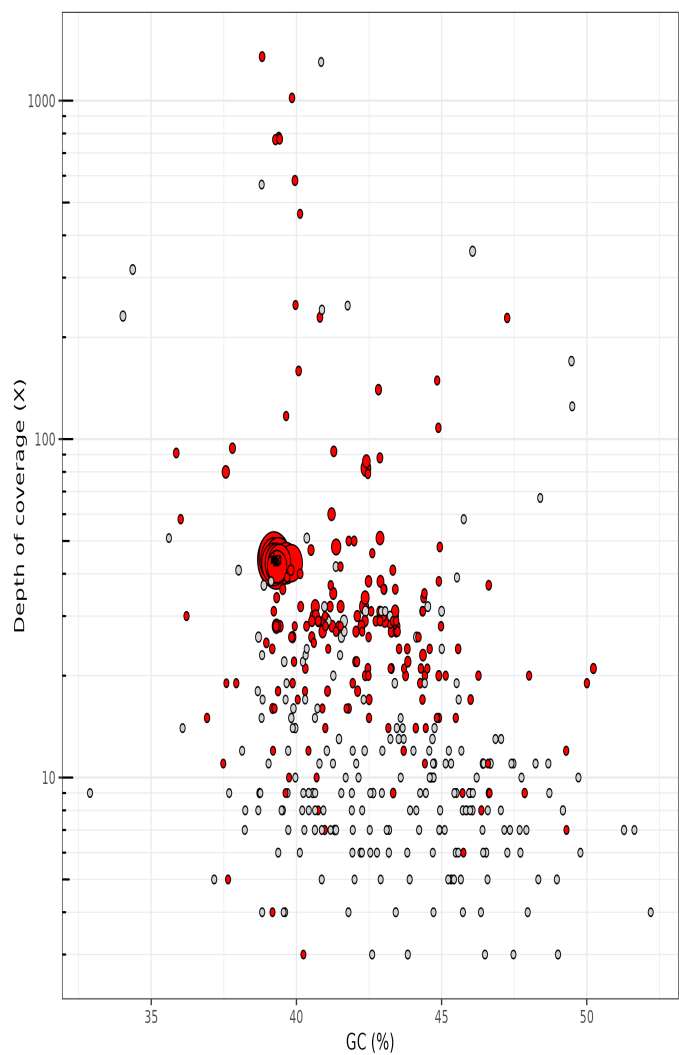


Distribution of k-mer counts coloured by their presence in reads/assemblies

# Post-curation contamination screening



TAPAs summary Graph



- Longest sequences (bp)
- wpSigMath1\_1 - 141950881 (Eukaryota)
  - ▲ wpSigMath1\_2 - 123039242 (Eukaryota)
  - wpSigMath1\_3 - 122805021 (Eukaryota)
  - + wpSigMath1\_4 - 108236363 (Eukaryota)
  - ▣ wpSigMath1\_6 - 83475863 (Eukaryota)

- Length (bp)
- 5.0e+07
  - 1.0e+08

- 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	Long reads	Arima
Coverage	45	144

# 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: Emilie Teodori

Affiliation: Genoscope

Date and time: 2025-11-09 07:21:26 CET