

ERGA Assembly Report

v24.10.15

Tags: ERGA-BGE

TxID	2712997
ToLID	icDenFora10
Species	Dendarus foraminosus
Class	Insecta
Order	Coleoptera

Genome Traits	Expected	Observed
Haploid size (bp)	557,593,036	593,991,149
Haploid Number	10 (source: ancestor)	11
Ploidy	2 (source: ancestor)	2
Sample Sex	Unknown	Unknown

EBP metrics summary and curation notes

Obtained EBP quality metric for collapsed: 7.7.Q62

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
- . Not 90% of assembly in chromosomes for collapsed

Curator notes

- . Interventions/Gb: 15
- . Contamination notes: ""
- . Other observations: "The assembly of Dendarus foraminosus (icDenFora10.1) is based on 58X PacBio data and OmniC Hi-C data generated as part of the European Reference Genome Atlas (ERGA, <https://www.erga-biodiversity.eu/>) via the Biodiversity Genomics Europe project (BGE, <https://biodiversitygenomics.eu/>). 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, 13 contigs were identified as contaminants (bacterial), totaling 840 Kb (with the largest being 358 Kb). Additionally, 255 regions totaling 30.5 Mb (with the largest being 1.6 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, no regions were tagged as allelic duplications or contaminants ; Scaffold_9 and Scaffold_11 with low coverage, and homology with Tenebrio molitor (X,Y) chromosomes, were renamed as X and Y respectively. Telomeric repeat pattern found is TCGGG. Chromosome-scale scaffolds

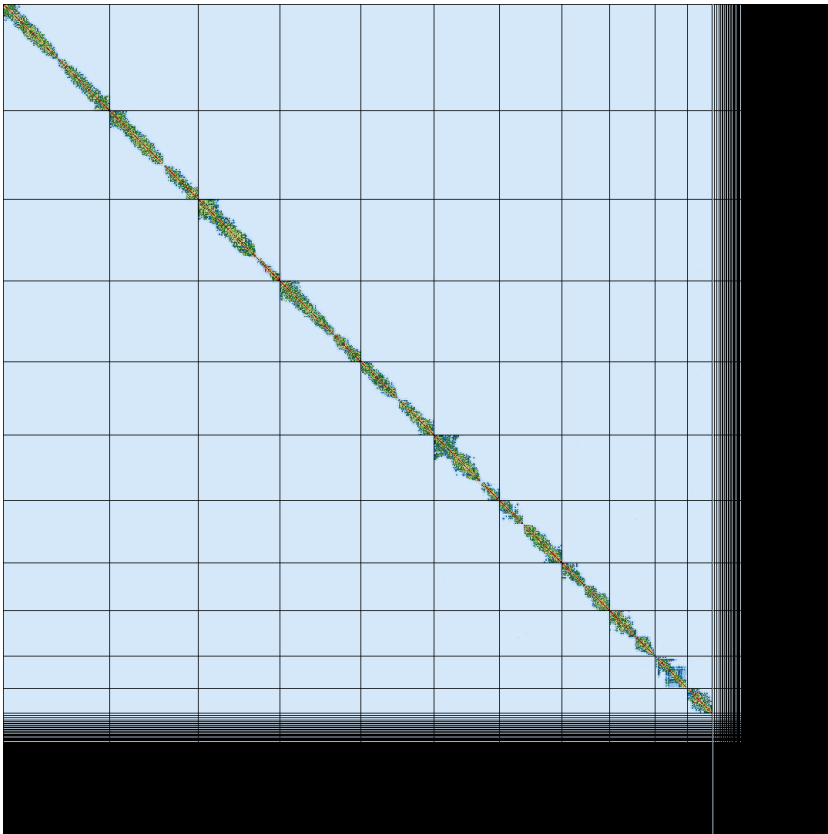
confirmed by Hi-C data were named in order of size. "

Quality metrics table

Metrics	Pre-curation collapsed	Curated collapsed
Total bp	593,969,168	593,991,149
GC %	35.92	35.92
Gaps/Gbp	16.84	25.25
Total gap bp	1,000	2,000
Scaffolds	419	417
Scaffold N50	41,134,000	51,873,200
Scaffold L50	6	5
Scaffold L90	52	49
Contigs	429	432
Contig N50	24,404,000	24,404,000
Contig L50	9	9
Contig L90	61	63
QV	62.2928	62.293
Kmer compl.	66.781	66.7811
BUSCO sing.	98.0%	98.0%
BUSCO dupl.	1.6%	1.6%
BUSCO frag.	0.4%	0.4%
BUSCO miss.	0.0%	0.0%

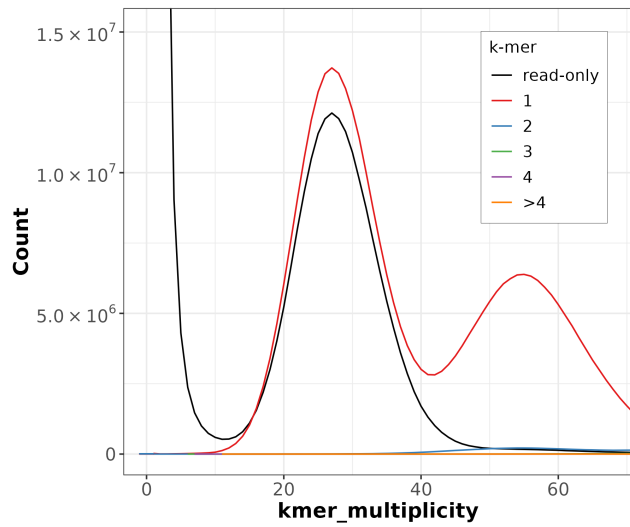
BUSCO: 5.4.3 (euk_genome_met, metaeuk) / Lineage: eukaryota_odb10 (genomes:70, BUSCOs:255)

HiC contact map of curated assembly

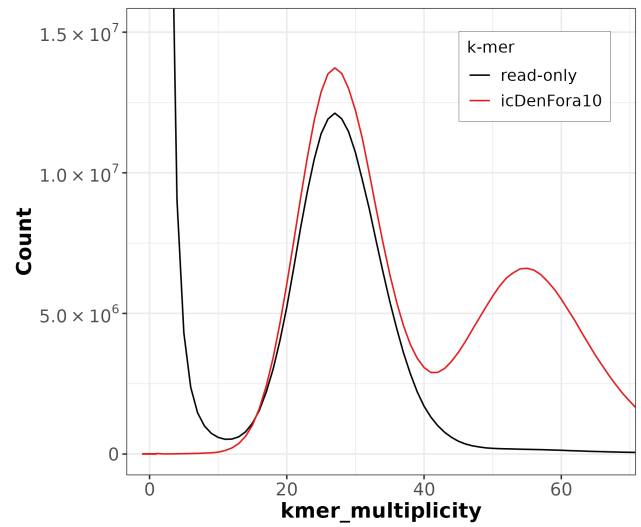


collapsed [\[LINK\]](#)

K-mer spectra of curated assembly

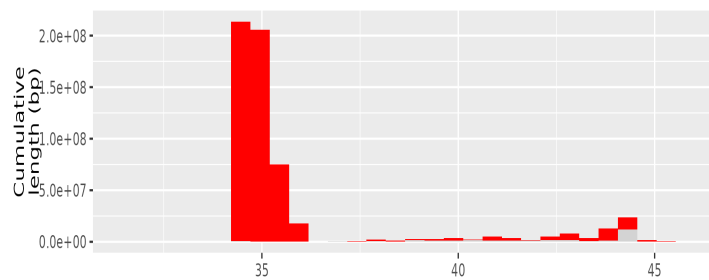


Distribution of k-mer counts per copy numbers found in asm



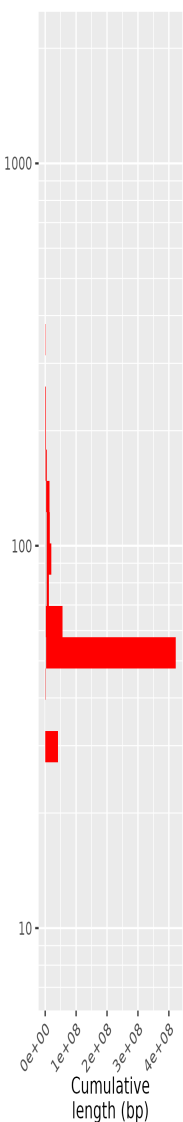
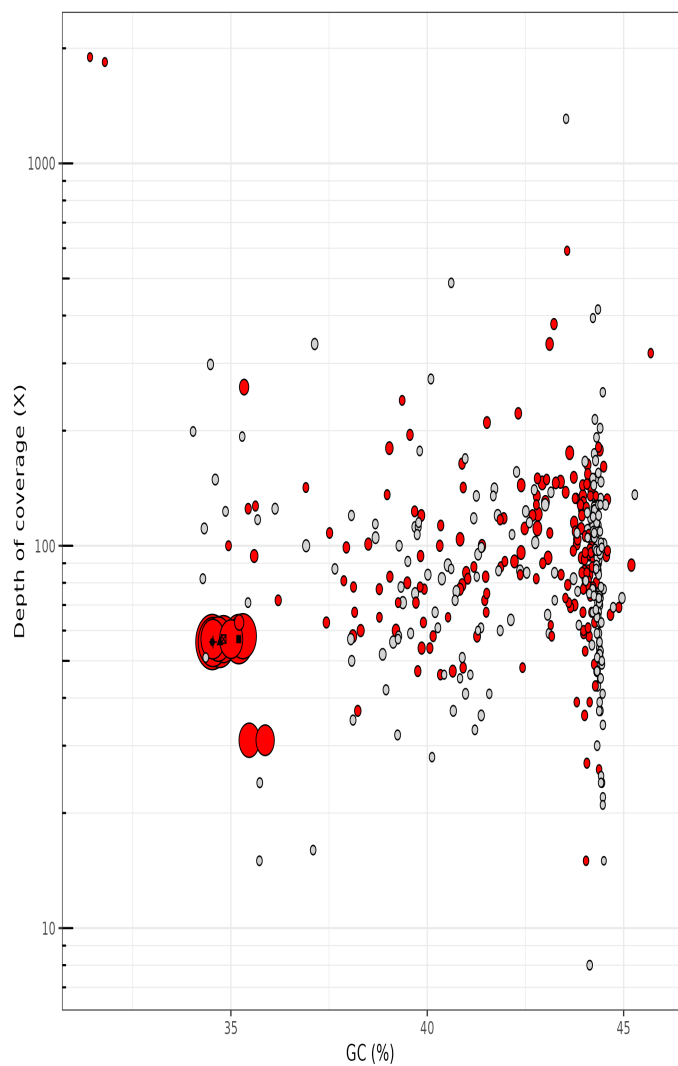
Distribution of k-mer counts coloured by their presence in reads/assemblies

Post-curation contamination screening

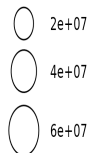


TAPAs summary Graph

(2 0X contigs have been hidden)



Length (bp)



Longest sequences (bp)

- icDenFora10_1 - 76254026 (Eukaryota)
- icDenFora10_2 - 63295100 (Eukaryota)
- icDenFora10_3 - 58385200 (Eukaryota)
- icDenFora10_4 - 57732603 (Eukaryota)
- icDenFora10_5 - 51873200 (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	Omnic
Coverage	58	82

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: S.Duprat

Affiliation: Genoscope

Date and time: 2025-01-24 03:08:38 CET