



# BRAKER2: A pipeline integrating data on genomic, RNA and protein sequences into inference of plant and animal genome annotation

Katharina J. Hoff<sup>1\*</sup>, Tomáš Brůna<sup>2\*</sup>, Alexandre Lomsadze<sup>2\*</sup>, Mario Stanke<sup>1</sup> and Mark Borodovsky<sup>2</sup>

1) Institute for Mathematics and Computer Science, University of Greifswald, Greifswald, GERMANY

2) Joint Georgia Tech and Emory Wallace H Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, U.S.A.

Contact: katharina.hoff@uni-greifswald.de, \*) Authors contributed equally

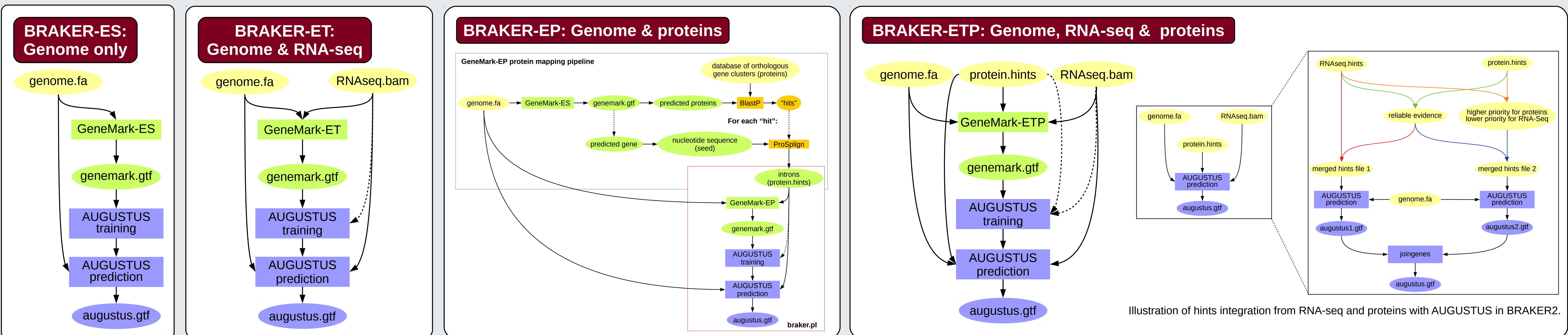


## Abstract

The rapidly growing number of sequenced eukaryotic genomes requires fully automated methods for accurate gene structure annotation. With this goal in mind, we had developed BRAKER1 [1], a combination of self-training GeneMark-ET [2] and AUGUSTUS [3], that uses genomic and RNA-seq data to automatically generate full gene structure annotations in novel genomes (including alternative isoforms).

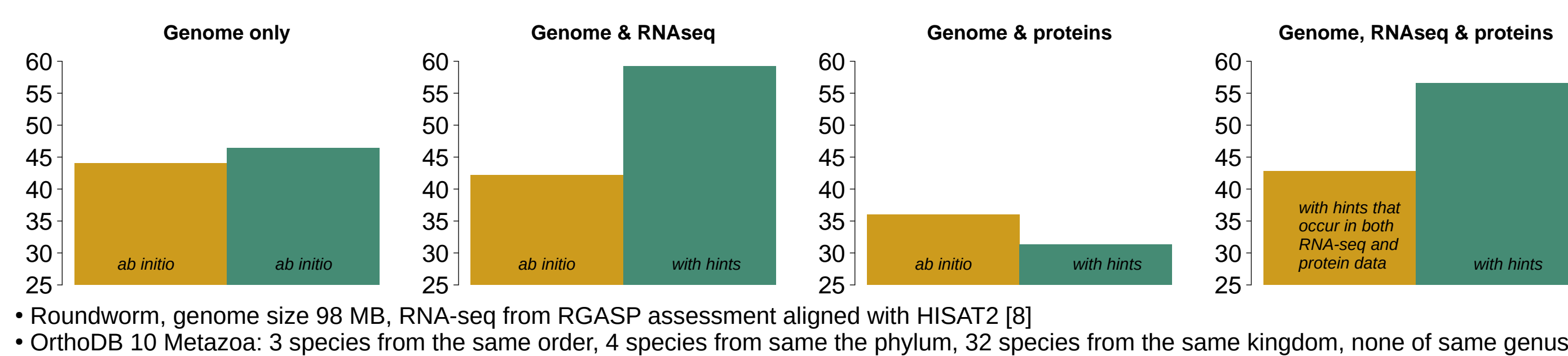
BRAKER2 [4] is an extension of BRAKER1. The new tool supports diverse annotation modes: i/ use of genome sequence only (ES), ii/ use of genome and RNA-seq data (ET), iii/ use of genome and proteins of possibly distant evolutionary origin (EP), iv/ use of genome, RNA-seq data and proteins. We have assessed gene prediction accuracy of BRAKER2 for two model organisms *Caenorhabditis elegans* and *Drosophila melanogaster*. In addition, we applied BRAKER2 to ten non-model organisms, using VARUS [5] as RNA-seq sampling tool and OrthoDB [6] as protein data resource.

BRAKER2 is available for download at <http://github.com/Gaius-Augustus/BRAKER>, the GeneMark-EP protein mapping pipeline is available for download at [http://exon.gatech.edu/GeneMark/Braker/protein\\_mapping\\_pipeline.tar.gz](http://exon.gatech.edu/GeneMark/Braker/protein_mapping_pipeline.tar.gz).

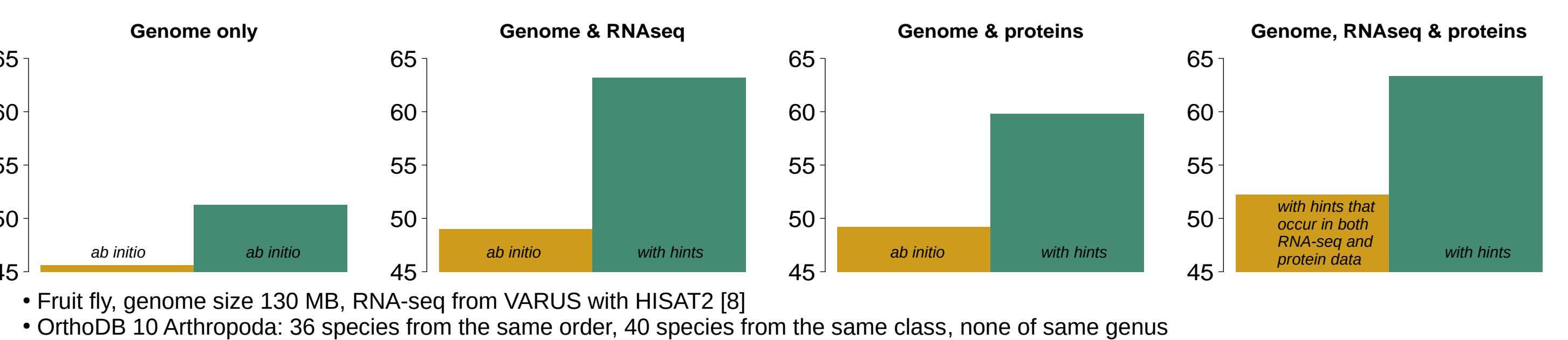


## Results for model species

### *Caenorhabditis elegans*



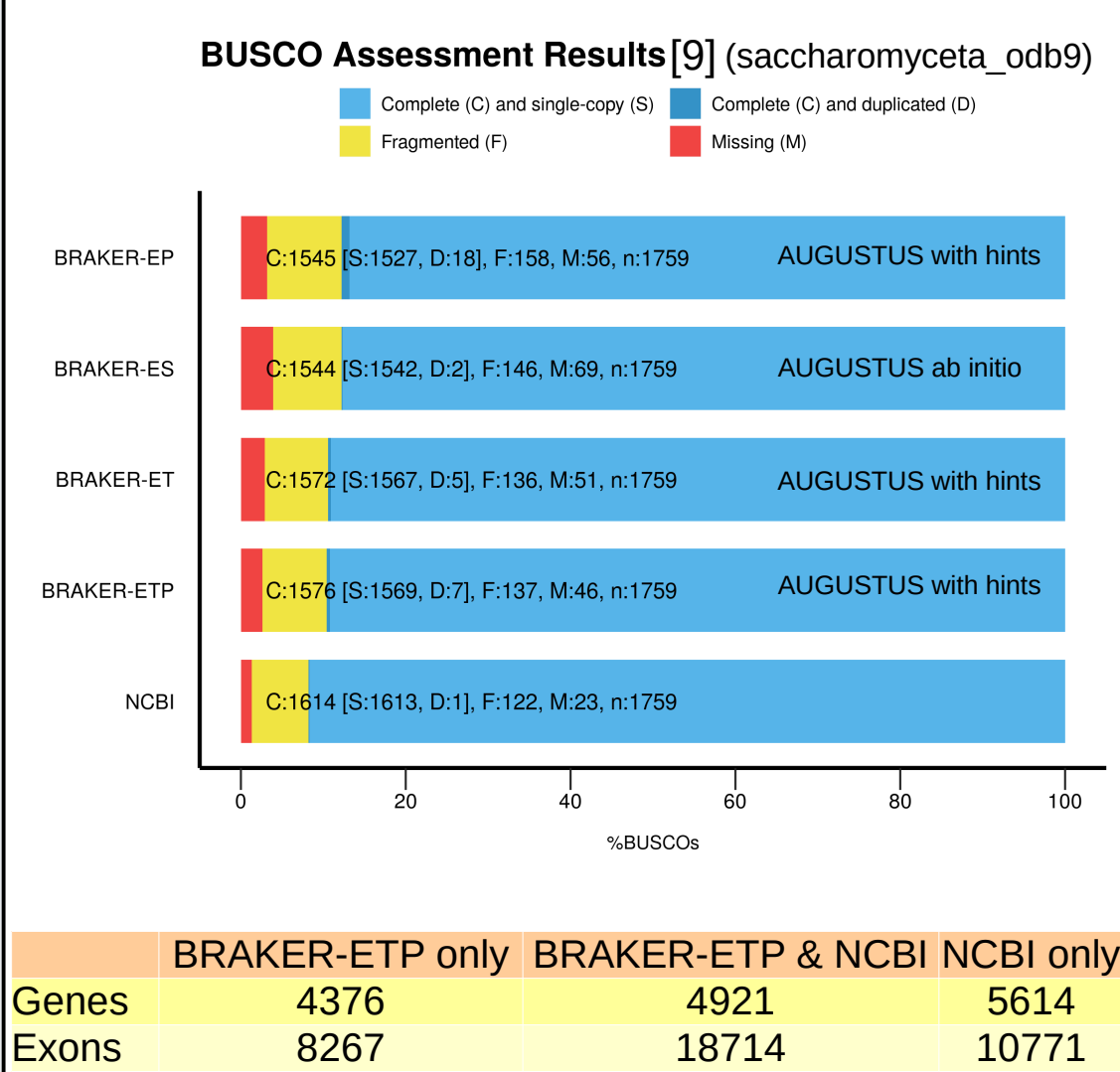
### *Drosophila melanogaster*



## Results for non-model species

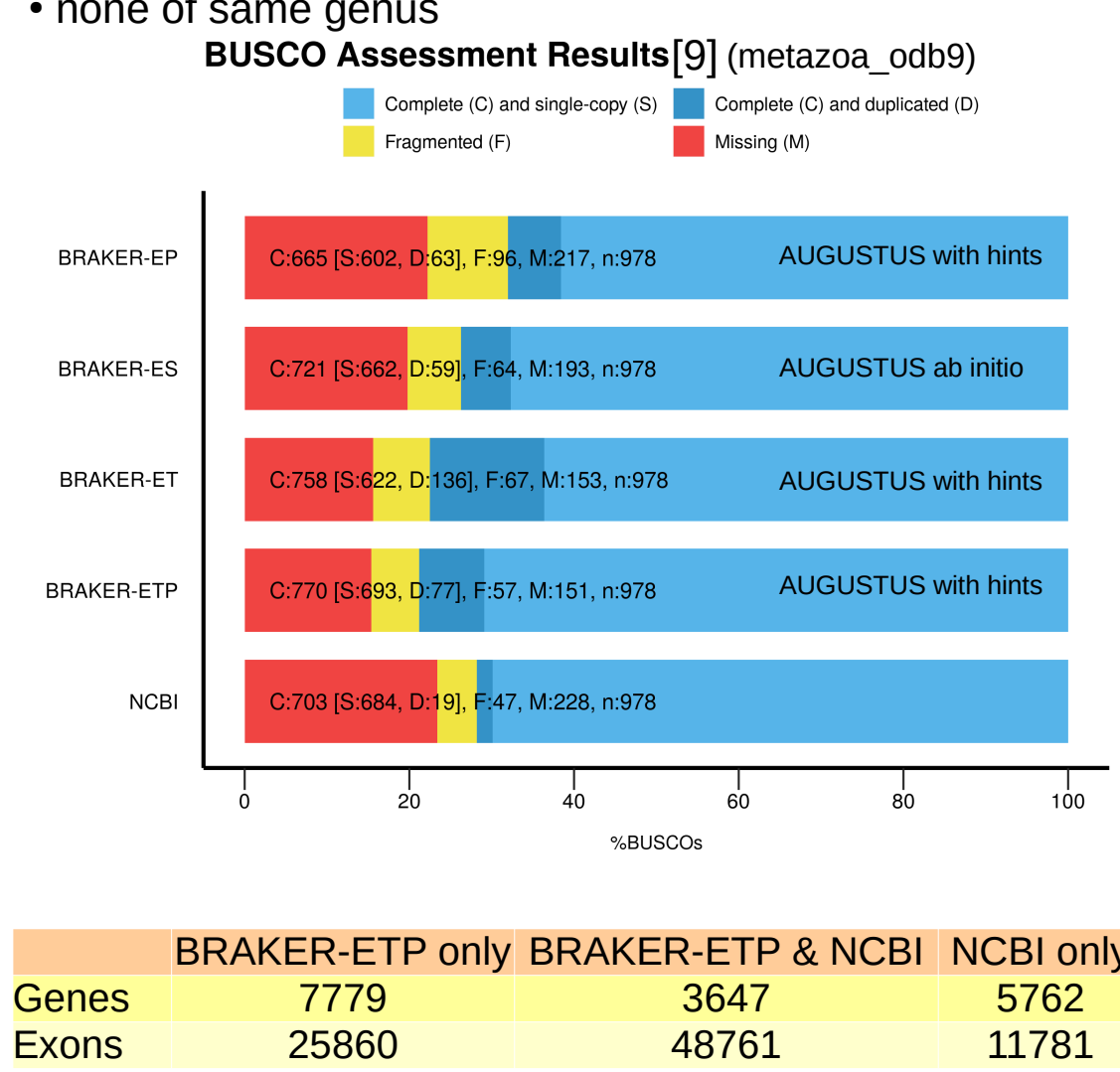
### *Verticillium dahliae*

Genome: GCF\_000150675.1\_ASM15067v2\_genomic.fna  
Genome size: 33MB  
RNA-seq: 50 M reads by VARUS with STAR [6]  
Protein mapping pipeline with OrthoDB v10 Fungi:  
• 12 species from the same order  
• 75 species from the same class  
• none of same genus



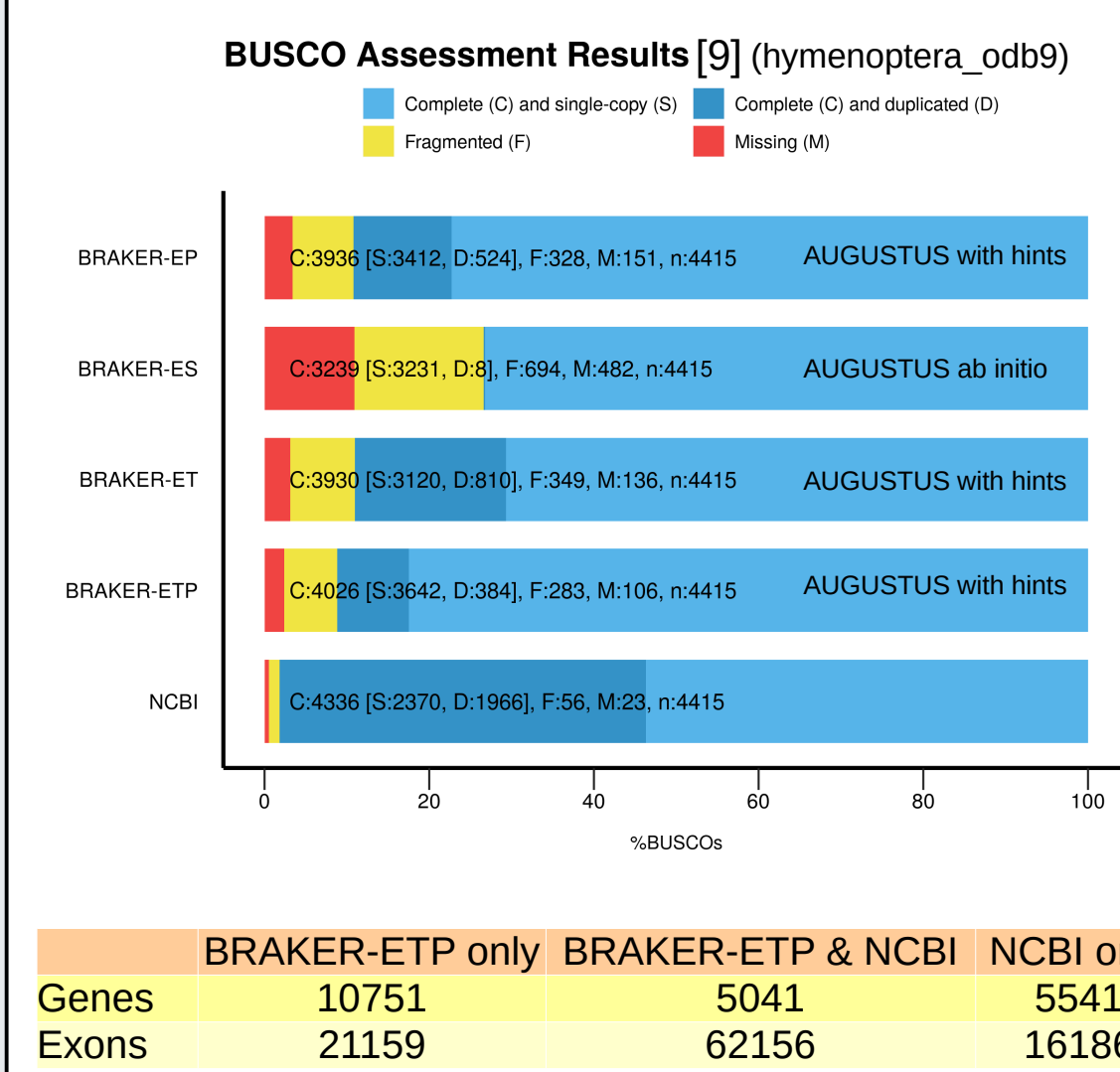
### *Hymenolepis microstoma*

Genome: GCA\_000469805.2\_HMIC002\_genomic.fna  
Genome size: 176 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Metazoa:  
• 1 species from the same order  
• 4 species from the same phylum  
• 28 species from the same kingdom  
• none of same genus



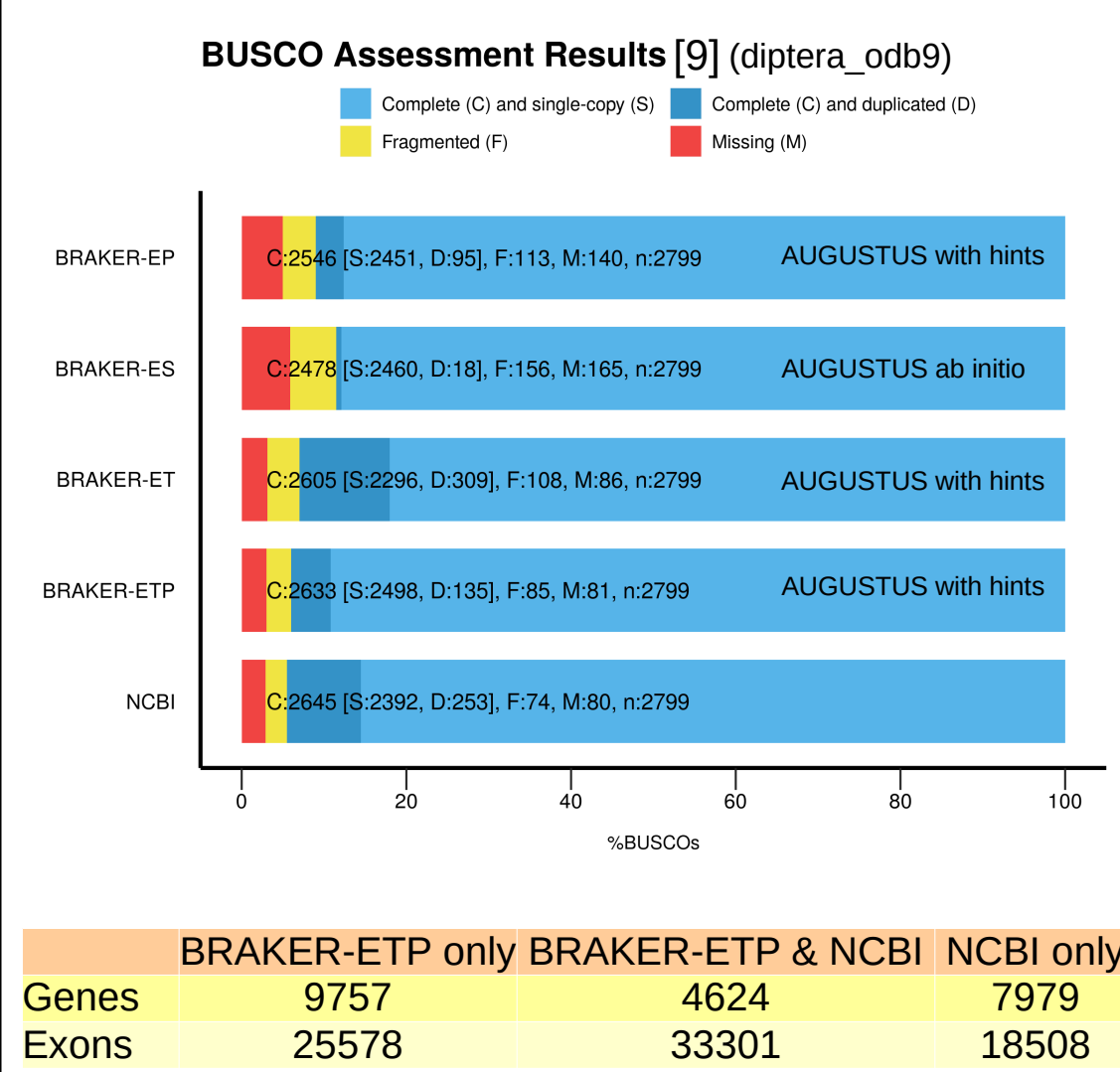
### *Bombus terrestris*

Genome: GCF\_000214255.1\_Bter\_1.0\_genomic.fna  
Genome size: 241 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Arthropoda:  
• 6 species from the same family  
• 39 species from the same order  
• none of same genus



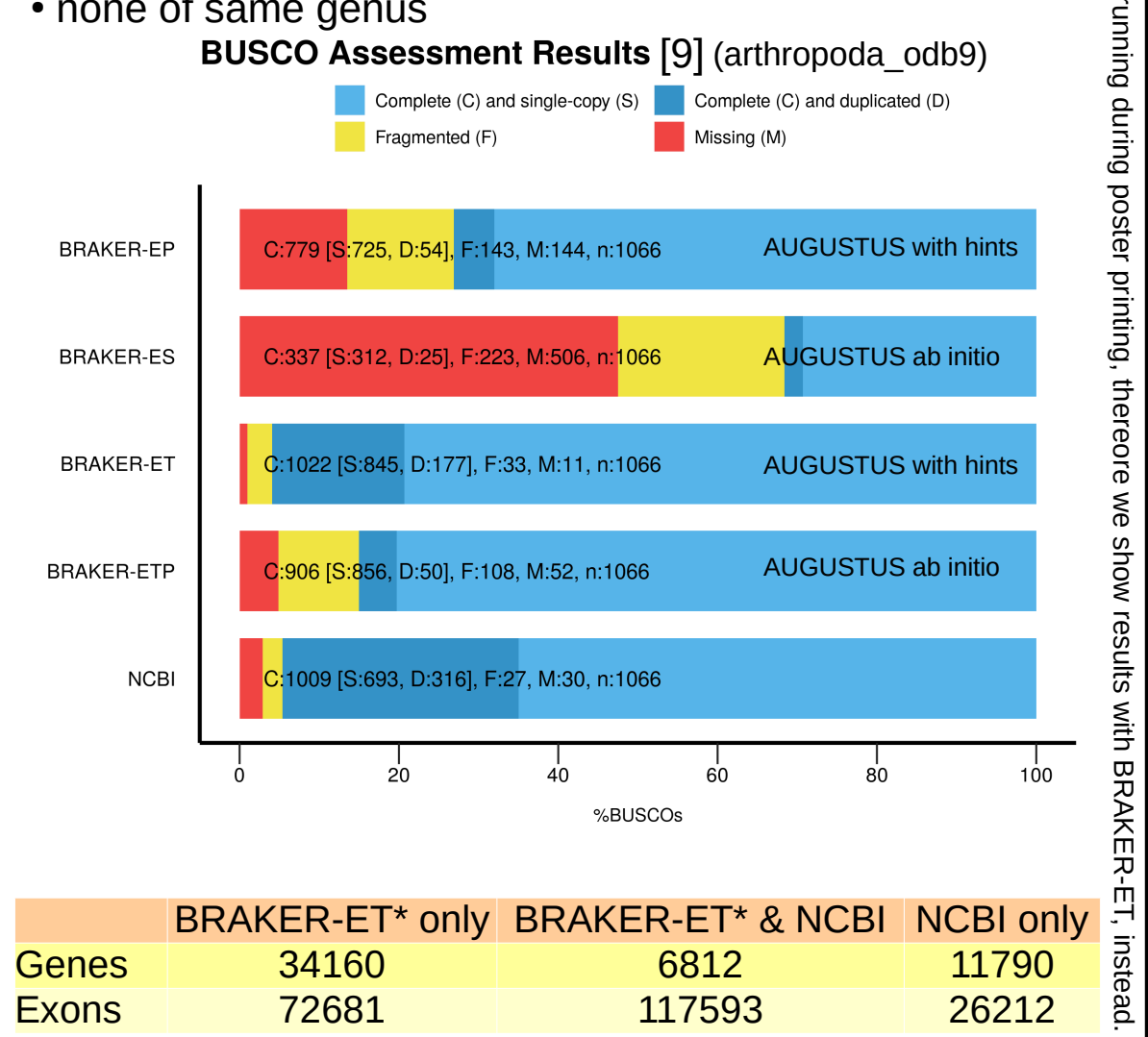
### *Anopheles gambiae*

Genome: GCF\_000005575.2\_AgamP3\_genomic.fna  
Genome size: 257 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Arthropoda:  
• 3 species from the same order  
• 42 species from the same order  
• none of same genus



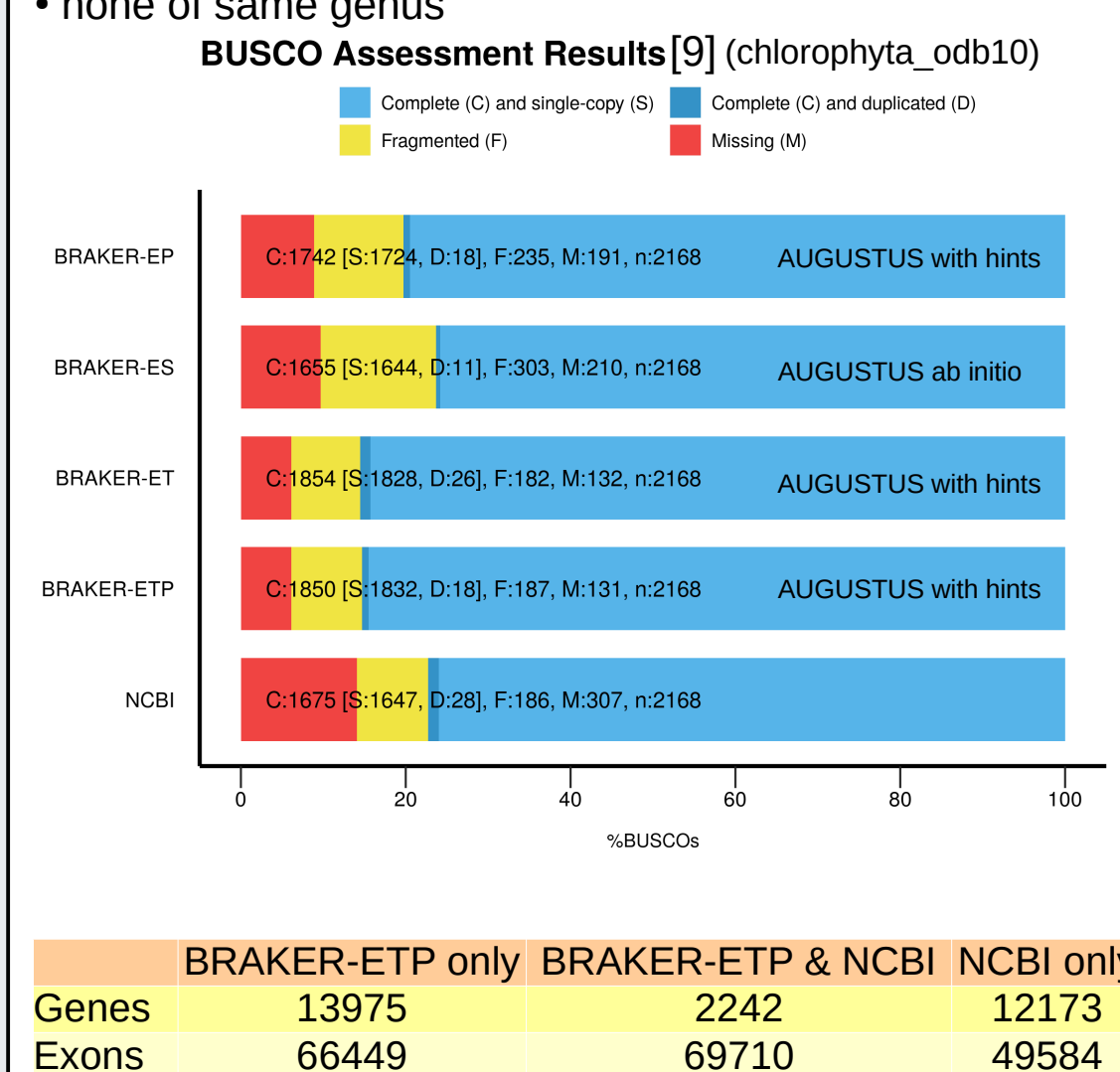
### *Parasteatoda tepidariorum*

Genome: GCF\_000365465.2\_Ptep\_2.0\_genomic.fna  
Genome size: 1.4 GB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Arthropoda:  
• 1 species from the same order  
• 10 species from the same class  
• 31 species from the same phylum  
• none of same genus



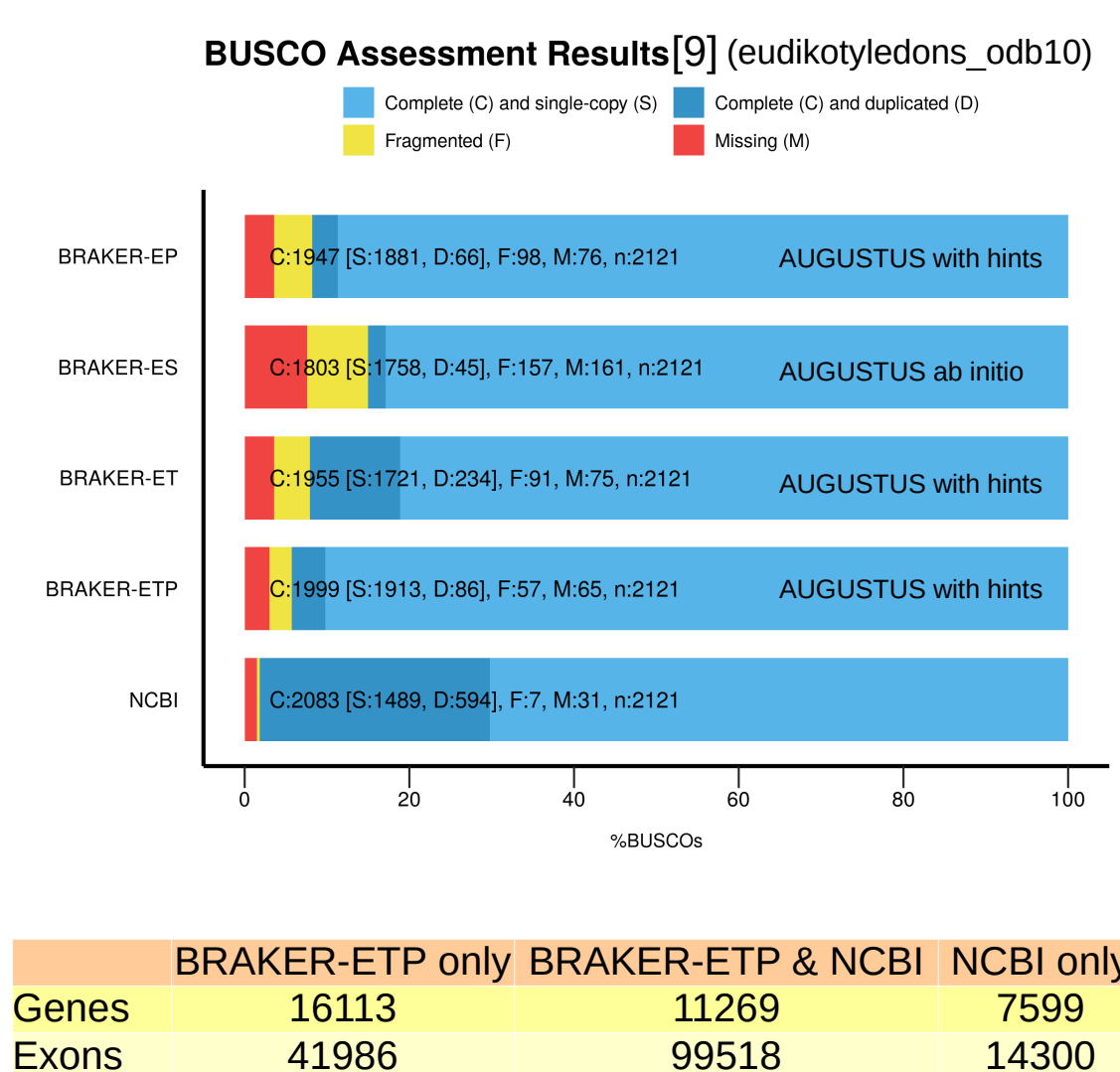
### *Chlamydomonas reinhardtii*

Genome: GCF\_000002595.1\_v3.0\_genomic.fna  
Genome size: 117 MB  
RNA-seq: 50 M reads by VARUS with STAR [6]  
Protein mapping pipeline with OrthoDB v10 Plants:  
• 3 species from the same order  
• 4 species from the same class  
• 14 species from the same phylum  
• none of same genus



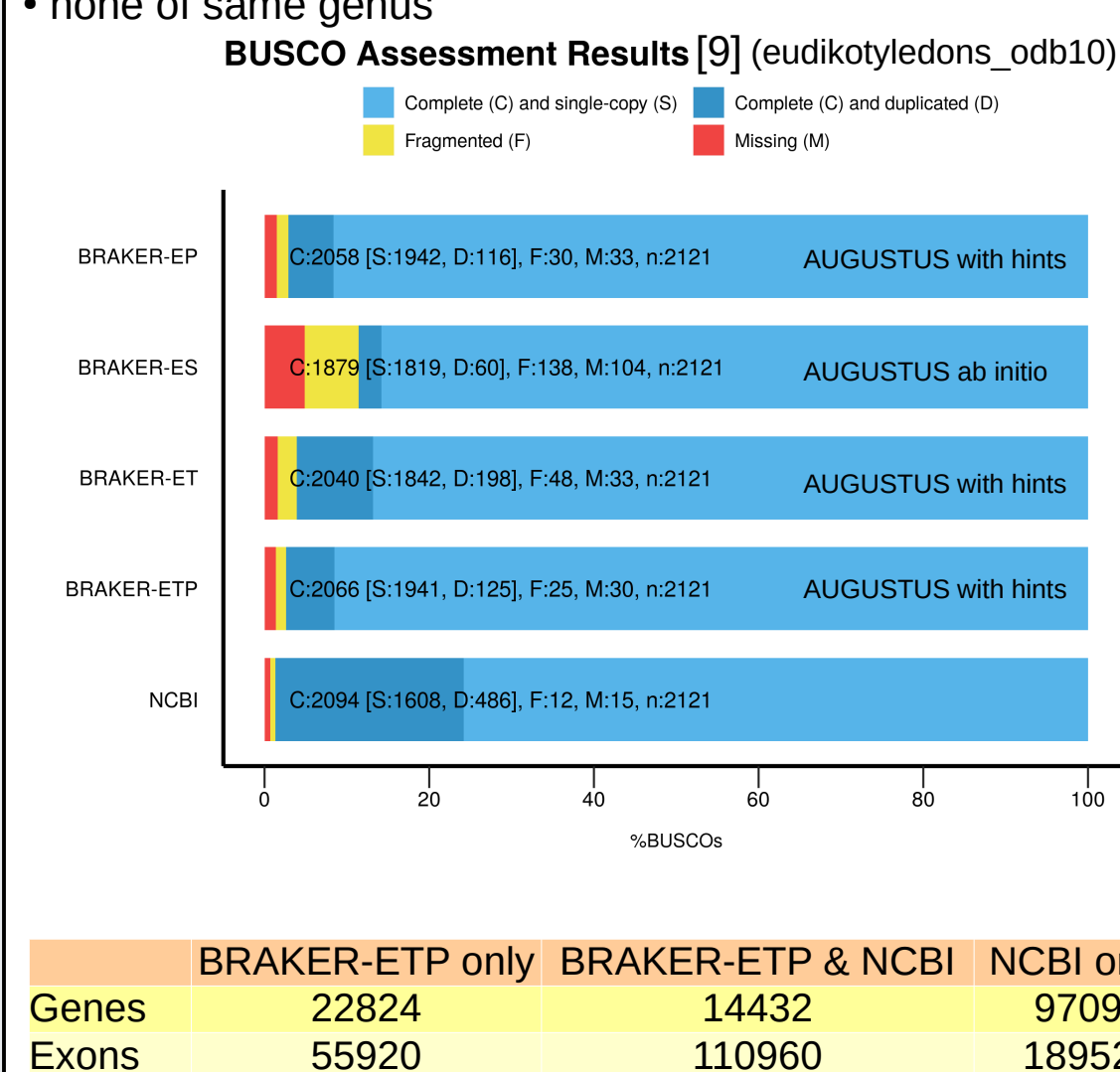
### *Cucumis sativus*

Genome: GCF\_000004075.2\_ASM407v2\_genomic.fna  
Genome size: 189 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Plants:  
• 2 species from the same family  
• 23 species from the same phylum  
• none of same genus



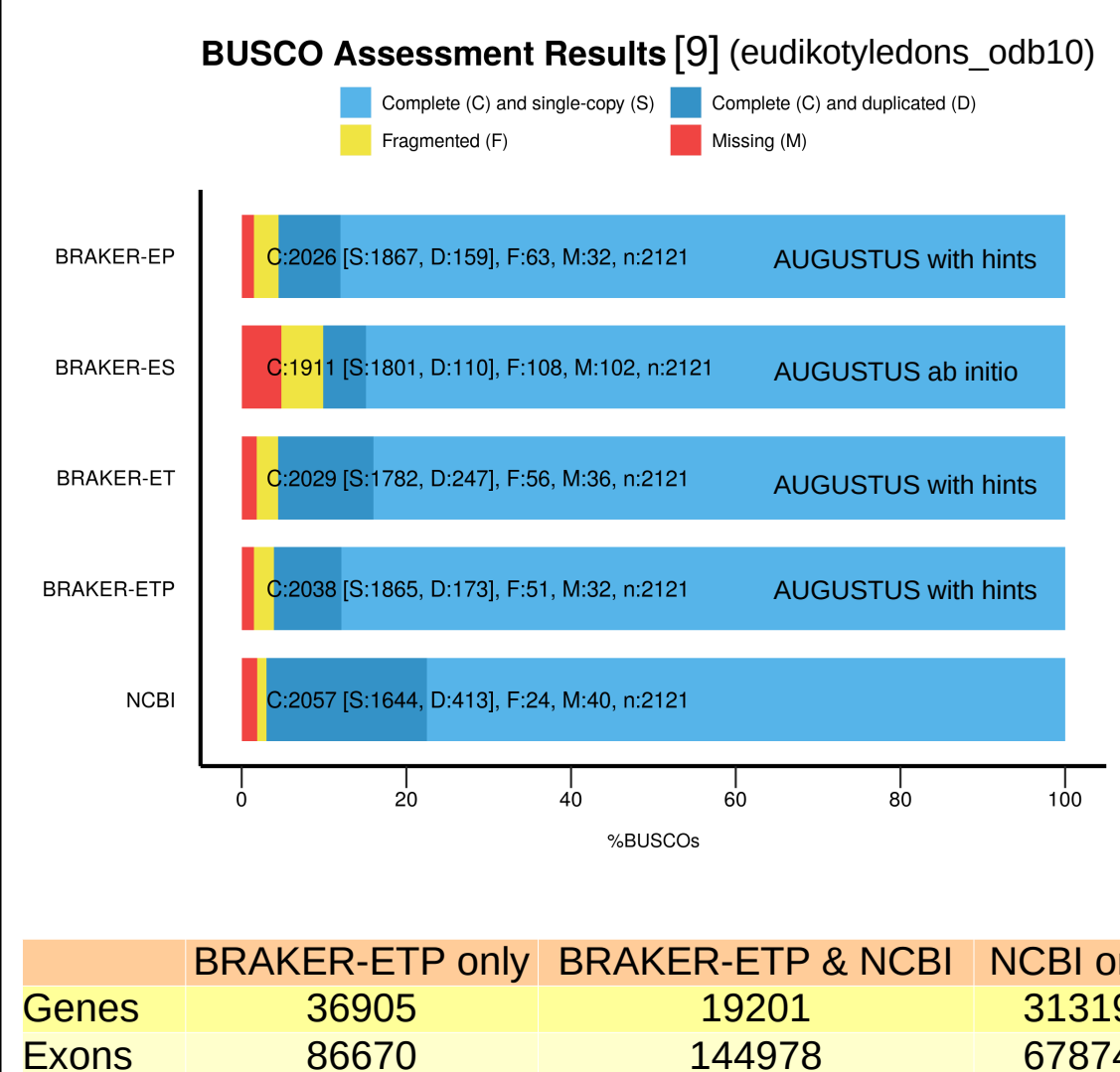
### *Fragaria vesca*

Genome: GCF\_000184155.1\_FraVesHawaii\_1.0\_genomic.fna  
Genome size: 208 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Plants:  
• 4 species from the same family  
• 8 species from the same order  
• 23 species from the same phylum  
• none of same genus



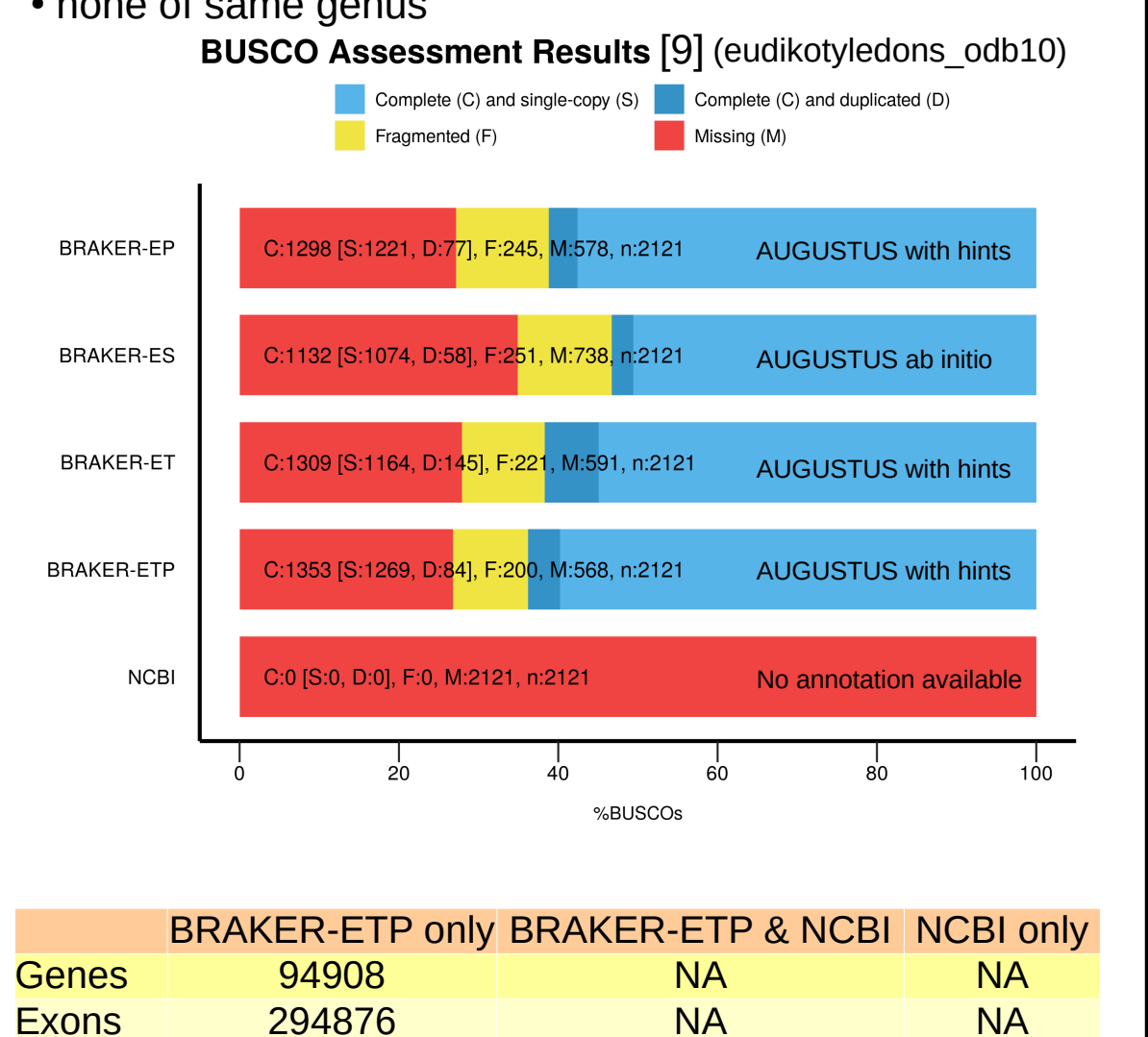
### *Medicago truncatula*

Genome: GCF\_000219495.3\_MedtrA17\_4.0\_genomic.fna  
Genome size: 399 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Plants:  
• 9 species from the same family  
• 12 species from the same phylum  
• none of same genus



### *Cannabis sativa*

Genome: GCA\_001865755.1\_ASM186575v1\_genomic.fna  
Genome size: 566 MB  
RNA-seq: 50 M reads by VARUS with HISAT2 [8]  
Protein mapping pipeline with OrthoDB v10 Plants:  
• 4 species from the same family  
• 8 species from the same order  
• 23 species from the same phylum  
• none of same genus



## References

[1] Hoff, Katharina J., et al. "BRAKER1: unsupervised RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS." Bioinformatics 32.5 (2015): 767-769.

[2] Lomsadze, Alexandre, Paul D. Burns, and Mark Borodovsky. "Integration of mapped RNA-Seq reads into automatic training of eukaryotic gene finding algorithm." Nucleic acids research 42.15 (2014): e119-e119.

[3] Stanke, Mario, et al. "Using native and syntenically mapped cDNA alignments to improve de novo gene finding." Bioinformatics 24.5 (2008): 637-644.

[4] Hoff, Katharina J., et al. "Whole-Genome Annotation with BRAKER" Springer Protocols (2019), in press.

[5] Stanke, Mario et al. "Automatic Genome Annotation Looping over Species" Poster PE0094 at PAG XXVII (2019).

[6] Kriventseva, Evgenia V., et al. "OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. Nucleic acids research (2018): doi.org/10.1093/nar/gky1053s.

[7] Dobin, Alexander, et al. "STAR: ultrafast universal RNA-seq aligner." Bioinformatics 29.1 (2013): 15-21.

[8] Daehwan, Kim, et al. "HISAT: a fast spliced aligner with low memory requirements." Nature methods 12.4 (2015): 357.

[9] Waterhouse, Robert M., et al. "BUSCO applications from quality assessments to gene prediction and phylogenomics." Molecular biology and evolution 35.3(2017):543-548