

## Genome analysis

# BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS

Katharina J. Hoff<sup>1,\*</sup>, Simone Lange<sup>1</sup>, Alexandre Lomsadze<sup>3</sup>, Mark Borodovsky<sup>2,3,4,\*</sup> and Mario Stanke<sup>1</sup>

<sup>1</sup>Ernst Moritz Arndt Universität Greifswald, Institute for Mathematics and Computer Science, 17487 Greifswald, Germany, <sup>2</sup>School of Computational Science and Engineering, Atlanta, GA 30332, USA, <sup>3</sup>Joint Georgia Tech and Emory University Wallace H Coulter Department of Biomedical Engineering, Atlanta, GA 30332, USA and <sup>4</sup>Moscow Institute of Physics and Technology, Dolgoprudny, Moscow Region, Russia

\*To whom correspondence should be addressed.

Associate Editor: Ivo Hofacker

Received on May 8, 2015; revised on October 2, 2015; accepted on October 26, 2015

## Abstract

**Motivation:** Gene finding in eukaryotic genomes is notoriously difficult to automate. The task is to design a work flow with a minimal set of tools that would reach state-of-the-art performance across a wide range of species. GeneMark-ET is a gene prediction tool that incorporates RNA-Seq data into unsupervised training and subsequently generates *ab initio* gene predictions. AUGUSTUS is a gene finder that usually requires supervised training and uses information from RNA-Seq reads in the prediction step. Complementary strengths of GeneMark-ET and AUGUSTUS provided motivation for designing a new combined tool for automatic gene prediction.

**Results:** We present BRAKER1, a pipeline for unsupervised RNA-Seq-based genome annotation that combines the advantages of GeneMark-ET and AUGUSTUS. As input, BRAKER1 requires a genome assembly file and a file in *bam*-format with spliced alignments of RNA-Seq reads to the genome. First, GeneMark-ET performs iterative training and generates initial gene structures. Second, AUGUSTUS uses predicted genes for training and then integrates RNA-Seq read information into final gene predictions. In our experiments, we observed that BRAKER1 was more accurate than MAKER2 when it is using RNA-Seq as sole source for training and prediction. BRAKER1 does not require pre-trained parameters or a separate expert-prepared training step.

**Availability and implementation:** BRAKER1 is available for download at <http://bioinf.uni-greifswald.de/bioinf/braker/> and <http://exon.gatech.edu/GeneMark/>.

**Contact:** [katharina.hoff@uni-greifswald.de](mailto:katharina.hoff@uni-greifswald.de) or [borodovsky@gatech.edu](mailto:borodovsky@gatech.edu)

**Supplementary information:** [Supplementary data](#) are available at *Bioinformatics* online.

## 1 Introduction

Transcriptome sequencing data, RNA-Seq reads, aligned to a genome sequence have great potential to improve the accuracy of structural genome annotation: spliced alignments may indicate intron positions, and coverage increase and decrease along genomic sequence may

indicate locations of exon-noncoding region borders. Nevertheless, RNA-Seq coverage alone is no reliable indicator of protein coding regions (Hoff and Stanke, 2015).

The prediction of protein coding regions in genomes is often accomplished by tools that use statistical models. Some gene prediction tools can additionally use RNA-Seq to improve prediction accuracy.



introns predicted *ab initio* and also supported by RNA-Seq read mapping) is not significantly affected by TEs masking since TEs have no anchored introns. However, at the prediction step TEs can corrupt gene prediction. For this reason, soft masking of genomic sequence is recommended before execution of BRAKER1. In this publication we used RepeatModeler to generate repeat library and RepeatMasker to mask sequence [Smit, A.F.A. and Hubley, R. (2008–2015) RepeatModeler Open-1.0, <http://www.repeatmasker.org/>].

### 3 Results and discussion

When comparing BRAKER1 to MAKER2 (details on the MAKER2 run are described in [Supplementary Materials](#)), BRAKER1 gains on average 15% points in accuracy on gene level (see [Table 1](#)). A gene may have several transcripts both in the reference and in the predicted gene set. When computing transcript level accuracy, each transcript variant is counted as a TP/FP/FN on its own. When computing gene level accuracy, a predicted gene is counted as TP if at least one of the predicted gene's transcripts matches correctly a reference transcript. (For details, see documentation of the EVAL package in [Keibler and Brent, 2003](#)). We should remind that in these runs of BRAKER1 and MAKER2 we use only RNA-Seq information as the source of external evidence.

Notably, [Reid et al. \(2014\)](#) developed a pipeline, SnowyOwl, with GeneMark-ES ([Ter-Hovhannissyan et al., 2008](#)) and AUGUSTUS to predict genes in fungal genomes. SnowyOwl attempts to improve prediction accuracy by selecting a gene variant with the highest homology score from a set of predicted gene variants in the same locus. This pipeline requires protein database information as additional external resource. We did not include SnowyOwl into comparisons since it cannot work without protein information.

Yet another recently developed automatic pipeline for fungal genome annotation utilizing RNA-Seq data is CodingQuarry ([Testa et al., 2015](#)). Tests of CodingQuarry on the *S.pombe* genome demonstrated that it makes an improvement in comparison to MAKER2, however, BRAKER1 is on average ~4% more accurate on gene level than CodingQuarry ([Table 1](#)).

In attempt to elucidate the roles and contributions of separate gene finding tools in BRAKER1 and MAKER2 as well as the role of repeat masking and incorporation of RNA-Seq information, we show the following results: values of *ab initio* accuracies of GeneMark-ET and AUGUSTUS for repeat masked and unmasked genomes are provided in [Supplementary Tables 1.1 and 1.2](#), respectively. These two tables show the BRAKER1 accuracies as well.

Given that BRAKER1 uses AUGUSTUS trained on genes most reliably predicted by GeneMark-ET, and since AUGUSTUS incorporates RNA-Seq into the prediction step, we expect to see an increase in accuracy when comparing BRAKER1 (the 'hints supported' AUGUSTUS) with GeneMark-ET and with *ab initio* AUGUSTUS. This is the case for *A.thaliana*, *C.elegans* and *D.melanogaster*; on the fungus *S.pombe*, GeneMark-ET shows even higher accuracy than the current formal output of BRAKER1 (see [Supplementary Tables S1.1 and S1.2](#)).

Repeat masking on genome scale is an optional pre-processing step for running BRAKER1; still, taking this step does not significantly affect prediction accuracy ([Supplementary Table S1.2](#)).

To quantify the accuracy that MAKER2 gains by combining predictions from SNAP, AUGUSTUS and GeneMark-ES, from masking and from RNA-Seq information, we show the *ab initio* accuracies of the three gene-finders on unmasked genomes ([Supplementary Table S1.3](#)). These results show that the unsupervised training of GeneMark-ES

allows to get accuracy close to or even better (*S.pombe*) than the one achieved by the MAKER2 training with utilization of RNA-Seq information.

Interestingly, we have observed ([Supplementary Table S1.4](#)) that the prediction accuracy of *ab initio* AUGUSTUS fully automatically trained by BRAKER1 is in most cases few percent lower than the *ab initio* accuracy of AUGUSTUS utilizing the packaged parameter files (obtained by supervised training). However, after adding RNA-Seq information, prediction accuracy of BRAKER1 ([Supplementary Table S1.2](#)) clearly exceeds accuracy of *ab initio* predictions made by 'expert trained' AUGUSTUS.

In summary, we have observed that when the transcript data (RNA-Seq) is used as a sole source of evidence, BRAKER1 predicts genes more accurately than MAKER2 and CodingQuarry. The gain of accuracy is due to i/ use of GeneMark-ET and generation of accurate training sets for AUGUSTUS as well as ii/ use of hints originated from mapping of RNA-Seq reads that AUGUSTUS incorporates in the final gene prediction step.

In contrast to running MAKER2, running BRAKER1 is a 'one step process', meaning that after starting it once, it will execute training and prediction in fully automated mode without manual command execution.

The example running time of BRAKER1 is ~17.5 hours on a single CPU for training and prediction on *D.melanogaster*; running time can be improved by use of parallel processors.

### Acknowledgement

We would like to thank Mark Yandell and Carson Holt for valuable advice on running MAKER2.

### Funding

This work is supported in part by the US National Institutes of Health grant HG000783 to MB and by the German Research Foundation (DFG) grant STA 1009/10-1 to MS.

*Conflict of Interest:* none declared.

### References

- Hoff,K.J. and Stanke,M. (2015) Current methods for automated annotation of protein-coding genes. *Curr. Opin. Insect Sci.*, 7, 8–14.
- Holt,C. and Yandell,M. (2011) MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics*, 12, 491.
- Keibler,E. and Brent,M.R. (2003) Eval: a software package for analysis of genome annotations. *BMC Bioinformatics*, 4, 50.
- Lomsadze,A. et al. (2014) Integration of mapped RNA-Seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.*, 42, e119.
- Reid,I. et al. (2014) SnowyOwl: accurate prediction of fungal genes by using RNA-Seq and homology information to select among *ab initio* models. *BMC Bioinformatics*, 15, 229.
- Stanke,M. et al. (2008) Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics*, 24, 637.
- Steijger,T. et al. (2013) Assessment of transcript reconstruction methods for RNA-seq. *Nat. Methods*, 10, 1177–1184.
- Ter-Hovhannissyan,V. et al. (2008) Gene prediction in novel fungal genomes using an *ab initio* algorithm with unsupervised training. *Genome Res.*, 18, 1979–1990.
- Testa,A.C. et al. (2015) CodingQuarry: highly accurate hidden Markov model gene prediction in fungal genomes using RNA-seq transcripts. *BMC Genomics* 16, 170.