

METHODOLOGY ARTICLE

Open Access



Combining RNA-seq data and homology-based gene prediction for plants, animals and fungi

Jens Keilwagen^{1*}, Frank Hartung¹, Michael Paulini², Sven O. Twardziok³ and Jan Grau⁴ 

Abstract

Background: Genome annotation is of key importance in many research questions. The identification of protein-coding genes is often based on transcriptome sequencing data, *ab-initio* or homology-based prediction. Recently, it was demonstrated that intron position conservation improves homology-based gene prediction, and that experimental data improves *ab-initio* gene prediction.

Results: Here, we present an extension of the gene prediction program GeMoMa that utilizes amino acid sequence conservation, intron position conservation and optionally RNA-seq data for homology-based gene prediction. We show on published benchmark data for plants, animals and fungi that GeMoMa performs better than the gene prediction programs BRAKER1, MAKER2, and CodingQuarry, and purely RNA-seq-based pipelines for transcript identification. In addition, we demonstrate that using multiple reference organisms may help to further improve the performance of GeMoMa. Finally, we apply GeMoMa to four nematode species and to the recently published barley reference genome indicating that current annotations of protein-coding genes may be refined using GeMoMa predictions.

Conclusions: GeMoMa might be of great utility for annotating newly sequenced genomes but also for finding homologs of a specific gene or gene family. GeMoMa has been published under GNU GPL3 and is freely available at <http://www.jstacs.de/index.php/GeMoMa>.

Keywords: Homology-based gene prediction, RNA-seq, Genome annotation

Background

The annotation of protein-coding genes is of critical importance for many fields of biological research including, for instance, comparative genomics, functional proteomics, gene targeting, genome editing, phylogenetics, transcriptomics, and phylostratigraphy. The process of annotating protein-coding genes to an existing genome (assembly) can be described as specifying the exact genomic location of genes comprising all (partially) coding exons. A difficulty in gene annotation is distinction between protein-coding genes, transposons and pseudogenes.

Genome annotation pipelines utilize three main sources of information, namely evidence from wet-lab transcriptome studies [1, 2], *ab-initio* gene prediction based on general features of (protein-coding) genes [3, 4], and homology-based gene prediction relying on gene models of (closely) related, well-annotated species [5–7].

Experimental data allow for inferring coverage of gene predictions and splice sites bordering their exons, which may assist computational *ab-initio* or homology-based approaches. Due to the progress in the field of next generation sequencing, RNA-seq has revolutionized transcriptomics [8]. Today, RNA-seq data is available for a wide range of organisms, tissues and environmental conditions, and can be utilized for genome annotation pipelines.

In recent years, several programs have been developed that combine multiple sources allowing for a more accurate prediction of protein-coding genes [9–11]. MAKER2

*Correspondence: jens.keilwagen@julius-kuehn.de

¹Institute for Biosafety in Plant Biotechnology, Julius Kühn-Institut (JKI) - Federal Research Centre for Cultivated Plants, D-06484, Quedlinburg, Germany
Full list of author information is available at the end of the article



is a pipeline that integrates support of different resources including *ab-initio* gene predictors and RNA-seq data [9]. CodingQuarry is a pipeline for RNA-Seq assembly-supported training and gene prediction, which is only recommended for application to fungi [10]. Recently, [11] published BRAKER1 a pipeline for unsupervised RNA-seq-based genome annotation that combines the advantages of GeneMark-ET [12] and AUGUSTUS [4].

Here, we present an extension of GeMoMa [7] that utilizes RNA-seq data in addition to amino acid sequence and intron position conservation. We investigate the performance of GeMoMa on publicly available benchmark data [11] and compare it with state-of-the-art competitors [9–11].

Subsequently, we demonstrate how combining homology-based predictions based on gene models from multiple reference organisms can be used to improve the performance of GeMoMa. Finally, we apply GeMoMa to four nematode species provided by Wormbase [13] and to the recently published barley reference genome [14], where GeMoMa predictions will be included into future versions of the corresponding genome annotations.

Methods

In this section, we describe recent extensions of GeMoMa to make use of evidence from RNA-seq data, the RNA-seq pipelines used and the data considered in the benchmark and application studies.

GeMoMa using RNA-seq

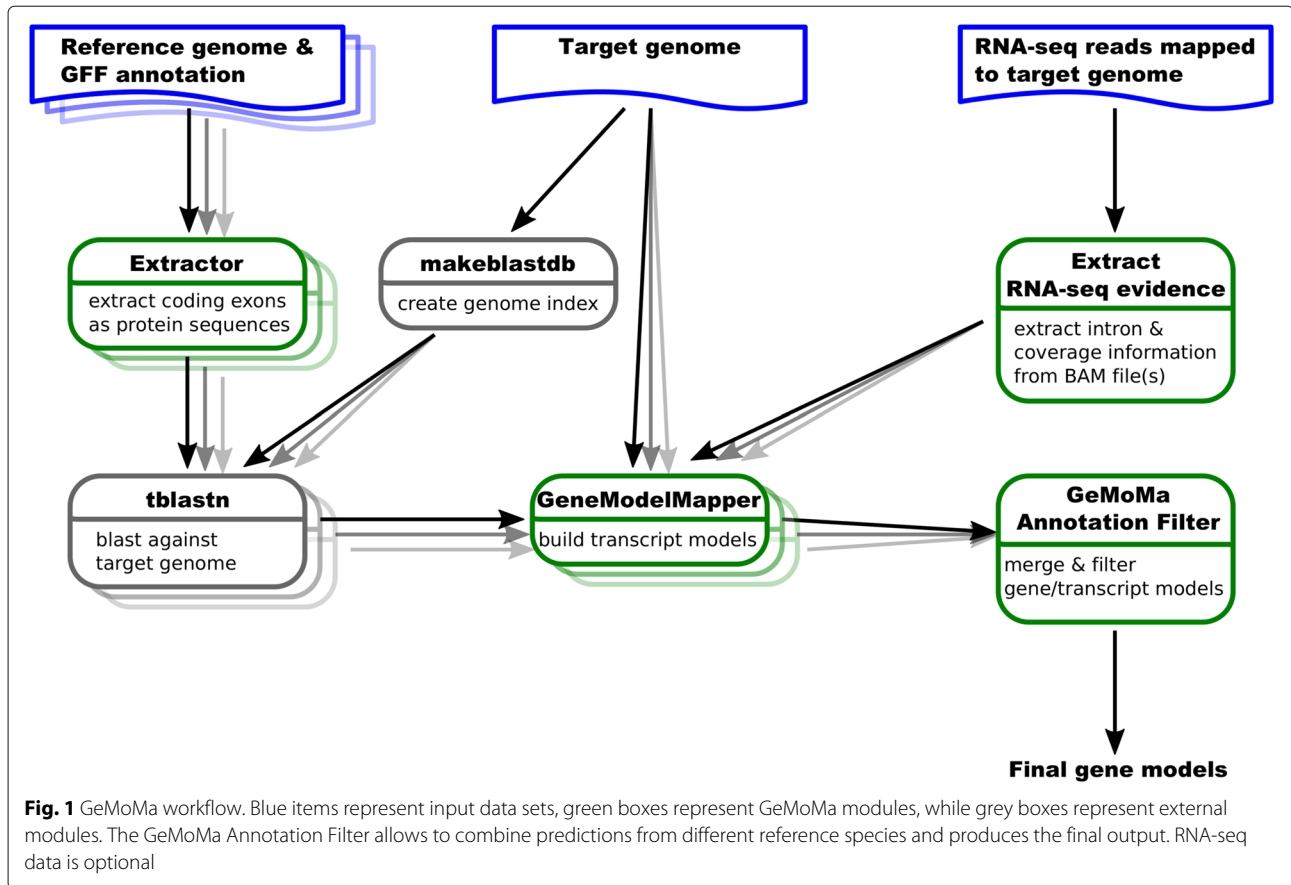
GeMoMa predicts protein-coding genes utilizing the general conservation of protein-coding genes on the level of their amino acid sequence and on the level of their intron positions, i.e., the locations of exon-exon boundaries in CDSs [7]. To this end, sequences of (partially) protein-coding exons are extracted from well-annotated reference genomes. Individual exons are then matched to loci on the target genome using tblastn [15], matches are adjusted for proper splice sites, start codons and stop codons, respectively, and joined to full, protein-coding genes models. In this process, the conserved dinucleotides GT and GC for donor splice sites, and AG for acceptor splice sites have been used for the identification of splice sites bordering matches to the (partially) protein-coding exons of the reference transcripts. The improved version of GeMoMa may now also include experimental splice site evidence extracted from mapped RNA-seq data to improve the accuracy of splice site and, hence, exon annotation. We visualize the extended GeMoMa pipeline in Fig. 1.

Starting from mapped RNA-seq data, the module *Extract RNA-seq evidence* (ERE) allows for extracting introns and, if user-specified, read coverage of genomic regions. GeMoMa filters these introns using a user-specified minimal number of split reads within the

mapped RNA-seq data. Introns passing this filter define donor and acceptor splice sites, which are treated independently within GeMoMa. If splice sites with experimental evidence have been detected in a genomic region with a good match to an exon of a reference transcript, these are collected for further use. If no splice sites with experimental evidence have been detected in a genomic region with a good match to an exon of a reference transcript, GeMoMa resorts to conserved dinucleotides allowing to identify gene models that are not covered by RNA-seq data due to, e.g., very specifically or lowly expressed transcripts. Combining two potential exons, all in-frame combinations using the collected donor and acceptor splice sites are tested and scored according to the reference transcript. The best combination is used for the prediction.

Based on this experimental evidence, the improved version of GeMoMa provides several new properties reported for gene predictions. The most prominent features are *transcript intron evidence* (tie) and *transcript percentage coverage* (tpc). The tie of a transcript varies between 0 and 1, and corresponds to the fraction of introns (i.e., splice sites of two neighboring exons) that are supported by split reads in the mapped RNA-seq data. In case of transcripts comprising a single coding exon, NA is reported. The tpc of a transcript also varies between 0 and 1, and corresponds to the fraction of (coding) bases of a predicted transcript that are also covered by mapped reads in the RNA-seq data. Further properties reported by GeMoMa are i) *tae* and *tde*, the percentages of acceptor and donor sites, respectively, with RNA-seq evidence, ii) *minCov* and *avgCov*, the minimum and average coverage, respectively, of the predicted transcript, and iii) *minSplitReads*, the minimum number of split reads supporting any of the predicted introns of a transcript. Optionally, GeMoMa reports *pAA* and *iAA*, the percentage of positive-scoring and identical amino acids in a pairwise alignment, if the reference protein is provided as input.

GeMoMa allows for computing and ranking multiple predictions per reference transcript, but does not filter these predictions. Predictions of different reference transcripts might be highly overlapping or even identical, especially if the reference transcripts are from the same gene family. Since GeMoMa 1.4, the default parameters for number of predictions and contig threshold have been changed which might lead to an increased number of highly overlapping or identical predictions. In addition, it might be beneficial to run GeMoMa starting from multiple reference species to broaden the scope of transcripts covered by the predictions. However, these may also result in redundant predictions for, e.g., orthologs or paralogs stemming from the different reference species considered. To handle such situations, the new module *GeMoMa annotation filter* (GAF) of the improved version



of GeMoMa now allows for joining and reducing such predictions using various filters. Filtering criteria comprise the relative GeMoMa score of a predicted transcript, filtering for complete predictions (starting with start codon and ending with stop codon), and filtering for evidence from multiple reference organisms. In addition, GAF also joins duplicate predictions that originate from different reference transcripts.

Initially, GAF filters predictions based on their relative GeMoMa score, i.e., the GeMoMa score divided by the length of the predicted protein. This filter removes spurious predictions. Subsequently, the predictions are clustered based on their genomic location. Overlapping predictions on the same strand yield a common cluster. For each cluster, the prediction with the highest GeMoMa score is selected. Non-identical predictions overlapping the high-scoring prediction with at least a user-specified percentage of borders (i.e., splice sites, start and stop codon, cf. *common border filter*) are treated as alternative transcripts. Predictions that have completely identical borders to any previously selected prediction are removed and only listed in the GFF attribute field *alternative*. All filtered predictions of a cluster are assigned to one gene with a generic gene name. Finally, GAF checks for nested

genes in the cluster looking for discarded predictions that do not overlap with any selected prediction, which are recovered. In the benchmark studies comparing GeMoMa with state-of-the-art competitors, we directly use the GAF results without any further filters on attributes reported by the GeMoMa pipeline.

In addition to the modules for annotating a genome (assembly) described above, we also provide two additional modules in GeMoMa for analyzing and comparing to prediction to a reference annotation. The module *CompareTranscripts* determines that CDS of the reference annotation with the largest overlap with the prediction utilizing the F_1 measure as objective function [7]. The module *AnnotationEvidence* computes tie and tpc of all CDSs of a given annotation. Hence, these two modules can be used to determine, whether a prediction is known, partially known or new and whether the overlapping annotation has good RNA-seq support.

MAKER2 predictions

Recently, we have shown that GeMoMa outperforms state-of-the-art homology-based gene predictors [7]. We are not aware of any homology-based gene prediction program that allows for incorporating of RNA-seq

data. Hence, we provide predictions of MAKER2 using the same reference proteins as GeMoMa for a minimal comparison. Internally, MAKER2 uses exonerate [5] for homology-based gene prediction. We run MAKER2 with default parameters except `protein2genome=1`, and `genome` and `protein` set to the respective input files. In addition, we run MAKER2 using (i) RNA-seq data in form of Trinity 2.4 transcripts (`-jaccard_clip`) [16], (ii) homology in form of proteins of one related reference species, and (iii) *ab-initio* gene prediction in form of Augustus 3.3 [4]. In this case, we run MAKER2 with default parameters except `genome`, `est`, `protein`, and `augustus_species`, which have been set to the corresponding species. For comparison, we run Maker2 with the same parameter settings but using the GeMoMa predictions for `protein_gff` instead of using `protein`.

RNA-seq pipelines

Computational pipelines have been used to infer gene annotation from RNA-seq data produced by next generation sequencing methods. Dozens of tools and tool combinations have been proposed. Here, we focus on the short read mapper TopHat2 [17], the transcript assemblers Cufflinks [1] and StringTie [2], and the coding sequence predictor TransDecoder [16]. Based on the transcript assemblers, we build two RNA-seq pipelines following the instructions in [11].

Data

For the benchmark studies, we consider target species and their genome versions as specified in the BRAKER1 supplement. For the homology-based prediction by GeMoMa, we choose one closely related reference species per target species that are sequenced and annotated [13, 18–20]. For these species, we consider the latest genome versions available (Additional file 1: Table S1). For the analysis of *C. elegans*, we use the manually curated gene set of *C. briggsae* provided by Wormbase. In addition, we use the experimental evidence from RNA-seq data referenced in the BRAKER1 publication.

For the analysis of the four nematode species, *C. brenneri*, *C. briggsae*, *C. japonica*, and *C. remanei*, we use the genome assembly and gene annotation of Wormbase WS257 [13]. We choose the model organism *C. elegans* as reference species (Additional file 1: Table S2). In addition to genome assembly and gene annotation, we also use publicly available RNA-seq data of these four nematode species, which have been mapped by Wormbase using STAR [21]. We used a minimum intron size of 25 bp, a maximum intron size of 15Kb, specify that only reads mapping once or twice on the genome are reported, and alignments are reported only if their ratio of mismatches to mapped length is less than 0.02. In accordance

with the previous benchmark study, we use the manually curated gene set of Wormbase.

For the analysis of barley, we use the latest genome assembly and gene annotation [14]. As reference species, we choose *A. thaliana* [22], *B. distachyon* [23], *O. sativa* [24], and *S. italica* [25] (Additional file 1: Table S2). In addition to genome assembly and gene annotation, we also used RNA-seq data from four different public available data sets (ERP015182, ERP015986, SRP063318, SRP071745). Reads were mapped and assembled using Hisat2 and StringTie [26]. As reference annotation, we used the union of high and low confidence annotation.

As independent evidence for validating GeMoMa predictions in the nematode species and barley, we use ESTs and cDNAs. While Wormbase provides coordinates for *best BLAT matches*, we adapt the pipeline and download all available EST from NCBI and map them to the genome using BLAT [27].

Results and discussion

Benchmark

The comparison of different software pipelines is often critical as a) specific parameters settings might be crucial for good results and b) different input might be used. For these reasons, we designed the benchmark as follows. First, we use publicly available gene predictions results. Second, we limit the number of reference species to one in the initial study.

We used GeMoMa for predicting the gene annotations of *A. thaliana*, *C. elegans*, *D. melanogaster*, and *S. pombe*. In Table 1, we summarize the performance of BRAKER1, MAKER2, and CodingQuarry as reported in Hoff et al. [11], as well as the performance of GeMoMa with and without RNA-seq evidence, purely RNA-seq-based pipelines and various MAKER2 predictions. The results of CodingQuarry reported by Hoff et al. [11] deviate substantially from those originally reported by Testa et al. [10]. We find that the performance of CodingQuarry is highly sensitive to RNA-seq processing, whereas the performance of GeMoMa is barely affected (Additional file 1: Table S5). For all comparisons, we provide sensitivity (Sn) and specificity (Sp) for the categories gene, transcript, and exon, respectively [28]. In addition, we compare CodingQuarry with GeMoMa for *S. cerevisiae* (Additional file 1: Table S6).

First, we compare the two purely homology-based predictions, namely on the one hand side MAKER2 using exonerate and on the other hand side GeMoMa without RNA-seq data. In all cases, we use the same reference species and reference proteins. We find that MAKER2 using only homologous proteins has a higher exon specificity than GeMoMa without RNA-seq data for *C. elegans*, while the opposite is true for all other categories and target species.

Table 1 Benchmark results on the BRAKER1 test sets

| | MAKER2 ⁺ (exonerate) | GeMoMa ⁺ without RNA-seq data | GeMoMa ⁺ with RNA-seq data | RNAseq- Cufflinks | RNAseq- StringTie | BRAKER1* | MAKER2* | CodingQuarry* | MAKER2 ⁺ (exonerate, Trinity, Augustus) | MAKER2 ⁺ (GeMoMa, Trinity, Augustus) |
|----------------------------------|------------------------------------|---|--|----------------------|----------------------|-------------|---------|---------------|---|---|
| <i>Arabidopsis thaliana</i> | | | | | | | | | | |
| (ref. <i>A. lyrata</i>) | | | | | | | | | | |
| Gene Sn | 44.0 | 61.3 | 66.5 | 28.9 | 35.9 | 64.4 | 51.3 | NA | 56.9 | 57.9 |
| Gene Sp | 47.8 | 65.7 | 71.3 | 47.9 | 59.1 | 52.0 | 52.5 | NA | 65.7 | 67.8 |
| Transcript Sn | 37.5 | 52.2 | 57.2 | 26.6 | 33.7 | 55.0 | 43.5 | NA | 48.3 | 49.1 |
| Transcript Sp | 47.8 | 65.7 | 65.3 | 35.6 | 48.3 | 50.9 | 52.5 | NA | 65.7 | 67.8 |
| Exon Sn | 70.0 | 79.3 | 80.6 | 58.1 | 60.8 | 82.9 | 76.1 | NA | 81.8 | 82.1 |
| Exon Sp | 81.9 | 86.6 | 87.5 | 81.9 | 87.1 | 79.0 | 76.1 | NA | 87.5 | 88.6 |
| <i>Caenorhabditis elegans</i> | | | | | | | | | | |
| (ref. <i>C. briggsae</i>) | | | | | | | | | | |
| Gene Sn | 26.2 | 39.6 | 49.1 | 18.7 | 22.6 | 55.0 | 41.0 | NA | 40.5 | 47.3 |
| Gene Sp | 38.0 | 49.9 | 63.8 | 29.1 | 36.1 | 55.2 | 30.8 | NA | 51.5 | 56.4 |
| Transcript Sn | 21.0 | 30.7 | 39.8 | 16.2 | 20.0 | 43.0 | 31.3 | NA | 31.4 | 36.2 |
| Transcript Sp | 38.0 | 49.9 | 58.7 | 24.1 | 30.1 | 53.2 | 30.8 | NA | 51.5 | 56.4 |
| Exon Sn | 50.3 | 64.2 | 67.1 | 54.4 | 59.1 | 80.2 | 69.4 | NA | 70.5 | 75.2 |
| Exon Sp | 82.6 | 81.5 | 87.5 | 81.3 | 84.1 | 85.3 | 62.3 | NA | 85.6 | 86.7 |
| <i>Drosophila melanogaster</i> | | | | | | | | | | |
| (ref. <i>D. simulans</i>) | | | | | | | | | | |
| Gene Sn | 64.3 | 78.2 | 83.1 | 55.7 | 55.2 | 64.9 | 55.2 | NA | 61.5 | 64.0 |
| Gene Sp | 69.2 | 81.6 | 87.1 | 71.3 | 73.5 | 59.4 | 46.3 | NA | 69.6 | 71.9 |
| Transcript Sn | 44.1 | 52.9 | 65.0 | 48.7 | 49.0 | 46.1 | 38.5 | NA | 42.7 | 44.3 |
| Transcript Sp | 69.2 | 81.6 | 81.2 | 60.1 | 65.7 | 57.9 | 46.3 | NA | 69.6 | 71.9 |
| Exon Sn | 69.0 | 76.3 | 80.0 | 67.8 | 66.2 | 75.0 | 66.5 | NA | 74.3 | 76.3 |
| Exon Sp | 89.1 | 92.0 | 93.3 | 85.4 | 88.3 | 81.7 | 66.9 | NA | 88.0 | 89.1 |
| <i>Schizosaccharomyces pombe</i> | | | | | | | | | | |
| (ref. <i>S. octosporus</i>) | | | | | | | | | | |
| Gene Sn | 49.2 | 76.4 | 79.2 | 69.0 | 65.8 | 77.4 | 42.8 | 79.7 | 71.6 | 74.6 |
| Gene Sp | 59.9 | 84.6 | 88.0 | 93.8 | 92.5 | 80.5 | 68.7 | 72.6 | 88.1 | 89.1 |
| Transcript Sn | 49.2 | 76.4 | 79.2 | 69.0 | 65.8 | 77.4 | 42.8 | 79.7 | 71.6 | 74.6 |
| Transcript Sp | 59.9 | 84.6 | 87.6 | 80.5 | 71.3 | 76.5 | 68.7 | 72.6 | 88.1 | 89.1 |
| Exon Sn | 56.1 | 81.6 | 83.1 | 77.2 | 77.7 | 83.2 | 50.1 | 79.6 | 79.2 | 81.2 |
| Exon Sp | 73.3 | 88.6 | 91.9 | 87.6 | 81.7 | 83.2 | 71.4 | 81.7 | 92.0 | 92.6 |

The target species are given in multi-column rows. The same reference species, which is given in brackets, is used for all tools using homology-based gene prediction indicated by plus. The asterisks indicates that the performance of BRAKER1, MAKER2 and CodingQuarry is given as reported in [11]. The highest value per line is depicted in bold-face

Second, we additionally consider RNA-seq data. MAKER2 does not allow for combining RNA-seq evidence and homology-based predictions without using any *ab-initio* gene predictor. In contrast, GeMoMa allows for additionally using intron position conservation and RNA-seq data. For this reason, we compare the performance of GeMoMa with and without RNA-seq evidence (Table 1). We find that sensitivity and specificity in almost all cases increases by up to 13.9 with only two exceptions for transcript specificity of *A. thaliana* and *D. melanogaster* which decreases by at most 0.4. Hence, we summarize that RNA-seq evidence improves the sensitivity and specificity of GeMoMa and should be used if available.

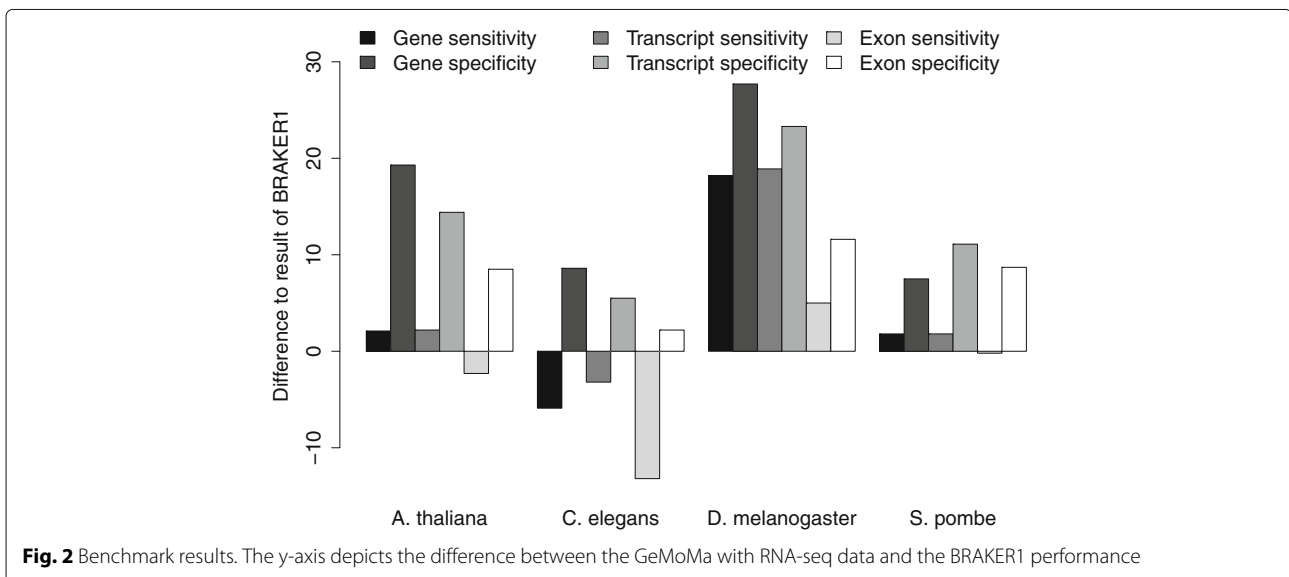
Third, we compare the performance of GeMoMa using RNA-seq evidence to that of purely RNA-seq-based pipelines, namely Cufflinks and StringTie (Table 1). We find for all four species that GeMoMa using RNA-seq evidence outperforms purely RNA-seq-based pipelines. Interestingly, purely RNA-seq-based pipelines also yield the worst gene/transcript sensitivity and specificity for *C. elegans*. Comparing the results based on different transcript assemblers, we find that the results based on StringTie are better than those based on Cufflinks for *A. thaliana* and *C. elegans*, while the opposite is true for *S. pombe*. For *D. melanogaster*, both pipelines perform comparably. Additional RNA-seq reads increasing the coverage might improve the performance of purely RNA-seq-based pipelines but could also improve the performance of GeMoMa.

Summarizing these three observations, we find that GeMoMa performs better than purely homology-based or purely RNA-seq-based pipelines and that including RNA-seq data improves the performance of GeMoMa.

Hence, we compare GeMoMa to combined gene prediction approaches. Specifically, we compare the performance of GeMoMa using RNA-seq evidence to BRAKER1 in Fig. 2, which provides the best overall performance in [11]. We find that GeMoMa performs better than BRAKER1 for the categories gene and transcript with the exception of gene and transcript sensitivity for *C. elegans*. Interestingly, we find the biggest improvements for *D. melanogaster* where gene/transcript sensitivity and specificity increases between 18.2 and 27.7. For the exon category, we find a less clear picture. In total, we observe the worst results for *C. elegans* where the sensitivity for all three categories decreases between 3.2 and 13.2, while the specificity increases only between 2.2 and 8.6. Notably, we generally find the worst gene/transcript sensitivity and specificity for *C. elegans* compared with the other target species considering the best performance of all tools.

In summary, we find that the gene predictors MAKER2, BRAKER1, CodingQuarry and GeMoMa, and the transcript assemblers Cufflinks and StringTie often perform quite well on exon level. The main difference becomes evident on transcript and gene level, where exons need to be combined correctly (Table 1) as reported earlier [29, 30]. Homology-based gene predictors might benefit from experimentally validated and manually curated reference transcripts guiding the prediction of transcripts in the target organism.

Although GeMoMa performed well, it is not able to predict genes that do not show any homology to a protein in the reference species, while *ab-initio* gene predictors might fail in other cases. As both types of approaches have their specific advantages, users will probably use combinations of different gene predictors in practice to obtain a comprehensive gene annotation.



In addition, we performed a small runtime study for the two main time-consuming steps of the pipeline to demonstrate that GeMoMa is reasonably fast (Additional file 1: Table S7).

Combined gene prediction pipelines

Combined gene prediction pipelines, as for instance MAKER2, use RNA-seq evidence, homology-based and *ab-initio* methods for predicting final gene models. MAKER2 uses exonerate by default for homology-based gene prediction. However, MAKER2 also provides the possibility to use other homology-based gene predictors instead of exonerate (cf. parameter protein_gff). For this reason, we compare the performance of MAKER2 using either exonerate or GeMoMa for homology based gene prediction (Table 1). In addition, we use Augustus as *ab-initio* gene prediction program and Trinity transcripts in MAKER2. We find that MAKER2 using GeMoMa performs better than MAKER2 using exonerate for all species and all measure. The improvement varies between 0.3% and 6.8% with clearly the biggest improvement for *C. elegans*.

In addition, we find that the MAKER2 performance is substantially improved compared to the performance of the the previously reported MAKER2 predictions, either purely based on proteins (cf. Table 1, column MAKER2⁺ (exonerate)) or as reported in [11] (cf. Maker2*). These other predictions do not utilize all available sources of information as they either ignore RNA-seq data and *ab-initio* gene prediction or homology to proteins of related species. Based on this observation, we agree that combined gene prediction pipelines benefit from

the inclusion of all available evidence and that performance is decreased if some important evidence is missed [9].

Furthermore, we compare GeMoMa using RNA-seq evidence with MAKER2 using RNA-seq evidence, homology-based and *ab-initio* gene prediction. In some cases, it is hard to compare these results as sensitivity of one tool is higher than the sensitivity of the other tool and the opposite is true for specificity. In machine learning, recall, also known as sensitivity, and precision, which is called specificity in the context of gene prediction evaluation [31], are combined into a single scalar value called F1 measure [32] that can be compared more easily. We combined sensitivity and specificity resulting in an F1 measure for each evaluation level gene, transcript and exon (Additional file 1 – Table S4) We find that in many cases GeMoMa using RNA-seq evidence outperforms MAKER2. The reason for this observation might be that RNA-seq data and homology based gene prediction is used in MAKER2 to train *ab-initio* gene predictors, in this case Augustus. With the recommended parameter setting, homology-based gene predictions are not directly used for the final prediction and doing so might further improve performance.

Influence of reference species

Utilizing different fly species from FlyBase [33], we scrutinize the influence of different or multiple reference species on the performance of GeMoMa using RNA-seq data (Additional file 1: Table S8). In Fig. 3, we depict gene sensitivity and gene specificity for eight different reference species indicated by points. We find that performance

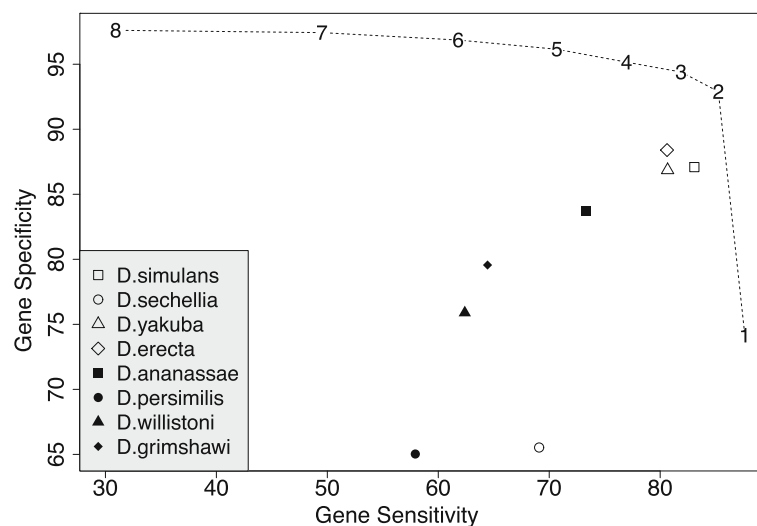


Fig. 3 Gene sensitivity and specificity for *D. melanogaster* using different or multiple reference species in GeMoMa. The points correspond to the eight reference species. In addition, the dashed line indicates the usage of multiple reference species. Using multiple reference species allows for filtering identical predictions from several reference as indicated by the numbers

varies with the reference species. In this specific case, *D. sechellia* and *D. persimilis* yield the worst results for single reference-based predictions. This observation might be related to the fact that genome assembly of *D. sechellia* and *D. persimilis* is of lower quality [34], while the genome of *D. simulans* has been updated [35] later. Besides these two outliers, the performance of the different fly species as reference species for *D. melanogaster* in GeMoMa correlates with their evolutionary distance [36]. Generally speaking, the closer a reference species is related to the target species *D. melanogaster*, the better is the performance in terms of gene sensitivity and specificity. Hence, we speculate that two requirements must be met to have a good reference species. First, the evolutionary distance between reference and target species should be small and second, the genome assembly and annotation of the reference species should be comprehensive and of high quality.

The new GAF module of GeMoMa allows for combining the predictions based on different reference organisms. The combined predictions may be filtered by number of reference species with perfect support (#evidence), as indicated by the dashed line. We find that combining multiple reference organisms improves prediction performance and stability. Depending on the number of supporting reference organisms required, gene specificity and gene sensitivity may be balanced according to the needs of a specific application. We observe that (i) gene sensitivity increases but specificity decreases when requiring support from at least one reference organism, whereas (ii) gene specificity increases but sensitivity decreases severely filtering for perfect support from all eight reference species. In summary, the inclusion of multiple reference species may yield an improved prediction performance for GeMoMa using the GAF module, where we suggest to filter predictions for support by at least two but not necessarily all reference species.

Furthermore, we check whether GeMoMa allows for identifying new transcripts in *D. melanogaster* that do not overlap with any annotated transcript but are supported by RNA-seq data. First, we check whether we could identify transcripts based on the GeMoMa predictions using *D. simulans* as reference organism. We find 35 multi-coding-exon predictions that do not overlap with any annotated transcript but have a tie of 1, i.e., all introns are supported by split reads in the RNA-seq data (see “Methods”). In addition, we find 15 single-coding-exon predictions that do not overlap with any annotated transcript but have a tpc of 1, i.e., that are fully covered by mapped RNA-seq reads. Second, we check whether we could identify transcripts that are supported by at least two of the eight reference species (cf. above). We find 14 multi-coding-exon predictions that do not overlap with any annotated transcript, obtain a tie of 1 and are

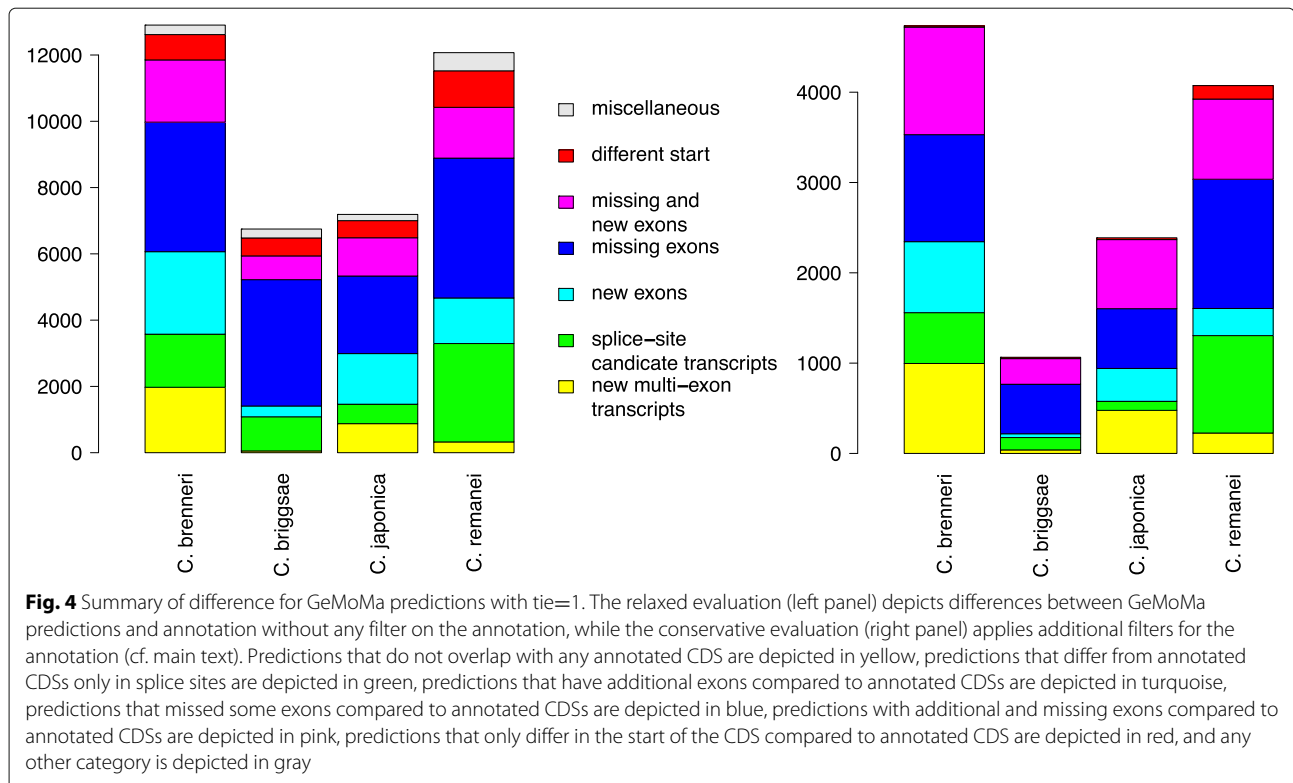
supported by at least two of the eight reference species. In addition, we find 9 single-coding-exon predictions that do not overlap with any annotated transcript, have a tpc of 1 and are supported by at least two of the five reference species. In summary, those genes supported by multiple reference organisms or additional RNA-seq data might be promising candidates for extending the existing genome annotation of *D. melanogaster*.

Analysis of nematode species

The relatively poor results for *C. elegans* in the benchmark study, might be due to insufficiencies in the current *C. briggsae* annotation. Hence, we decided to scrutinize the Wormbase annotation of four nematode species comprising *C. brenneri*, *C. briggsae*, *C. japonica*, and *C. remanei* based on the model organism *C. elegans*. We compare GeMoMa predictions with manually curated CDS from Wormbase. Based on RNA-seq evidence, we collect multi-coding-exon predictions of GeMoMa with tie=1 and compare these to the annotation as depicted in Fig. 4.

In summary, we find between 6749 differences for *C. briggsae* and 12903 for *C. brenneri* (cf. Fig. 4). The most interesting category are new multi-coding-exon predictions, which vary between 53 for *C. briggsae* and 1974 for *C. brenneri*. The largest category are GeMoMa predictions that missed exons compared to annotated CDSs, which vary between 2340 for *C. japonica* and 4220 for *C. remanei*.

We additionally filter the transcripts showing differences to obtain a smaller, more conservative set of high-confidence predictions. First, we filter new multi-coding exon GeMoMa predictions for tpc=1 obtaining between 39 and 996 for *C. briggsae* and *C. brenneri*, respectively. Second, we filter GeMoMa predictions that have different splice sites compared to highly overlapping annotated transcripts, contain new exons, have missing exons, or have new and missing exons for tie<1 of the overlapping annotation. We obtain between 100 and 1079 predictions with different splice-site, between 42 and 786 predictions containing new exons, between 548 and 1431 predictions with missing exons, and between 284 and 1191 predictions with new and missing exons. Finally, for GeMoMa predictions that differ in the start codon compared to the annotation, we filter for tpc=1 of the GeMoMa prediction and tpc<1 for the annotation obtaining between 14 and 149 for *C. brenneri* and *C. remanei*, respectively. In summary, we obtain between 1065 predictions differing from the annotation for *C. briggsae* and 4735 predictions for *C. brenneri*, respectively (cf. Fig. 4) using these strict criteria. Despite the overall reduction of transcripts considered, GeMoMa predictions that missed exons compared to annotated CDSs are the largest category for all four nematode species.



For both evaluations, we find that the predictions for *C. briggsae* are in better accordance with the annotation than the predictions of the remaining three nematode species. One possible explanation might be that the annotation of *C. briggsae* has recently been updated using RNA-seq data (Gary Williams, personal communication), while the annotation of *C. japonica* is based on Augustus (Erich Schwartz, personal communication) and the annotation of the other two nematodes are NGASP sets from multiple *ab-initio* gene prediction programs [37]. For *C. japonica*, we find the second best results, although *C. japonica* is phylogenetically more distantly related to *C. elegans* than the remaining two nematodes [38]. This is additional evidence that the annotation pipeline employed has a decisive influence on the quality and completeness of the annotation.

In addition, we checked for *C. brenneri* whether the GeMoMa predictions partially overlap with cDNAs or ESTs mapped to the *C. brenneri* genome. In 472 cases, the prediction overlaps with a cDNA or EST, but not with the annotation. In 364 out of these 472 cases, the prediction has tie=1. To evaluate the predictions, we manually checked about 9% (43) of the predicted missing genes with tie=1. Based on RNA-seq data, protein homology, cDNA/ESTs and manual curation, 95% were genuine new isoforms which have been missed in the original *C. brenneri* gene set. This shows that GeMoMa is valuable

in finding isoforms missed by traditional prediction methods.

Analysis of barley

Complementary to the studies in animals in the last subsection, we used GeMoMa to predict the annotation of protein-coding genes in barley (*Hordeum vulgare*). Based on the benchmark results for *D. melanogaster*, we used several reference organisms to predict the gene annotation using GeMoMa and GAF and finally obtain 75 484 transcript predictions. Most of the predictions showed a good overlap with the annotation ($F_1 \geq 0.8$). Nevertheless, 27 204 out of these 75 484 predictions had little ($F_1 < 0.8$) or no overlap with high or low confidence gene annotations. However, thousands of the transcripts contained in the official annotation do not have start or stop codons [14], which renders an exact comparison of predictions with perfect or at least very good overlap unreasonable.

Hence, we focus on 19 619 predictions with no overlap with any annotated transcript (Table 2). Scrutinizing these predictions, we find 1 729 single-coding-exon predictions that are completely covered by RNA-seq reads (tpc=1) but that are not contained in the annotation. Out of these, 367 are partially supported by best BLAT matches of ESTs to the genome. In addition, we analyzed multi-coding-exon predictions and find 2 821 predictions that obtain tie=1,

Table 2 Predictions that do not overlap with any high or low confidence annotation

| a) Single-coding-exon predictions | | | |
|-----------------------------------|--------------|-------------|---------------|
| #evidence | tpc = 0 | 0 < tpc < 1 | tpc = 1 |
| 1 | 1 971 (11) | 878 (14) | 1 005 (137) |
| 2 | 204 (19) | 158 (8) | 299 (55) |
| 3 | 200 (16) | 126 (5) | 257 (92) |
| 4 | 91 (17) | 43 (9) | 168 (83) |
| Σ | 2 466 (63) | 1 205 (36) | 1 729 (367) |
| b) Multi-coding-exon predictions | | | |
| #evidence | tie = 0 | 0 < tie < 1 | tie = 1 |
| 1 | 9 671 (287) | 942 (211) | 1 681 (775) |
| 2 | 283 (36) | 86 (32) | 456 (196) |
| 3 | 155 (31) | 64 (43) | 382 (223) |
| 4 | 142 (57) | 55 (37) | 302 (196) |
| Σ | 10 251 (411) | 1 147 (323) | 2 821 (1 390) |

The numbers in parenthesis depict those predictions that are partially supported by any best BLAT hit of ESTs

stating that each predicted intron is supported by at least one split read from mapped RNA-seq data. Out of these, 1 390 are partially supported by best BLAT matches of ESTs to the genome.

Besides predictions that are well supported by RNA-seq data, we also observe thousands of predictions that are not ($tpc = 0$ or $tie = 0$) or only partially ($0 < tpc < 1$ or $0 < tie < 1$) supported by RNA-seq. Despite no or only partial RNA-seq support, we find that 833 are partially supported by best BLAT matches of ESTs to the genome.

Alternatively, we can utilize the number of reference organisms that support a prediction (#evidence) to filter the predictions as noted for *D. melanogaster*. This approach will decrease sensitivity, but increase specificity obtaining predictions with a high confidence. Although, we find the most predictions with #evidence = 1, we also find about 3 500 predictions with #evidence > 1, more than 1 100 of these predictions are additionally supported by RNA-seq data or ESTs.

Conclusions

Summarizing the methods and results, we present an extension of GeMoMa that allows for the incorporation of RNA-seq data into homology-based gene prediction utilizing intron position conservation. Comparing the performance of GeMoMa with and without RNA-seq evidence, we demonstrate for all four organism included in the benchmark that RNA-seq evidence improves the performance of GeMoMa. GeMoMa performs equally well or better than BRAKER1, MAKER2, CodingQuarry, and purely RNA-seq-based pipelines on the benchmark data sets including plants, animals and fungi.

We also find that the performance depends on the evolutionary distance between reference and target organism. However, prediction performance also depends on several further aspects including i) the quality of the target genome (assembly), ii) the number of reference organisms available and iii) especially the quality of the reference annotation(s) itself. Hence, we recommend to balance between evolutionary distance and (expected) quality of the reference annotation when selecting reference species for GeMoMa.

The integration of RNA-seq data into GeMoMa might help to overcome wrongly annotated splice sites in the reference species in some cases. However, missing or wrongly additional annotated exons in the reference annotation might still lead to partially wrong gene model predictions in the target species. The benefit of RNA-seq data, however, also depends on the quality and amount of sequenced reads, on the diversity (tissues, conditions) of the sequenced samples, and on the library type, where stranded libraries should be more informative than non-stranded ones. In addition, GeMoMa uses RNA-seq data currently only to refine homologous genes models and not to identify transcribed gene models that do not show any homology. Hence, GeMoMa should be used in combination with other gene predictors allowing for purely RNA-seq-based or *ab-initio* gene predictions. Exemplarily, we demonstrate that GeMoMa helps to improve the performance of combined gene predictor pipelines as for instance MAKER2.

Notably, model organisms have been used as target organisms in this benchmark, whereas they would typically be used as reference organisms in real applications. Hence, the performance of homology-based gene prediction programs might be underestimated. In summary, we recommend to use homology-based gene prediction using RNA-seq data as implemented in GeMoMa whenever high-quality gene annotations of related species are available.

Interestingly, we find that GeMoMa works especially well for *D. melanogaster* in the benchmark study compared to the performance of its competitors. One possible reason could be that Flybase used homology and RNA-seq data besides other evidence to infer the gene annotation [19]. In contrast, we find the worst results in *C. elegans* in the benchmark study, which might be related to the fact that the *C. elegans* gene set contains many rare isoform community submissions whereas *C. briggsae* was annotated by a large scale gene predictions effort based on RNA-seq.

Scrutinizing the annotation in Wormbase, we predicted protein-coding transcripts for four nematode species based on the annotation of the model organism *C. elegans*. We find that a substantial part of the GeMoMa predictions is either missing, marked as modification

of annotated transcripts or alternative transcripts. Especially for the three nematodes, *C. brenneri*, *C. japonica* and *C. remanei*, that are annotated solely using *ab-initio* gene prediction, we find a large part of the annotation that is marked as questionable or missing. This may give an indication, why homology-based gene prediction for *C. elegans* shows less good performance in the benchmark study. The GeMoMa predictions of the four nematodes will be included in Wormbase in the upcoming releases. Furthermore, GeMoMa will be included in the WormBase gene curation process and trialled for strain annotation.

Furthermore, we predicted protein-coding transcripts for barley using four reference species and find several hundreds of predictions that are not included in the reference annotation but are supported by RNA-seq data, ESTs or multiple reference species. Hence, we conclude that these are valuable predictions harboring additional barley genes. These predictions will be incorporated in the new barley annotation.

GeMoMa provides a user-friendly documentation and can be used as command line tool or through the Galaxy workflow management system [39] providing its own Galaxy integration (Fig. 1). GeMoMa is freely available under GNU GPL3 at <http://www.jstacs.de/index.php/GeMoMa>.

Additional files

Additional file 1: Supplementary Tables and Figures. **Table S1:** Data used for the BRAKER1 benchmark. **Table S2:** Data for Wormbase study. **Table S3:** Data for barley annotation. **Table S4:** F1-measure on the BRAKER1 test sets. **Table S5:** CodingQuarry and GeMoMa results for *S. pombe* using the original read mappings from [10, 11]. **Table S6:** CodingQuarry and GeMoMa results for *S. cerevisiae*. **Table S7:** GeMoMa runtime. **Table S8:** GeMoMa performance for *D. melanogaster*. (PDF 147 kb)

Abbreviations

avgCov: Average coverage; cDNA: Complementary DNA; CDS: Coding sequence; ERE: Extract RNA-seq evidence; EST: Expressed sequence tag; GAF: GeMoMa annotation filter; GeMoMa: Gene Model Mapper; iAA: Identical amino acids; tae: Transcript acceptor evidence; tde: Transcript donor evidence; tie: Transcript intron evidence; tpc: Transcript percentage coverage; minCov: Minimal coverage; minSplitReads: Minimum number of split reads; pAA: Positive-scoring amino acids

Acknowledgements

We thank Katharina Hoff for providing the BRAKER1 benchmark data sets, Carson Holt for assisting the MAKER2 comparison, Gil dos Santos for his comments on the quality of the *Drosophila* genome assemblies, Erich Schwartz for his comments on *C. japonica*, Alison Testa for her help with CodingQuarry, Gary Williams for his comments on *C. briggsae*, and Thomas Berner for technical assistance. We thank the open access publication fund of Martin Luther University Halle–Wittenberg for funding article-processing charges.

Availability of data and materials

Not applicable. Accession numbers of publicly available data used in this study are listed in Additional file 1.

Authors' contributions

JK and JG implemented the software, designed and performed benchmark studies. JK, FH, MP, SOT and JG contributed to data analysis and interpretation, and to writing the manuscript. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Institute for Biosafety in Plant Biotechnology, Julius Kühn-Institut (JKI) - Federal Research Centre for Cultivated Plants, D-06484, Quedlinburg, Germany. ²European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SD, Cambridge, UK. ³Plant Genome and Systems Biology, Helmholtz Center Munich - German Research Center for Environmental Health, D-85764, Neuherberg, Germany. ⁴Institute of Computer Science, Martin Luther University Halle–Wittenberg, D-06120, Halle (Saale), Germany.

Received: 20 November 2017 Accepted: 14 May 2018

Published online: 30 May 2018

References

1. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL, Wold BJ, Pachter L. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol.* 2010;28(5):511–5. <https://doi.org/10.1038/nbt.1621>.
2. Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotech.* 2015;33(3):290–5. <https://doi.org/10.1038/nbt.3122>.
3. Solovyev V, Kosarev P, Seledsov I, Vorobyev D. Automatic annotation of eukaryotic genes, pseudogenes and promoters. *Genome Biol.* 2006;7(1):10. <https://doi.org/10.1186/gb-2006-7-s1-s10>.
4. Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics.* 2008;24(5):637. <https://doi.org/10.1093/bioinformatics/btn013>.
5. Slater G, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics.* 2005;6(1):31. <https://doi.org/10.1186/1471-2105-6-31>.
6. She R, Chu JS-C, Uyar B, Wang J, Wang K, Chen N. genBlastG: using BLAST searches to build homologous gene models. *Bioinformatics.* 2011;27(15):2141–3. <https://doi.org/10.1093/bioinformatics/btr342>. <http://bioinformatics.oxfordjournals.org/content/27/15/2141.full.pdf+html>.
7. Keilwagen J, Wenk M, Erickson JL, Schattat MH, Grau J, Hartung F. Using intron position conservation for homology-based gene prediction. *Nucleic Acids Res.* 2016;44(9):89. <https://doi.org/10.1093/nar/gkw092>.
8. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet.* 2009;10(1):57–63. <https://doi.org/10.1038/nrg2484>.
9. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-database management tool for second-generation genome projects. *BMC Bioinformatics.* 2011;12(1):491. <https://doi.org/10.1186/1471-2105-12-491>.
10. Testa AC, Hane JK, Ellwood SR, Oliver RP. CodingQuarry: highly accurate hidden Markov model gene prediction in fungal genomes using RNA-seq transcripts. *BMC Genomics.* 2015;16(1):170. <https://doi.org/10.1186/s12864-015-1344-4>.
11. Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: Unsupervised RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics.* 2016;32(5):767. <https://doi.org/10.1093/bioinformatics/btv661>.
12. Lomsadze A, Burns PD, Borodovsky M. Integration of mapped RNA-seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.* 2014;42(15):119. <https://doi.org/10.1093/nar/gku557>.
13. Howe KL, Bolt BJ, Cain S, Chan J, Chen WJ, Davis P, Done J, Down T, Gao S, Grove C, Harris TW, Kishore R, Lee R, Lomax J, Li Y, Muller H-M, Nakamura C, Nuin P, Paulini M, Raciti D, Schindelman G, Stanley E,

- Tuli MA, Van Auken K, Wang D, Wang X, Williams G, Wright A, Yook K, Berriman M, Kersey P, Schedl T, Stein L, Sternberg PW. Wormbase 2016: expanding to enable helminth genomic research. *Nucleic Acids Res.* 2016;44(D1):774. <https://doi.org/10.1093/nar/gkv1217>.
14. Mascher M, Gundlach H, Himmelbach A, Beier S, Twardziok SO, Wicker T, Radchuk V, Dockter C, Hedley PE, Russell J, Bayer M, Ramsay L, Liu H, Haberer G, Zhang X-Q, Zhang Q, Barrero RA, Li L, Taudien S, Groth M, Felder M, Hastie A, Šimková H, Staňková H, Vrána J, Chan S, Muñoz-Amatrián M, Ounit R, Wanamaker S, Bolser D, Colmsee C, Schmutzer T, Aliyeva-Schnorr L, Grasso S, Tanskanen J, Chailyan A, Sampath D, Heavens D, Clissold L, Cao S, Chapman B, Dai F, Han Y, Li H, Li X, Lin C, McCooke JK, Tan C, Wang P, Wang S, Yin S, Zhou G, Poland JA, Bellgard MI, Borisjuk L, Houben A, Doležel J, Ayling S, Lonardi S, Kersey P, Langridge P, Muehlbauer GJ, Clark MD, Caccamo M, Schulman AH, Mayer KFX, Platzer M, Close TJ, Scholz U, Hansson M, Zhang G, Braumann I, Spannagl M, Li C, Waugh R, Stein N. A chromosome conformation capture ordered sequence of the barley genome. *Nature.* 2017;544(7651):427–33. <https://doi.org/10.1038/nature22043>.
 15. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol.* 1990;215(3):403–10. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2).
 16. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D, Li B, Lieber M, MacManes MD, Ott M, Orvis J, Pochet N, Strozzi F, Weeks N, Westerman R, William T, Dewey CN, Henschel R, LeDuc RD, Friedman R, Regev A. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protocols.* 2013;8(8):1494–512. <https://doi.org/10.1038/nprot.2013.084>.
 17. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.* 2013;14(4):R36. <https://doi.org/10.1186/gb-2013-14-4-r36>.
 18. Rawat V, Abdelsamad A, Pietzenek B, Seymour DK, Koenig D, Weigel D, Pecinka A, Schneeberger K. Improving the annotation of arabidopsis lyrata using rna-seq data. *PLOS ONE.* 2015;10(9):1–12. <https://doi.org/10.1371/journal.pone.0137391>.
 19. Matthews BB, dos Santos G, Crosby MA, Emmert DB, St Pierre SE, Gramates LS, Zhou P, Schroeder AJ, Falls K, Strelets V, Russo SM, Gelbart WM. The FlyBase Consortium. Gene model annotations for drosophila melanogaster: Impact of high-throughput data. *G3: Genes Genomes Genet.* 2015;5(8):1721–36. <https://doi.org/10.1534/g3.115.018929>. <http://www.g3journal.org/content/5/8/1721.full.pdf>.
 20. Rhind N, Chen Z, Yassour M, Thompson DA, Haas BJ, Habib N, Wapinski I, Roy S, Lin MF, Heiman DI, Young SK, Furuya K, Guo Y, Pidoux A, Chen HM, Robbertse B, Goldberg JM, Aoki K, Bayne EH, Berlin AM, Desjardins CA, Dobbs E, Dukaj L, Fan L, FitzGerald MG, French C, Gujja S, Hansen K, Keifenheim D, Levin JZ, Mosher RA, Müller CA, Pfiffner J, Priest M, Russ C, Smialowska A, Swoboda P, Sykes SM, Vaughn M, Vengrova S, Yoder R, Zeng Q, Allshire R, Baulcombe D, Birren BW, Brown W, Ekwall K, Kellis M, Leatherwood J, Levin H, Margalit H, Martienssen R, Nieduszynski CA, Spatafora JW, Friedman N, Dalgaard JZ, Baumann P, Niki H, Regev A, Nusbaum C. Comparative functional genomics of the fission yeasts. *Science.* 2011;332(6032):930–6. <https://doi.org/10.1126/science.1203357>. <http://science.sciencemag.org/content/332/6032/930.full.pdf>.
 21. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29(1):15. <https://doi.org/10.1093/bioinformatics/bts635>.
 22. Lamesch P, Berardini TZ, Li D, Swarbreck D, Wilks C, Sasidharan R, Muller R, Dreher K, Alexander DL, Garcia-Hernandez M, Karthikeyan AS, Lee CH, Nelson WD, Ploetz L, Singh S, Wensel A, Huala E. The Arabidopsis Information Resource (TAIR): improved gene annotation and new tools. *Nucleic Acids Res.* 2012;40(D1):1202. <https://doi.org/10.1093/nar/gkr1090>.
 23. International Brachypodium Initiative. Genome sequencing and analysis of the model grass *Brachypodium distachyon*. *Nature.* 2010;463(5):763–8. <https://doi.org/10.1038/nature08747>.
 24. Ouyang S, Zhu W, Hamilton J, Lin H, Campbell M, Childs K, Thibaud-Nissen F, Malek R, Lee Y, Zheng L, Orvis J, Haas B, Wortman J, Buell CR. The tigr rice genome annotation resource: improvements and new features. *Nucleic Acids Res.* 2007;35(suppl_1):883. <https://doi.org/10.1093/nar/gkl976>.
 25. Bennetzen JL, Schmutz J, Wang H, Percifield R, Hawkins J, Pontaroli AC, Estep M, Feng L, Vaughn JN, Grimwood J, Jenkins J, Barry K, Lindquist E, Hellsten U, Deshpande S, Wang X, Wu X, Mitros T, Triplett J, Yang X, Ye C-Y, Mauro-Herrera M, Wang L, Li P, Sharma M, Sharma R, Ronald PC, Panaud O, Kellogg EA, Brutnell TP, Doust AN, Tuskan GA, Rokhsar D, Devos KM. Reference genome sequence of the model plant *Setaria*. *Nat Biotechnol.* 2012;30(6):555–61. <https://doi.org/10.1038/nbt.2196>.
 26. Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat Protocols.* 2016;11(9):1650–67. <https://doi.org/10.1038/nprot.2016.095>.
 27. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res.* 2002;12(4):656–64. <https://doi.org/10.1101/gr.229202>. Article published online before March 2002.
 28. Keibler E, Brent MR. Eval: A software package for analysis of genome annotations. *BMC Bioinformatics.* 2003;4(1):50. <https://doi.org/10.1186/1471-2105-4-50>.
 29. Steijger T, Abril JF, Engström PG, Kokocinski F, Hubbard TJ, Guigó R, Harrow J, Bertone P, Consortium R, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods.* 2013;10(12):1177–84.
 30. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szczesniak MW, Gaffney DJ, Elo LL, Zhang X, Mortazavi A. A survey of best practices for RNA-seq data analysis. *Genome Biol.* 2016;17(1):13. <https://doi.org/10.1186/s13059-016-0881-8>.
 31. Burset M, Guigó R. Evaluation of gene structure prediction programs. *Genomics.* 1996;34(3):353–67. <https://doi.org/10.1006/geno.1996.0298>.
 32. Powers DMW. Evaluation: From precision, recall and F-measure to ROC, informedness, markedness & correlation. *J Mach Learn Technol.* 2011;2(1):37–63.
 33. Gramates LS, Marygold SJ, Santos Gd, Urbano J-M, Antonazzo G, Matthews BB, Rey AJ, Tabone CJ, Crosby MA, Emmert DB, Falls K, Goodman JL, Hu Y, Ponting L, Schroeder AJ, Strelets VB, Thurmond J, Zhou P. FlyBase at 25: looking to the future. *Nucleic Acids Res.* 2017;45(D1):663–71. <https://doi.org/10.1093/nar/gkw1016>.
 34. Clark AG, Eisen MB, Smith DR, Bergman CM, Oliver B, Markow TA, Kaufman TC, Kellis M, Gelbart W, Iyer VN, et al. Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature.* 2007;450(7167):203–18.
 35. Hu TT, Eisen MB, Thornton KR, Andolfatto P. A second-generation assembly of the *Drosophila simulans* genome provides new insights into patterns of lineage-specific divergence. *Genome Res.* 2013;23(1):89–98. <https://doi.org/10.1101/gr.141689.112>. <http://genome.cshlp.org/content/23/1/89.full.pdf+html>.
 36. Singh ND, Larracuent AM, Sackton TB, Clark AG. Comparative genomics on the *drosophila* phylogenetic tree. *Annu Rev Ecol Evol Syst.* 2009;40(1):459–80. <https://doi.org/10.1146/annurev.ecolsys.110308.120214>.
 37. Coghlan A, Fiedler TJ, McKay SJ, Flicek P, Harris TW, Blasiar D, nGASP Consortium, Stein LD. ngasp—the nematode genome annotation assessment project. *BMC Bioinformatics.* 2008;9:549. <https://doi.org/10.1186/1471-2105-9-549>.
 38. Kiontke KC, Félix M-A, Ailion M, Rockman MV, Braendle C, Pénigault J-B, Fitch DH. A phylogeny and molecular barcodes for caenorhabditis, with numerous new species from rotting fruits. *BMC Evol Biol.* 2011;11(1):339. <https://doi.org/10.1186/1471-2148-11-339>.
 39. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Cech M, Chilton J, Clements D, Coraor N, Eberhard C, Grünig B, Guerler A, Hillman-Jackson J, Von Kuster G, Rasche E, Soranzo N, Turaga N, Taylor J, Nekrutenko A, Goecks J. The galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids Res.* 2016;44(W1):3. <https://doi.org/10.1093/nar/gkw343>.