

GeneMark-EP+: eukaryotic gene prediction with self-training in the space of genes and proteins

Tomáš Brůna^{1,†}, Alexandre Lomsadze^{2,†} and Mark Borodovsky^{1,2,3,*}

¹School of Biological Sciences, Georgia Institute of Technology, Atlanta, GA 30332, USA, ²Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA and ³School of Computational Science and Engineering, Georgia Institute of Technology, Atlanta, GA 30332, USA

Received December 13, 2019; Revised March 10, 2020; Editorial Decision April 01, 2020; Accepted May 12, 2020

ABSTRACT

We have made several steps toward creating a fast and accurate algorithm for gene prediction in eukaryotic genomes. First, we introduced an automated method for efficient *ab initio* gene finding, GeneMark-ES, with parameters trained in iterative *unsupervised* mode. Next, in GeneMark-ET we proposed a method of integration of unsupervised training with information on intron positions revealed by mapping short RNA reads. Now we describe GeneMark-EP, a tool that utilizes another source of external information, a protein database, readily available prior to the start of a sequencing project. A new specialized pipeline, ProtHint, initiates massive protein mapping to genome and extracts hints to splice sites and translation start and stop sites of potential genes. GeneMark-EP uses the hints to improve estimation of model parameters as well as to adjust coordinates of predicted genes if they disagree with the most reliable hints (the -EP+ mode). Tests of GeneMark-EP and -EP+ demonstrated improvements in gene prediction accuracy in comparison with GeneMark-ES, while the GeneMark-EP+ showed higher accuracy than GeneMark-ET. We have observed that the most pronounced improvements in gene prediction accuracy happened in large eukaryotic genomes.

INTRODUCTION

One of the major challenges of gene prediction in eukaryotes is finding an optimal way to combine sources of information extrinsic and intrinsic to genome of interest. External information could be transferred from RNA transcripts as well as from cross-species proteins derived from annotated genomes. Integration of transcript information, e.g. RNA-Seq reads, with *ab initio* gene prediction was implemented in several algorithms and software tools, e.g. AUGUSTUS (1), GeneMark-ET (2), EuGene (3,4) and

mGene.ngs (5). Also, a few other tools made use of protein sequences. Complexity of a task of leveraging cross-species protein sequence information for gene identification in a newly sequenced genome is growing with increase of evolutionary distance. Therefore, mapping a protein to genomic locus where a homologous protein is expected to be encoded was a subject for developing specialized tools known as tools for protein spliced alignment [e.g. currently available GeneWise (6), GenomeThreader (7), ProSplign (8) and Spaln (9)]. Beyond a single reference protein, a reference family of homologous proteins could be used to map elements of gene structure conserved in evolution; for instance, AUGUSTUS-PPX (10) uses protein profiles derived from conserved protein domains. Information about intron position, conserved in protein primary structures of multiple homologs, was used in another tool, GeMoMa (11). Notably, an attempt to combine protein profiles with intron position profiles for refinement of predicted genes was made by yet another method, GSA-MPSA (12).

Weakness of methods solely relying on mapping homologous proteins lies in the patchiness of the evidence they generate; a sizable fraction of a whole complement of genes may code for proteins with few or no orthologs. Another weakness is that protein spliced alignments become less accurate as the distance between the two species increases. Therefore, *ab initio* gene finders [e.g. GENSCAN (13), GeneMark.hmm (14), AUGUSTUS (15) or GeneID (16)] have been a necessary part of genome annotation tools and pipelines [e.g. GNOMON (17), PASA (18) and Ensembl (19)].

Application of *ab initio* algorithms for genome-wide eukaryotic gene prediction was for long time hampered by the need of tedious and time-consuming training. To address this issue, we have earlier developed an *ab initio* gene finder GeneMark-ES (20,21) with model parameters estimated by iterative unsupervised training. This algorithm did not require expert-based training or hints for building a training set. GeneMark-ET (2) was developed to make GeneMark-ES able to integrate into the training process available tran-

*To whom correspondence should be addressed. Tel: +1 404 894 8432; Email: borodovsky@gatech.edu

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

script information, raw RNA-Seq reads spliced aligned to genome in question.

Here, we describe GeneMark-EP, an algorithm and software tool that integrated into training information extracted from a reference set of cross-species protein sequences. To generate protein hints for a given genomic locus, we first identify a set of proteins, homologous to the protein likely encoded in the genomic locus. Then, a specialized pipeline, ProtHint, computes the hints, a set of mapped splice sites (intron borders) and translation start and stop sites along with the scores characterizing hint confidence. The most reliably determined elements of spliced alignment could be used to directly identify elements of exon–intron structures; we call this mode of algorithm execution with direct gene structure correction as GeneMark-EP+.

A key question is how to find an optimal method of hint incorporation into the *ab initio* algorithm. Unsupervised training implemented in GeneMark-ES carries a risk of convergence to a biased set of model parameters. On the other hand, giving too much weights to protein hints may generate parameters dictated by a narrow set of conserved genes and proteins (22). By design, the GeneMark-EP algorithm combines strong features of both methods: (i) ability of unsupervised iterative training of an *ab initio* gene finder to create a set of training sequences with a size beyond reach of conventional supervised training and (ii) ability to correct model parameters as well as (the -EP+ mode) structures of newly discovered genes with the hints derived from homologous cross-species proteins. The new method falls into category of gene prediction methods with semi-supervised training.

MATERIALS AND METHODS

For assessment of GeneMark-EP as well as ProtHint accuracy, we selected annotated genomes from diverse clades: fungi, worms, plants, insects and vertebrates (Table 1). The genome length varied from <100 Mb (*Neurospora crassa*) to >1.3 Gb (*Danio rerio*). With exception of *Solanum lycopersicum*, a species representing large genome plants important for economy, all selected species are model organisms whose genomes presumably have high-quality annotation. To assess accuracy of gene prediction made for model species, we compared genes predicted and annotated on a whole genome scale. In case of *S. lycopersicum*, we used a limited set of genes, validated by available RNA-Seq data. In all genomic datasets, contigs not assigned to any chromosome were excluded from the analysis as well as genomes of organelles.

We used OrthoDB v10 protein database (23) as an all-inclusive source of protein sequences. Still, for generating protein hints for particular species we used subsets of OrthoDB: plant proteins for gene prediction in *Arabidopsis thaliana*, arthropod proteins for gene prediction in *Drosophila melanogaster*, etc. (Table 2).

As an additional test set, we used annotation of major protein isoforms available in the APPRIS database (24); this assessment was done for *C. elegans*, *D. melanogaster* and *D. rerio* (Supplementary Table S1). Arguably, accuracy of prediction of major isoforms is of significant interest, since in a

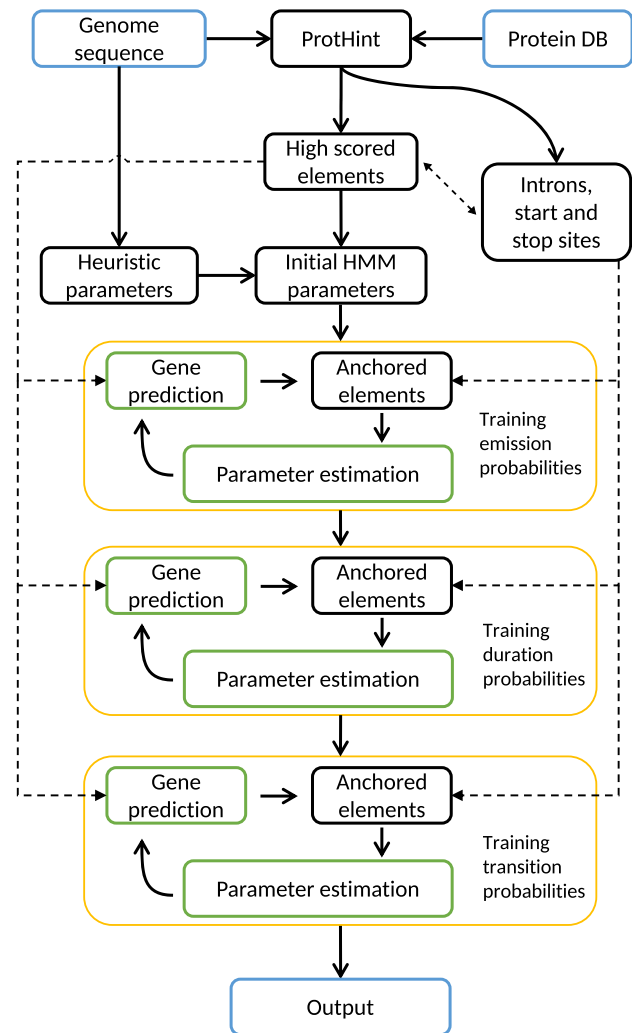


Figure 1. A flowchart of the GeneMark-EP and -EP+ iterative training.

gene locus the major isoform was observed to be expressed in higher volume than other (minor) isoforms (24).

Integration of genomic sequence patterns and protein homology into gene prediction

The GeneMark-EP and -EP+ algorithm goes step by step through the following tasks: (i) selection of genomic regions, *seed regions*, containing gene candidates (*seed genes*); (ii) identification for each seed gene a set of homologous proteins; (iii) constructing spliced alignments of homologous proteins to each seed region and generating hints for exon–intron structures; (iv) running iterative semi-supervised training with selection of most reliable (*anchored*) elements of predicted genes in each iteration; and (v) final gene prediction without (-EP mode) or with additional option (-EP+ mode) to enforce high-confidence hints in predicted exon–intron structures (Figure 1).

Tasks (i)–(iii) are devoted to generating protein hints and are solved by the ProtHint pipeline (Figure 2). To determine *seed regions* within a long genomic sequence [task (i)], we run unsupervised training of parameters of statistical

Table 1. Genomes used for assessment of GeneMark-EP and GeneMark-EP+ performance

Species	Assembly version	Genome size (Mb)	Annotation version	Number of genes in annotation	Introns per gene
<i>Neurospora crassa</i>	GCA_000182925	40	Broad Institute (2013)	10 785	1.7
<i>Caenorhabditis elegans</i>	GCA_001483305	100	WormBase WS271 (May 2019)	20 172	5.7
<i>Arabidopsis thaliana</i>	GCF_000001735	119	Tair Araport11 (June 2016)	27 445	4.9
<i>Drosophila melanogaster</i>	GCA_000001215	134	FlyBase R6.18 (June 2019)	13 929	4.3
<i>Solanum lycopersicum</i>	SL4.0	773	Consortium ITAG4.0 (September 2019)	33 562	3.5
<i>Danio rerio</i>	GCF_000002035	1345	Ensembl GRCz11.96 (May 2019)	25 254	8.2

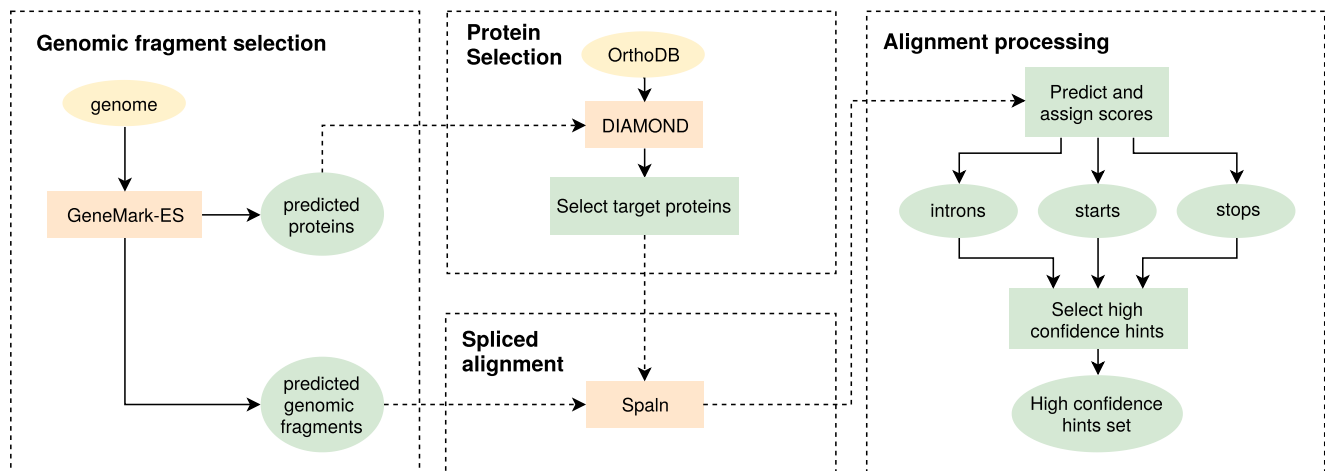
Introns per gene values were computed with respect to the whole number of genes, including single-exon genes.

Table 2. Characteristics of the OrthoDB v10 taxonomic space for each species we tested

Number of species in the same taxonomic unit	Genus	Family	Order	Class	Phylum	Kingdom	OrthoDB root used for tests	Number of proteins in the root
<i>N. crassa</i>	0	1	7	96	364	548	Fungi	5 850 648
<i>C. elegans</i>	2	2	4	5	6	447	Metazoa	8 266 016
<i>A. thaliana</i>	1	7	9	–	99	116	Plantae	3 510 742
<i>D. melanogaster</i> ^a	19	19	55	147	169	447	Metazoa	8 266 016
<i>S. lycopersicum</i>	1	9	10	–	99	116	Plantae	3 510 742
<i>D. rerio</i> ^a	0	4	4	49	245	447	Metazoa	8 266 016

The number of species is naturally the largest in the kingdom section of the database.

^aFor tests in the genus-, family- and order-excluded modes for *D. melanogaster* and *D. rerio*, the phylum set was used as the largest set of reference proteins.

**Figure 2.** An overview of the ProtHint pipeline.

models by GeneMark-ES (20) and generate *ab initio* gene predictions. To create a *seed region*, each predicted gene, a *seed gene*, is expanded upstream and downstream by adding 2000 nt margins. To identify proteins homologous to a *seed protein*, task (ii), we run DIAMOND similarity search (25) with a *seed protein* as a query against a protein sequence database (e.g. a section of OrthoDB). A set of proteins with statistically significant hits defines a set of target proteins presumed to be homologous to the query, the seed protein. Task (iii) is done to generate spliced alignments of multiple protein targets to the seed region [done by either Spaln (9) or ProSpign (8)] and to process the results of alignments in order to infer elements of exon–intron structures

(introns, splice sites, translation starts and stops) characterized by reliability scores. Mapped gene elements with reliability scores exceeding certain thresholds are designated as high-confidence hints. Final tasks (iv) and (v) correspond to training and prediction steps of GeneMark-EP and -EP+. At these steps, we use the hints to exon–intron structure coordinates as an input to an expectation-maximization-type algorithm that finds models of compositional patterns of protein-coding and non-coding regions simultaneously with the most likely parse of genomic sequence into coding and non-coding regions.

Iterative training of the GeneMark-EP statistical models [tasks (iv) and (v)] works as follows. In the first itera-

tion, full-length introns mapped by ProtHint with scores exceeding a stringent threshold (high-confidence elements) are used to estimate parameters of splice site models as well as branch point site models (particularly important for intron models of fungal genomes). The splice site models together with the heuristic models of protein-coding and non-coding regions make a complete set of models of a hidden Markov model (HMM) (20). The models are used in the first run of the Viterbi algorithm [see (14)] that generates a maximum likely parse of genomic sequence into coding and non-coding regions, the parse delineating the first set of genes predicted by GeneMark-EP. Next, we analyze available data to make updated training sets and re-estimate model parameters. We compare coordinates of exons predicted by GeneMark.hmm and intron hints determined by ProtHint within the *seed* regions. This comparison leads to selection of *anchored* elements, the exons with at least one splice site identified by both GeneMark.hmm and ProtHint. A set of anchored exons along with a set of predicted single-exon genes (with length >800 nt) comprises an updated training set for the three-periodic Markov chain model of protein-coding region (26). Sequences of introns bounded by two anchored splice sites as well as intergenic sequences bordered by anchored terminal and initial exons of adjacent genes (Figure 3) are used for updating parameters of the non-coding region model. The set of updated models is used by the Viterbi algorithm to generate a new set of predicted genes. A new update of anchored elements and the next round of parameter re-estimation follows.

Several probability distributions used in GeneMark-EP, such as length distributions of exon, intron and intergenic regions, are initially defined as uniform ones. More accurate estimation of these distributions is done in subsequent steps of iterative training (Figure 1). Also, in the later steps we estimate parameters of the three-phase models of splice sites indexed by a nucleotide position after which the intron divides a codon triplet. Prior to the final iteration, we update estimates of the HMM transition probabilities that affect frequencies of genes with given number of introns. In experimental runs for genomes with different lengths, we have verified that seven iterations were sufficient for GeneMark-ES and six iterations for GeneMark-EP and -EP+ to reach convergence in coordinates of predicted genes and values of model parameters.

Gene predictions made in the final iteration are reported as output of GeneMark-EP. The Viterbi algorithm could be run with enforcement of high-confidence elements mapped by ProtHint. Particularly, it is done by modifying components of the object function of the Viterbi algorithm associated with chosen hidden states. The sites that must be enforced receive high values of objective function to ensure their addition to a path selected by the optimization algorithm seeking the maximum value of the log Viterbi objective function. This mode of gene prediction produces the GeneMark-EP+ output.

Note that GeneMark-ES, -ET, -EP, and -EP+ algorithms are designed to predict non-overlapping genes with no alternative isoforms. This design suits the paradigm that each gene locus encodes a major (expressed in most tissues) protein isoform (24).

ProtHint: generating footprints (hints) of multiple homologous proteins for a genomic locus

General logic. The ProtHint role (Figure 2) in GeneMark-EP and -EP+ is 2-fold. This pipeline generates two sets of protein hints. The smaller one, the set of high-confidence hints, includes hints with high scores that ensure their high specificity. The larger one includes hints that have scores exceeding a liberally set threshold; thus, these hints have lower specificity but larger sensitivity. In the process of hint generation, ProtHint takes a *seed protein* and uses it as a query in similarity search for homologs of a true protein presumably encoded in the seed region. Next, ProtHint constructs spliced alignments of the detected homologs (target proteins) to the seed region. The whole set of multiple spliced alignments is then processed together to identify protein hints, mapped coordinates of the candidate splice sites and translation start and stop sites. The hint scoring system is discussed in detail in Supplementary Data.

Technically, for a given *seed protein*, ProtHint runs DIAMOND (25) against a relevant section of the OrthoDB database and retains in the output up to 25 target proteins (with hit *E*-value better than 0.001). Next, the target proteins are spliced aligned by Spaln (9) back to the seed region. Notably, the hints are defined by ProtHint processing of the Spaln raw pairwise alignments instead of using annotation of exons in the Spaln output. Rather frequently, multiple target proteins aligned to the original seed region may map out one and the same sequence fragment as an intron. Such an outcome would define an intron hint with a higher confidence than if an intron candidate is mapped only once.

Score system for introns. As described earlier, an expected evolutionary conservation between primary structures of target proteins and a protein encoded in the seed region has to be quantified and used for accurate identification of a new gene. To facilitate this quantification, we define three types of scores for exons and introns (AEE, IBA and IMC; see below) and two types of scores for candidate gene starts and stops (SMC and BAQ; see below).

Alignment of entire exon (AEE) score is defined as a score of the Spaln (or ProSplign) alignment of exon translation and a target protein (for more details, see Supplementary Data).

Intron border alignment (IBA) score is computed from kernel modified alignment scores of two adjacent exons with larger weights given to parts close to the donor and acceptor splice sites. An IBA score (within a window of length w , being 10 amino acids by default) is computed as follows.

For downstream (and upstream) exon defined in the Spaln spliced alignment, we compute S_d (and S_u) as

$$S_d = \sum_{i=1}^w S_a(G_i, P_i) \times W(i). \quad (1)$$

Here, $S_a(G_i, P_i)$ is a substitution score defined for target protein amino acid P_i and a codon defined amino acid G_i (see Suppl. Materials), and $W(i)$ is the weight function. For

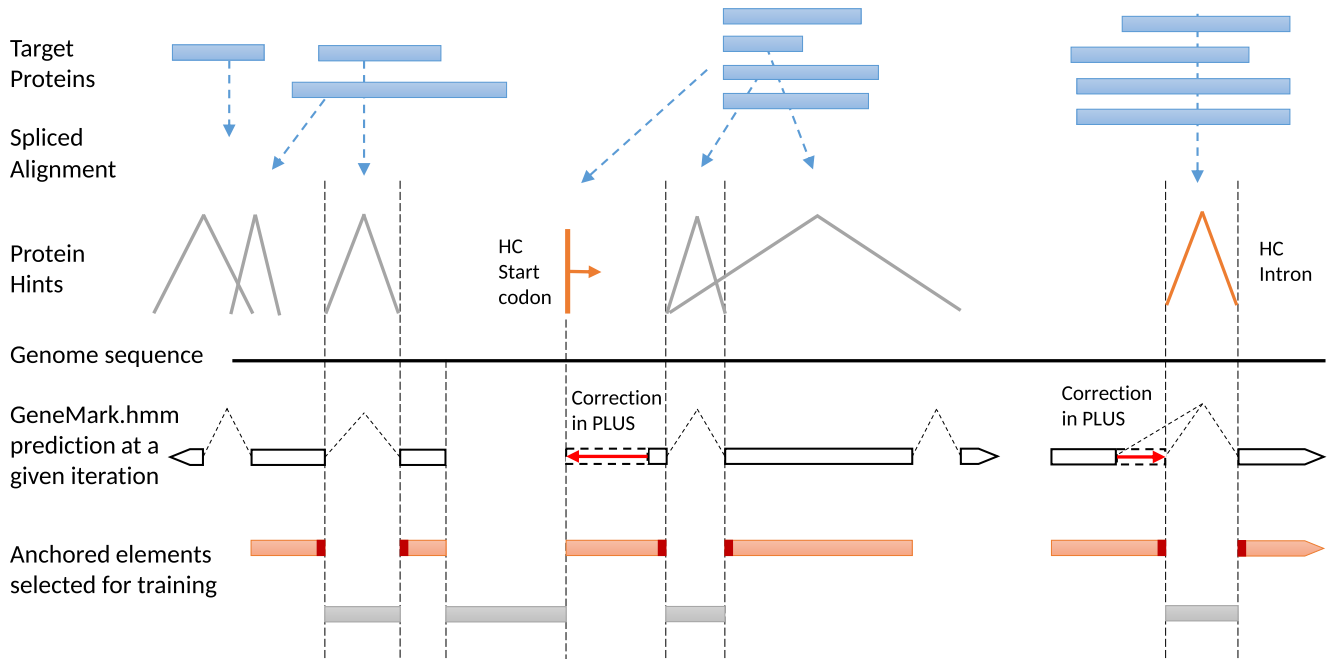


Figure 3. Selection of sequence regions for GeneMark-EP+ training with enforcement of high-confidence (HC) hints.

instance, for a downstream exon S_d ,

$$W_i = \frac{K(i)}{\sum_{i=1}^w K(i)}, \quad (2)$$

where $K(i)$ is the kernel value for position i counting in codons from an acceptor site. In a linear kernel,

$$K(i) = 1 - \frac{|i| - 1}{w}. \quad (3)$$

Then, we take a geometric mean of values of S_d and S_u :

$$S_{\text{intron}} = \begin{cases} \sqrt{S_u \times S_d}, & \min(S_u, S_d) > 0, \\ 0, & \text{otherwise.} \end{cases} \quad (4)$$

Finally, the IBA score is obtained by normalizing the S_{intron} score into $(0, 1)$ range: IBA score = $S_{\text{intron}} / \max(S_a)$, where $\max(S_a)$ is a maximum score among elements of the amino acid substitution matrix.

Intron mapping coverage (IMC) score is a count of how many times a given intron was exactly mapped by spliced alignments of target proteins. The IMC score is computed only from the set of the introns that IBA score exceeds certain level.

Notably, instead of all introns with identical coordinates and different IBA scores related to different target proteins, we deal with a single intron characterized by the maximum of individual IBA scores among all collapsed introns.

Application of the intron scores. For a particular seed region, we use the following method to define hints to exon–intron structure from a set of spliced alignments of target proteins.

First, we select introns whose two adjacent exons have scores $AEE \geq E_t$, where E_t is a chosen threshold. For $E_t = 25$, in a modeling on known genomes, we observed relatively

high S_n value of the candidate introns (Supplementary Figure S1). Further increase of E_t eliminated true introns while not significantly improving S_p value.

Next, to reduce number of false positives in the obtained set of introns, we selected a subset with the IBA score $> I_t$, where I_t is another chosen threshold. Our modeling has shown an increase in the S_p value of the candidate introns for $I_t = 0.1$ that occurred without noticeable change in S_n (Supplementary Figure S1).

Thus identified sets of introns for all the seed regions represent a set of *all mapped introns*; it is used as external evidence to generate *anchored introns* for GeneMark-EP training steps as described earlier.

Finally, within the set of *all mapped introns* we select a narrower set of *high-confidence introns*. These introns must have canonical GT–AG splice sites, an IMC score ≥ 4 and an IBA score ≥ 0.25 (Figure 4, Supplementary Figure S2).

We use high-confidence introns to estimate initial parameters of the GeneMark-EP intron model. Besides, these introns are enforced in the prediction step of the GeneMark-EP+ mode.

Score system for translation starts and stops. Similarly to scores introduced for intron hint generation, we define a *border alignment quality (BAQ)* score for translation starts and stops. This score is computed for w amino acids downstream (upstream) of start (stop) codon, weighted by a kernel-dependent function (Equation 1).

The second type of score is the *site mapping coverage (SMC)* score. This score is a count of N-terminals (C-terminals) of target proteins aligned to a particular start (stop) codon position of a candidate gene. The SMC scores are computed only from the sets of initial (terminal) exons whose BAQ scores exceed certain level.

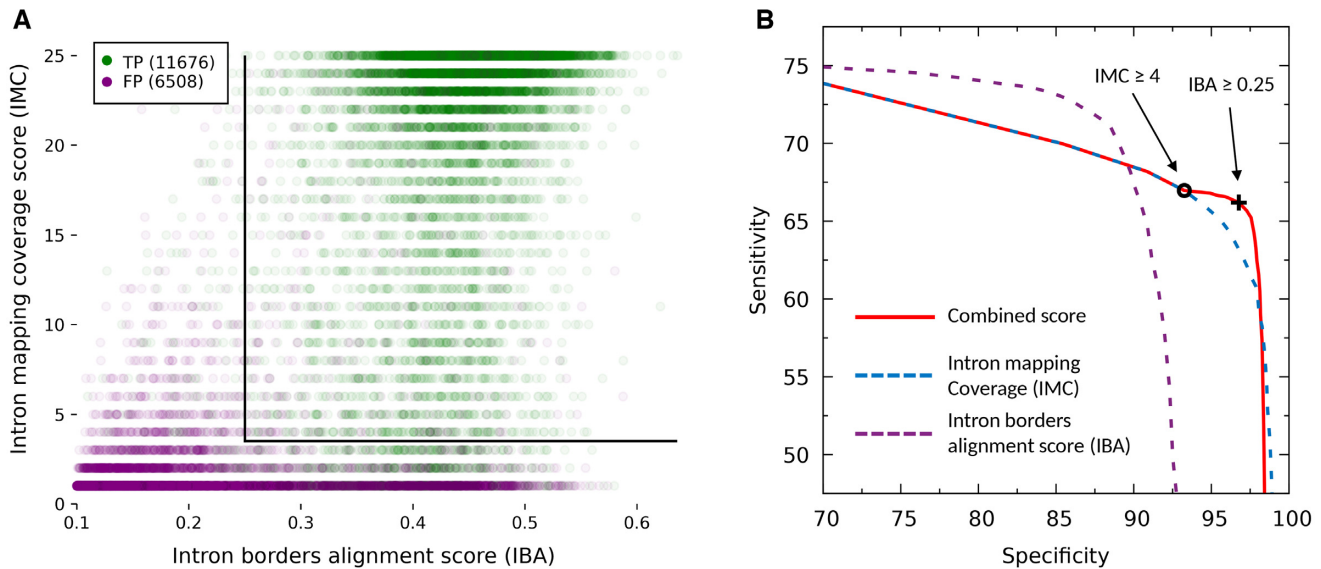


Figure 4. ProtHint intron processing in case of *N. crassa*. Introns were generated by spliced alignments of target proteins from species beyond *Neurospora* genus. (A) Distribution of the score vectors (IBA, IMC) of true positive (green) and false positive (purple) introns. The black lines represent cutoffs at $IMC = 4$ and $IBA = 0.25$. Total numbers of false and true positives are shown in the upper left corner. (B) Sn and Sp of intron sets selected by thresholds on IBA score and IMC score. IMC score is computed for introns that have IBA score ≥ 0.1 and exon AEE score ≥ 25 . The left branch of the curve reflects (Sp, Sn) values of the sets of introns selected by using IMC threshold from 0 to 4. The one with the IMC threshold = 4 is recorded as set A—the set corresponding to the black circle in the red curve. Then, the right branch of the curve reflects (Sp, Sn) of the set of introns generated by applying to set A an IBA score threshold changing from 0 to 0.25 and up to 1.0. Set B corresponds to the black cross in the red curve; introns in this set have $IMC \geq 4$ and $IBA \geq 0.25$. Separate curves for IMC score change (dashed blue) and IBA score change (dashed purple) are shown as well.

If a set of target proteins for a given seed region generates footprints situated more upstream than others, alternative start candidates situated downstream are removed from consideration (Supplementary Figure S3, details in Supplementary Data). We have observed that using these rules leads to increase in the hints' accuracy (Table 3, Supplementary Table S2).

Application of BAQ and SMC scores. Altogether, selection of a set of *all* translation start and stop hints is done by the following method.

A start codon candidate is an ATG codon present in a mapped initial exon and aligned to N-terminal methionine in a target protein; a stop codon candidate is a stop codon in a mapped terminal exon.

A candidate initial (terminal) exon containing candidate gene start (stop) should have AEE score ≥ 25 and BAQ score for candidate start (stop) codon ≥ 0 .

To select a subset of *high-confidence hints*, we choose stop codon candidates with SMC score ≥ 4 as well as start codon candidates with SMC score ≥ 4 and no overlap by longer target proteins (Supplementary Figure S3). The set of *high-confidence hints* to translation starts and stops is used to estimate parameters of GeneMark-EP models of translation initiation and termination sites. Also, the *high-confidence hints* are directly enforced in the prediction step of GeneMark-EP+.

Do introns mapped by ProtHint tend to occur in gene regions coding for conserved domains?

To address this question, we use the following procedure. Annotated genes are translated to proteins and used as

queries in RPS-BLAST (27) to search (E -value = 0.01) against NCBI Conserved Domains Database (28). Results of the RPS-BLAST searches are processed with *rpsbproc* utility (28) to generate a map of conserved domains for each RPS-BLAST query. Finally, coordinates of the conserved domains are mapped back to a *seed region* of genomic DNA and compared with the ProtHint output to find out how many introns are mapped into regions coding for conserved domains. We conducted this analysis for genes of *D. melanogaster*, *C. elegans* and *D. rerio* genomes annotated in the APPRIS database (24) as genes coding for principal protein isoforms (see 'Results' section).

Assessment of GeneMark-EP and -EP+ gene merging and gene splitting errors

Gene merging and splitting errors are expected to be reduced by the use of homologous protein-derived hints to gene translation starts and stops. This expected improvement in prediction accuracy of GeneMark-EP+ could be more accurately observed on properly prepared test sets. Prior to evaluation of gene splitting, we had to exclude from the test sets (i) genes fully overlapping shorter genes present inside introns in any strand; (ii) genes with larger isoforms combining or including shorter alternative components (Supplementary Figure S4); and (iii) genes with introns longer than 10 000 nt (the default maximum intron length). For genes with annotated multiple alternative isoforms, we used the longest one as a representative. Prior to evaluation of gene merging, overlapping genes present in annotation (e.g. a gene within an intron) were merged

Table 3. Case of *A. thaliana*: sensitivity and specificity of all gene start hints created by ProtHint as well as of high-confidence start hints

		All reported starts	Filtered with SMC ≥ 4	Filtered with SMC ≥ 4 and exon overlap = 0
<i>A. thaliana</i>	Sn	69.3	62.9	61.4
	Sp	70.9	89.8	94.4

High specificity was achieved through filtering by SMC scores as well as by removal of candidate starts overlapped by at least one target protein suggesting an alternative start upstream. Sensitivity was defined with respect to a full complement of starts, including alternative ones as given in annotation. The numbers were generated in tests with reference proteins from species outside a relevant genus. Results for all test species are shown in Supplementary Table S2.

into a single gene in order to exclude such cases from being counted as merged genes.

RESULTS

We have compared gene prediction accuracy of GeneMark-EP and -EP+ with accuracy of GeneMark-ES and GeneMark-ET. In addition, we made an accuracy assessment of ProtHint. We worked with genomes of six species *N. crassa*, *C. elegans*, *A. thaliana*, *D. melanogaster*, *S. lycopersicum* and *D. rerio* (Table 1). Genomes of model organisms (all the species but *S. lycopersicum*) were expected to have sufficiently accurate annotation. For model organisms, we only made comparisons between predicted and annotated gene coordinates on whole genome scales. In case of *S. lycopersicum*, we also built a test set of genes validated by RNA-Seq data.

In genomes of *C. elegans*, *A. thaliana*, *D. melanogaster* and *D. rerio* all regions of annotated pseudogenes were excluded from comparisons. In case of *D. rerio* we excluded annotated partial exons (ubiquitous in this genome) from exon-level accuracy assessment; we computed gene-level sensitivity only for genes having in annotation complete alternative transcripts.

We used OrthoDB v10 (23) as a source of protein sequences partitioned into relevant taxonomic divisions; particularly, we used plant division for *A. thaliana*, arthropod division for *D. melanogaster*, etc. (Table 2).

A principal feature of the new method is use of *multiple homologous proteins* for hint generation. We had to model practical situations when the *closest* evolutionary distance from a seed protein to the *homologs* detected by similarity search in a protein database could vary significantly depending on the evolutionarily distance to most closely related species with sequenced and annotated genome. To simulate these variations in our tests, we introduced restrictions on evolutionarily distance to the closest species from which the target proteins could be recruited. These restrictions were implemented by removing from the protein database (i) proteins encoded in the genome of a given species; (ii) proteins from all species from the same subgenus; (iii) proteins from the same genus; (iv) proteins from the same family; (v) proteins from the same order; and (vi) proteins from the same phylum. Notably, distributions of numbers of species within a genus, family, etc. defined by a given species are species specific (Table 2).

Assessment of accuracy of GeneMark-EP and -EP+

For each species (Table 1), we determined how the accuracy of GeneMark-EP and -EP+ at gene level (Figure 5)

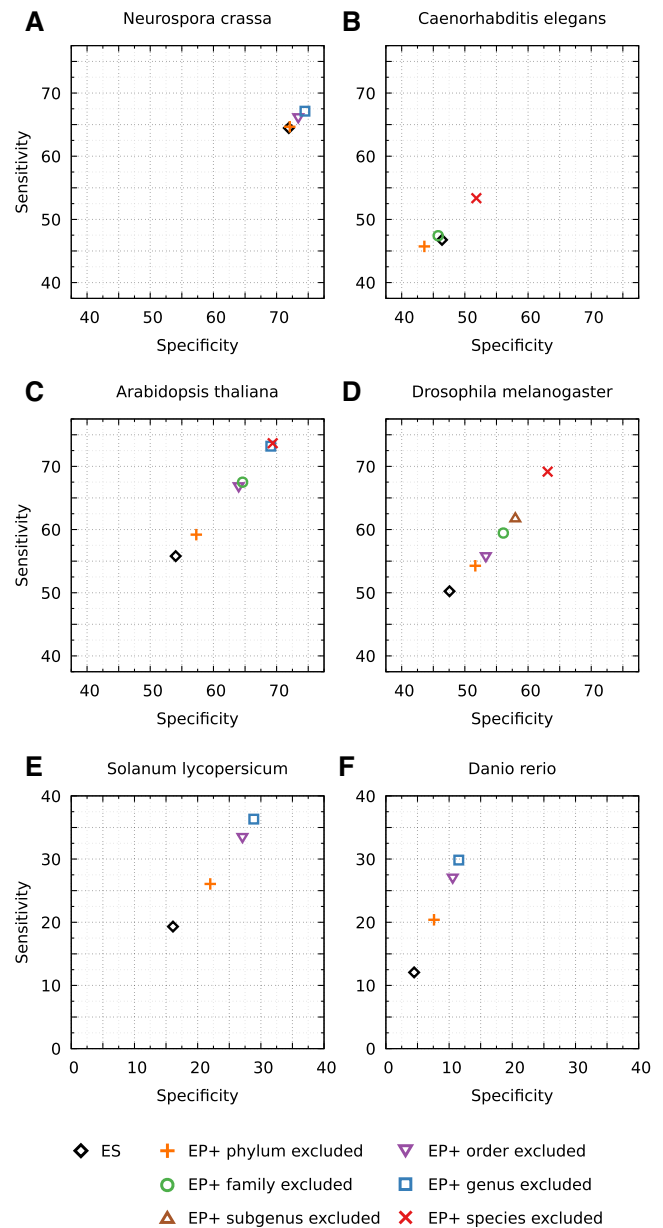


Figure 5. Comparison of GeneMark-ES and GeneMark-EP+ accuracy on gene level. Accuracy of GeneMark-EP+ is shown for cases when ProtHint works with different size sets of reference OrthoDB proteins: from the largest (only proteins from the same species are excluded) to the smallest (proteins of the whole phylum excluded). A gene prediction is considered to be correct if it matches one of the annotated isoforms. For *D. rerio*, gene-level Sn was computed only with respect to complete genes.

and exon level (Supplementary Figure S5) depended on the choice of a set of reference proteins. The pattern of accuracy change at *gene level* (Figure 5) was similar to the one observed at *exon level*; therefore, we show the results of accuracy assessment at gene level in the main text, while the results for exon level are provided in Supplementary Data (Supplementary Figure S5). Even more details on accuracy assessment of GeneMark-EP (running without enforcement of high-confidence hints) and of GeneMark-EP+ are given in Supplementary Data (Supplementary Table S3).

We present results for the three groups of genomes: fungal genomes, compact eukaryotic genomes and large eukaryotic genomes.

Fungal genomes. *N. crassa*. Accuracy of GeneMark-ES was high, as it has been typical for fungal genomes (21). Even with hints originated from the largest set of reference proteins, those outside genus/order, GeneMark-EP+ improved Sn value of GeneMark-ES only by ~2% (Figure 5a). With smaller set of more remote reference proteins that originated from the species outside fungal phylum, the accuracy of GeneMark-EP+ matched the accuracy of GeneMark-ES (Figure 5a). This result went in line with previous observations that GeneMark-ES was highly efficient *ab initio* gene finder for fungal genomes (21). We observed earlier that for fungal genomes a support of training with data from splice-aligned RNA-Seq reads makes accuracy of GeneMark-ET at best slightly higher than that of GeneMark-ES.

Compact eukaryotic genomes. *C. elegans*, *A. thaliana* and *D. melanogaster*. When GeneMark-EP+ used the largest set of reference proteins (just without proteins from the same species), we saw for *A. thaliana* and *D. melanogaster* an improvement of ~20% in both Sn and Sp in comparison with GeneMark-ES (Figure 5c and d). As target proteins were coming from larger and larger evolutionary distances, the accuracy did steadily decrease. When target proteins were selected outside the same phylum, there was an increase of only 5% in gene-level Sn and Sp in comparison with GeneMark-ES. For *C. elegans*, when the set of reference proteins excluded just proteins of the same species, GeneMark-EP+ improved the accuracy of GeneMark-ES by ~6% (Figure 5b). We observed almost no difference between GeneMark-EP+ and GeneMark-ES when the reference proteins were only from species outside the *C. elegans* family and a slight decrease in accuracy (~2%) for reference proteins outside of the taxonomical phylum. Notably, the gene-level accuracy for *C. elegans* was lower than that for other species with compact genomes.

Large eukaryotic genomes. *S. lycopersicum* and *D. rerio*. The gene-level accuracy of GeneMark-ES was low for these genomes (between 5% and 20%). GeneMark-EP+ improved the accuracy for *S. lycopersicum* by ~15%, when it used a protein reference set from species outside of the tomato genus or order (Figure 5e). For *D. rerio*, having a reference set of proteins without those from same genus or the same order as *D. rerio* led to Sn and Sp improvement of ~20% and ~5%, respectively (Figure 5f). However, the improve-

ments were twice as low when reference proteins were available only outside the *S. lycopersicum* or *D. rerio* phylum.

Relatively low gene prediction accuracy in large genomes could be partially attributed to incorrect and/or incomplete gene annotations. Therefore, we made additional effort to refine test sets in *D. rerio* and *S. lycopersicum* by selecting genes supported by RNA-Seq data.

We observed that if annotated genes of *S. lycopersicum* genome were supported by RNA-Seq, they were significantly better predicted by GeneMark-EP+ (Supplementary Table S4). To generate intron hints from RNA-Seq, we used VARUS (29). We divided annotated tomato multi-exon genes into two groups: (a) genes with all introns predicted by VARUS and (b) all other genes. GeneMark-EP+ sensitivity (for a GeneMark-EP+ run having reference proteins outside the *S. lycopersicum* genus) was 40% better in set (a) than in set (b), on gene, exon and intron levels. It is important to emphasize that RNA-Seq information was not used in GeneMark-EP+. Sensitivity defined for the set of introns mapped by ProtHint was also better in set (a) by ~40% (Supplementary Table S4).

We already mentioned that *D. rerio* annotation had many partial exons that in turn would be parts of incomplete transcripts. We evaluated exon-level Sn separately for exons within complete and incomplete transcripts (Supplementary Table S5) and observed 75.1% exon Sn in the 'complete' group versus 67.6% in the 'incomplete' group. Similarly, gene-level sensitivity was better by 6% in predicting genes with complete transcripts compared to all genes (Supplementary Table S5).

Altogether, we observed that for majority of the considered species, the accuracy of GeneMark-EP+ was better than accuracy of GeneMark-ES, regardless of how large a set of reference proteins was used for spliced alignments (Supplementary Table S3, Figure 5, Supplementary Figure S5). For the fungal genome, *N. crassa*, an improvement was negligible due to ability of GeneMark-ES to deliver high accuracy for fungal genomes; we also observed a small decrease of accuracy in the *C. elegans* test with phylum-excluded reference set of proteins (Figure 5).

Comparison with GeneMark-ET

In addition, we compared GeneMark-EP with GeneMark-ET (2) that uses RNA-Seq short reads to provide external information (hints to intron coordinates) to select anchored gene elements for the GeneMark-ET algorithm parameter estimation. GeneMark-ET does not have an '-ET+' mode in which predictions are directly guided by high-confidence hints. We ran GeneMark-ET with hints to coordinates of introns mapped by VARUS from RNA-Seq reads. VARUS automatically sampled, downloaded and aligned reads from NCBI's Sequence Read Archive (SRA) with time stamp of 22 January 2020 (30). The time stamp is important for the reproduction of results, since the VARUS outcome depends on the amount of RNA-Seq data deposited to SRA. As one could see (Supplementary Table S3), the accuracy of GeneMark-ET with training guided by hints derived from mapped RNA-Seq reads is very close to the accuracy of GeneMark-EP with training guided by hints derived from mapped proteins.

To get an idea to what extent a whole complement of genes is covered by hints originated from protein or transcript data, we determined fractions of genes with at least one hint or at least with one high-confidence hint generated by ProtHint (Supplementary Table S6) as well as the fractions of genes with at least one hint generated by VARUS (Supplementary Table S7).

When we used protein reference sets with largest sizes, >80–85% of annotated genes of *A. thaliana*, *C. elegans*, *D. melanogaster*, *N. crassa* and *D. rerio* harbored protein hints; this percentage was 69% in *S. lycopersicum* (Supplementary Table S6a). These numbers went down to a range of 40–60% when the sets of reference proteins decreased to their minimal sizes, which were reached when all proteins from the same phylum species were excluded. It was interesting to make comparison of just cited data with the percentage of annotated genes receiving high-confidence protein hints. For the largest reference sets, we observed that percentage of genes with high-confidence hint was lower than respective previous figures by just 3–5% in *A. thaliana*, *D. melanogaster*, *S. lycopersicum* and *D. rerio*; the drop, however, was 13% in *N. crassa* and 24% in *C. elegans* (Supplementary Table S6b). This large drop for *C. elegans* corresponds to the lowest accuracy of GeneMark-EP+ among all the species considered (Figure 5b). For the smallest reference sets, proteins from the same phylum excluded, the drop in high-confidence hints coverage was within 10–15% range for all species but *N. crassa* (22%).

The coverage by all protein hints (high confidence or not) in comparison with coverage by RNA-Seq-derived hints (Supplementary Table S7) was higher by 10–15% for all species but *C. elegans* (lower by 10%) and *D. rerio* (lower by 5.5%).

Finally, we saw that annotated genes that did not harbor any hints, either from ProtHint or from VARUS, made a sizable fraction only in *S. lycopersicum*—24.4% (Supplementary Table S8), while in all other species these fractions were rather small: from 5% (*C. elegans*) to 12.3% (*N. crassa*).

In our data, the percentage of annotated genes covered in a given genome by the largest set of protein hints is comparable to or higher than the percentage of genes covered by the transcript-derived hints. Also, we saw that a vast majority of genes in the six species receive one or another type of external support. The genes that do not have external support belong to intersections of sets of genes that code for orphan or unique proteins and sets of genes that did not show detectable expression in the experiments measuring gene expression. Still, we have to make correction for the fact that RNA-Seq-derived hints were not defined for single-exon genes even if they were expressed.

Sources of improvements in gene prediction

Better performance of GeneMark-EP+ in comparison with GeneMark-ES is expected due to two factors: (a) model parameterization on a better validated training set as the training process becomes semi-supervised instead of unsupervised and (b) enforcement of high-confidence hints in gene prediction steps. Notably, even when direct corrections are not made [GeneMark-EP mode where factor (b) is absent], for all the species but fungi GeneMark-EP showed im-

provement over GeneMark-ES. Surprisingly, GeneMark-EP showed only small fluctuations in accuracy when the size of the reference set of protein increased by including more evolutionarily close species (Supplementary Table S3).

The accuracy of GeneMark-EP+ was about the same as the accuracy of GeneMark-EP when the smallest reference set of proteins was used (proteins from species outside the phylum of the species in question). Accuracy of GeneMark-EP+ increases significantly when reference proteins from more evolutionarily close species are included, while accuracy of GeneMark-EP stays about the same. The only exception was *C. elegans* in which GeneMark-EP gene-level accuracy dropped by ~4% for the reference set of species outside the same phylum in comparison with GeneMark-ES (while GeneMark-EP+ shows the accuracy close to the level of GeneMark-ES; Supplementary Table S3).

These observations suggest that even a relatively small number of anchored introns play a critical role in parameter estimation in GeneMark-EP. Further increase in the number of anchored introns does not improve parameters of GeneMark-EP. For the case of *C. elegans*, one could argue that the sufficient minimal number of anchored introns was not found when proteins of the reference set were limited to ones from the species outside the *C. elegans* phylum.

To differentiate contributions into GeneMark-EP+ performance, we compared runs that used only high-confidence intron hints with runs that used only high-confidence hints for gene starts and stops (Supplementary Table S9). This experiment showed that enforceable hints of both kinds contributed equally to overall accuracy improvement. However, these hints contribute unequally to reduction of different types of errors. Enforcement of high-confidence intron hints led to higher prediction accuracy of internal exons, while enforcement of high-confidence hints to gene starts and stops led to reduction of errors in initial and terminal exons.

We observed that GeneMark-ES was more likely to generate gene merging than gene splitting errors (Table 4); for instance, comparison of the *A. thaliana* gene predictions and annotation showed 360 split genes and 743 merged genes. Use of GeneMark-EP (with reference proteins outside the same genus) decreased frequency of errors in gene merging (a ~15% decrease in all species); however, it also caused a slight increase in gene splitting (Table 4). Transition to GeneMark-EP+ (the last column in Table 4) reduces gene merging dramatically.

Enforcement of only high-confidence intron hints in GeneMark-EP+ reduced the number of split genes (by enforcing introns in place of incorrectly predicted intergenic regions). Still, these hints have little or no effect on the gene merging (Table 4). The most significant effect on gene splitting was observed for *D. rerio*—2010 split genes in the -EP+ mode compared to 2976 in the -EP mode.

Enforcement of high-confidence hints to gene starts and stops significantly reduced number of merged genes and caused a slight increase in number of split genes. For instance, the number of merged genes dropped by ~500 in *A. thaliana* between GeneMark-ES and GeneMark-EP+, a ~66% improvement; ~50% improvement was observed for the other species in our tests, except *C. elegans*. Altogether, GeneMark-EP+ (Table 4, last column, bold font) achieved

Table 4. Numbers of merged and split genes in predictions of GeneMark-ES, -EP and -EP+ with enforcement of (a) only high-confidence hints to introns, (b) only high-confidence hints to gene starts and stops, or (c) both (a) and (b)

	Genes	ES	EP	EP+ introns (a)	EP+ starts/stops (b)	EP+ full (c)
<i>N. crassa</i>	Merged	129	89	92	64	74
	Split	83	96	89	106	91
<i>C. elegans</i>	Merged	1120	1076	1090	1019	1029
	Split	588	725	614	731	622
<i>A. thaliana</i>	Merged	743	634	629	215	251
	Split	360	385	242	478	277
<i>D. melanogaster</i>	Merged	544	464	462	311	313
	Split	285	297	204	324	221
<i>S. lycopersicum</i>	Merged	2304	1871	1793	1165	1192
	Split	1550	1644	1139	1962	1252
<i>D. rerio</i>	Merged	1921	1415	1351	883	884
	Split	2553	2976	2018	3058	2010

All the numbers were obtained for reference sets of target proteins from the species outside of relevant genus.

significant reduction in numbers of both merged and split genes in comparison with GeneMark-ES and GeneMark-EP.

Comparison of GeneMark-EP+ predictions with genome annotations defined by the APPRIS database

We compared GeneMark-EP+ gene predictions with annotations of major protein isoforms in *C. elegans*, *D. melanogaster* and *D. rerio* genomes defined by the APPRIS database (24). This test did show (Supplementary Figure S6) an increase in exon-level sensitivity (by ~4% for *C. elegans* and *D. rerio*, by ~7% for *D. melanogaster*) and a decrease in exon-level specificity (by ~1.5% for *C. elegans*, by 3% for *D. melanogaster* and by ~8% for *D. rerio*) in comparison with our previous assessment results using genome annotations made by respective genomic communities (Table 1). The decrease in Sp could be expected since the APPRIS annotation contains smaller number of exons. The increase in Sn is a positive news indicating that GeneMark-EP+ when making prediction of a single isoform per locus is likely to predict genes for major protein isoforms. At gene level (Supplementary Figure S7), both Sn and Sp were reduced slightly in *C. elegans* and *D. rerio*, and by 5% in *D. melanogaster*. To correctly interpret this result, we have to remind the definition of gene-level accuracy—a gene is counted as correctly predicted if the prediction matches all exons in at least one alternative transcript. Thus, a gene is considered to be predicted correctly if just one of the isoforms (major or not) is correctly predicted (Figure 5, Supplementary Table S3). This is a rather liberal way of computing an Sn value on gene level.

Assessment of accuracy of ProtHint

The main role of ProtHint is generation of a list of coordinates as well as confidence scores of potential borders between coding and non-coding regions in a novel genome. Specific thresholds on confidence scores are defined to select subsets of hints (e.g. high-confidence set). The GeneMark-EP training procedure can tolerate a high number of false

positive intron hints since only a subset, the anchored introns, is used in training. It is important that the set of *all mapped hints* would have high Sn with respect to true gene elements, while the Sp level could be lower. On the other hand, the *high-confidence hints*—those utilized in initial GeneMark-EP+ parameter estimation as well as in the hints' enforcement—have to have high Sp, as these hints are directly enforced in predictions.

How large are fractions of correct hints among hints generated by ProtHint?

When the set of reference proteins had the maximum size (all proteins in a relevant OrthoDB division except those from the same species), the set of intron hints generated by ProtHint had Sn > 75% for exact introns and Sn ~ 70% for gene starts and stops (Table 5, Supplementary Table S10). The value of Sn was dropping down steadily as evolutionary distance to reference proteins was increasing. Particularly, when the proteins from species of the same *order* were excluded, Sn was, on average, ~65% for intron hints and ~40% for gene start and stop hints.

The largest reduction in the volume of the protein reference set—exclusion of proteins from the same phylum—decreased Sn of all reported intron hints down to ~40% on average (Supplementary Table S10). Here, the largest Sn value (the fraction of correct intron hints) was observed for *N. crassa* (60%), and the lowest one for *C. elegans* (26%). At the same time, the value of Sn of gene start and stop hints generated from the smallest reference set of proteins varied significantly between the species, from 8% for *C. elegans* to 30% for *N. crassa* (Supplementary Table S10).

How reliable are 'high-confidence' hints generated by ProtHint?

The sets of high-confidence hints were observed to have high specificity, averaging over 95% (5% of false positives) over the six species. This level remained high even for the smallest sets of reference proteins, proteins from the species outside

Table 5. Accuracy of ProtHint for the *D. melanogaster* genome: sensitivity and specificity of hints to introns, start and stop codons

	Level of exclusion of reference proteins									
	Species		Subgenus		Family		Order		Phylum	
	All reported	High confidence	All reported	High confidence	All reported	High confidence	All reported	High confidence	All reported	High confidence
Intron Sn	79.8	74.6	72.8	62.6	66.2	54.3	49.7	34.4	35.8	20.9
Intron Sp	83.5	98.9	79.6	98.8	79.5	98.8	80.5	99.0	88.4	99.5
Start Sn	70.3	60.7	49.8	36.5	37.7	29.2	22.3	15.9	14.1	9.7
Start Sp	79.5	97.4	75.6	96.7	71.6	95.6	73.4	94.5	75.0	93.5
Stop Sn	75.3	68.4	56.7	45.2	44.7	36.9	26.7	19.8	15.8	11.2
Stop Sp	94.8	99.3	94.2	98.8	92.8	98.5	94.5	98.9	95.8	99.2

The results are shown for all reported hints or just high-confidence hints. The Sn and Sp values are computed based on genome annotation of a full complement of introns, gene starts and stops, including alternative isoforms. Results for all tested species are shown in Supplementary Table S7.

the phylum of interest (Table 5, Supplementary Table S10). In case of *C. elegans*, along with high Sp, we observed low Sn value of high-confidence hints (for all the reference protein sets—larger or smaller), which is explained by the presence of just a few species with sequenced genomes in the *C. elegans* phylum (Table 2). For all other species, a decrease in Sn upon transition from *all mapped* to *high-confidence* hints was small in comparison with the simultaneous increase in Sp.

Distributions of (IMC, IBA) vectors representing intron hints generated for *N. crassa* (both false and true as compared with annotation) are shown in Figure 4a (for the genus-excluded protein reference set). The Sp–Sn curves are generated for sets of intron hints obtained by filtering with IMC and IBA thresholds (Figure 4b).

The distribution of the score vectors (Figure 4a) as well as the behavior of Sp–Sn curves (Figure 4b) depends on selection of the set of reference proteins (genus or order or phylum excluded; Supplementary Figure S2, left and middle panels). A choice of IBA threshold selecting high-confidence intron hints affects accuracy of GeneMark-EP+. We assessed the extent of this effect for *A. thaliana*, *N. crassa* and *S. lycopersicum* (Supplementary Figure S2, right panels). It was shown that the best average prediction accuracy was achieved with IBA threshold set to 0.25. Similar analysis produced necessary thresholds for high-confidence hints to gene starts and stops.

More intron hints are generated in regions encoding conserved protein domains

We found that ~50% of the whole set of introns annotated in the APPRIS set of principal isoforms is located within conserved protein domains (Supplementary Table S11).

In *D. melanogaster*, high-confidence intron hints generated by ProtHint from the ‘species-excluded’ reference set of proteins fell into regions coding for conserved domains in 55.9% of cases (Table 6). This fraction increased significantly as more proteins were excluded from the reference set (e.g. proteins from species outside of the *D. melanogaster* genus). This fraction reached 84.6% when only proteins originated from species outside the *D. melanogaster* phylum were considered (Table 6). Similar trends were observed for *C. elegans* and *D. rerio* (Supplementary Table S12). In the set of *all reported* intron hints, the fraction of introns mapped to regions coding for conserved domains was lower

than that in the set of high-confidence intron hints (Supplementary Table S12); however, the proportion of introns mapped into conserved domain regions also increased upon removing proteins from closely and moderately closely related species.

The same type of fractions computed for ‘high-confidence’ and ‘all reported’ intron hints were almost identical for the species between *D. melanogaster* and *D. rerio* (Supplementary Table S12). Still, for *C. elegans*, however, the fraction of true high-confidence introns was lower (Supplementary Table S12) apparently due to having fewer target proteins from close relatives in the protein database (a factor significantly affecting IMC score).

DISCUSSION

The main reason to develop GeneMark-EP and -EP+ was a clear need to leverage abundant protein sequence data available in public databases for improving accuracy of automatic gene prediction. It was well expected that iterative *ab initio* parameterization of statistical models (as done in GeneMark-ES) would become more precise, especially for large genomes, if we find an efficient method to add data on protein footprints into training and prediction steps. This project has grown into development of a whole GeneMark-EP pipeline, with GeneMark-EP+ mode that directly integrates most confident evidence into predicted exon–intron structures. In this respect, the new pipeline features a new method, ProtHint, developed to find multiple proteins homologous to a gene initially predicted in a genomic locus and then to derive reliable hints to the true gene exon–intron structure by constructing and processing multiple protein footprints. GeneMark-EP and -EP+ should become a universal extension of GeneMark-ES, as its application to a novel eukaryotic genome will be facilitated by use of a vast volume of protein sequences.

Another earlier developed method, GeneMark-ET (2), extended GeneMark-ES to use external evidence generated from transcriptome sequence data, when it is available along with a newly assembled genome.

Existing methods, such as GenomeThreader (7), rely on mapping proteins from closely related species as well as mapping gene elements from aligned genomic sequence of the close species to produce predicted exon–intron structures. However, its prediction accuracy is dropping fast with increase of evolutionary distance between species (6).

Table 6. Fractions of *D. melanogaster* high-confidence intron hints mapped into regions coding for conserved protein domains

Exclusion level	High-confidence intron hints matching APPRIS annotated introns	
	All high-confidence intron hints	High-confidence intron hints that fell into domains
Species	33 894	18 934 (55.9%)
Subgenus	28 437	17 475 (61.5%)
Family	24 670	16 057 (65.1%)
Order	15 829	11 984 (75.7%)
Phylum	9719	8222 (84.6%)

The hints were generated from sets of reference proteins having different evolutionary distance to *D. melanogaster*. Out of 41 010 *D. melanogaster* introns listed in the APPRIS annotation, 21 562 (52.6%) are located in regions encoding conserved protein domains.

Use of multiple homologous proteins proved to be important for keeping decent accuracy of prediction with increase of evolutionary distance between species with known genomes and the species of interest. Particularly, due to corroboration of footprints originated from multiple homologous proteins, we observed enrichment of high-confidence introns in regions coding for conserved domains (Table 6).

Use of anchored elements of gene structure was important for integration of signals originated from different sources (sites predicted from genomic sequence alone and sites identified by protein footprints). The logic of selection of anchored elements enabled filtering out ‘one-sided’ noises present in one or another source. Use of partial protein footprints, when a target protein mapping could contribute less than full exon–intron structure, was another important feature of the new method. Partial footprints were useful for improving training sets; they also added confident corrections at gene prediction steps (Supplementary Figure S8).

Use of anchored elements was most beneficial for large genomes (*S. lycopersicum* and *D. rerio*) where GeneMark-ES alone generated an elevated rate of random false positive errors within long intergenic regions.

Mapping of N- and C-terminals of target proteins allowed for better discrimination between introns and intergenic regions than it could be done by an *ab initio* algorithm. This improvement led to significant reduction of errors in gene merging (when intergenic regions were predicted as introns) though reduction in error rate of gene splitting (when introns were predicted as intergenic regions) was smaller (Table 4).

The most significant improvement in comparison with GeneMark-ES, observed in all species but fungi, *N. crassa*, occurred when GeneMark-EP+ used the largest possible set of reference proteins (Figure 5, Supplementary Figure S5). For *N. crassa*, use of protein evidence never led to noticeable improvement over GeneMark-ES whose high accuracy for fungal genomes was demonstrated earlier as well (21). We assume that relative drop in GeneMark-EP and -EP+ performance for *C. elegans* in comparison with *Arabidopsis* and *Drosophila* was related to a lower number of reference proteins within the *C. elegans* phylum. In tomato and fish genomes that have longer on average intergenic regions than other species, we saw low exon-level specificity (~55–60%) related to elevated false positive prediction of protein-coding genes in long intergenic regions (Supplementary Figure S5). Gene-level accuracy for *D. rerio*, ~30% Sn and ~12% Sp, for any set of reference proteins beyond

the *D. rerio* genus, was difficult to improve. Notably, the genes in fish genome have a rather large, 8.2, average number of introns per gene. Under independence of error assumption, a gene with a large number of introns would be improbable target for accurate prediction. Even though the independence assumption does not hold in the presence of external evidence, the gene error rate increases with the increase in number of introns (data not shown).

Annotation of genes encoding principal protein isoforms is available for *D. melanogaster*, *C. elegans* and *D. rerio* in the APPRIS database (24). GeneMark-EP+ comparison with the APPRIS annotation shows better Sn than in comparison with annotations containing all possible isoforms.

A question could be raised, how pseudogenes affect training of GeneMark-EP and -EP+. This question is difficult to address in a general setting. Still, since pseudogenes could have different ages, let us consider just groups of ‘young’ and ‘old’. Young pseudogenes with one or two mutations that make them dysfunctional still have all the sequence patterns that could be used in training. Old pseudogenes that accumulated many mutations would harm statistical models if included in training. We argue that old pseudogenes will not be predicted by GeneMark-ES in the course of self-training and therefore they have little or no chance to be included in a training set of anchored elements. On the other hand, elements of young pseudogenes could be identified by GeneMark-ES while the frameshifted exons from spliced alignments will be detected and scored unfavorably by ProHint. Therefore, the young pseudogenes could contribute to parameter training as their ‘intact’ parts will appear in both training and prediction. Addressing full complexity of this issue goes beyond the scope of this project; therefore, currently, GeneMark-EP and -EP+ do not collect information on frameshifts and potential pseudogenes.

Interestingly, the second run of full GeneMark-EP+ (when we took as seed genes the results of gene predictions made in the first full run) had a small but positive effect on the final gene prediction accuracy. This additional run is recommended if an increase in run-time is not a concern.

Running GeneMark-EP and -EP+ requires a protein database as well as tools searching for target proteins and for protein spliced alignments. We used OrthoDB as a database of reference proteins, DIAMOND (25) for the database search for proteins (targets) homologous to the seed proteins and Spaln (9) for spliced alignment of target proteins to genome. To accelerate the pipeline run, we limited the DIAMOND output by 25 target proteins per seed protein (Supplementary Figure S9); choice of Spaln

was also practical from the standpoint of run-time reduction. Additionally, we verified that use of GeneMark-ES for generating seeds was a faster and more efficient method in comparison with the six-frame translation with Procomp and ProSplign tools (8).

This discussion section would be incomplete if we do not mention limitations of the new method. GeneMark-EP does not support multiple models needed for genomes with heterogeneous nucleotide composition, like genomes of mammals and some plants (grasses, e.g. rice). While the current version of GeneMark-EP and -EP+ would outperform GeneMark-ES when running on such genomes, the overall accuracy could be significantly improved with more accurate modeling of genome heterogeneity.

We realize that use of taxonomic divisions for selecting or out-selecting of reference proteins is just the first step in accurate modeling of real-life distributions of evolutionary distances to database orthologs for genes and proteins existing in a novel species. Arguably, there is room for improvement of both intron and gene start/stop hints when modeling of sets of reference proteins is done based on evolutionary distance measures. Similarly, one would expect that such modeling would lead to improvement in selecting thresholds for introns and site mapping.

Another limitation of the current method is the search for a single optimal genomic sequence parse that leads to prediction of a single gene and a single protein isoform in each locus. Importance of genes with alternative splicing has been debated recently, as the evidence was accumulated that alternative splicing mainly operates with UTR regions rather than with translated regions of pre-mRNA. Moreover, the claims were made that when a translated region could be alternatively spliced, then only one among the protein isoforms, the major one, is expressed in the largest number of tissues (24). If gene prediction by GeneMark-EP and -EP+ is viewed as prediction of the major isoform, then the result should be naturally assessed in comparison with annotation of the major isoforms. Such comparison, done for *C. elegans*, *D. melanogaster* and *D. rerio*, used annotation provided by the APPRIS database, and showed improved sensitivity in predicting genes of major protein isoforms. Nonetheless, general tools able to predict all alternative isoforms are of significant interest for community. When external information representing alternative isoforms is provided at RNA level, an earlier developed pipeline, BRAKER1 (31), uses GeneMark-ET and AUGUSTUS to make predictions of alternative isoforms. A new pipeline, BRAKER2 (paper in preparation), combines GeneMark-EP and -EP+ with AUGUSTUS to identify a set of alternative protein isoforms when alternative variants of cross-species proteins are given among references. A new tool, GeneMark-ETP, will combine into gene prediction protein and transcript data (paper in preparation).

DATA AVAILABILITY

Full GeneMark-EP and -EP+ package, including ProHint, is available at http://topaz.gatech.edu/GeneMark/license_download.cgi. Software is compiled for Linux and Mac OS operating systems. All scripts and data used to generate figures and tables in this manuscript are available

at <https://github.com/gatech-genemark/GeneMark-EP-ProtHint-exp>. To give an example, the overall run-time of ProtHint and GeneMark-EP in -EP+ mode on the *D. melanogaster* genome (having ~14 000 genes in 134 Mb sequence) with target proteins selected from species outside Drosophilidae family was ~5 h on 8 CPU/8 GB RAM machine. In our experiments, the run-time grew linearly with respect to both genome length and number of genes.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

FUNDING

National Institutes of Health [GM128145 to M.B.]. Funding for open access charge: National Institutes of Health [GM128145].

Conflict of interest statement. None declared.

REFERENCES

- Hoff, K.J. and Stanke, M. (2019) Predicting genes in single genomes with AUGUSTUS. *Curr. Protoc. Bioinformatics*, **65**, e57.
- Lomsadze, A., Burns, P.D. and Borodovsky, M. (2014) Integration of mapped RNA-Seq reads into automatic training of eukaryotic gene finding algorithm. *Nucleic Acids Res.*, **42**, e119.
- Foissac, S., Gouzy, J., Rombauts, S., Mathe, C., Amselem, J., Sterck, L., Van de Peer, Y., Rouze, P. and Schiex, T. (2008) Genome annotation in plants and fungi: EuGene as a model platform. *Curr. Bioinformatics*, **3**, 87–97.
- Sallet, E., Gouzy, J. and Schiex, T. (2019) EuGene: an automated integrative gene finder for eukaryotes and prokaryotes. *Methods Mol. Biol.*, **1962**, 97–120.
- Behr, J., Bohnert, R., Zeller, G., Schweikert, G., Hartmann, L. and Rättsch, G. (2010) Next generation genome annotation with mGene.ngs. *BMC Bioinformatics*, **11**, O8.
- Birney, E., Clamp, M. and Durbin, R. (2004) GeneWise and Genomewise. *Genome Res.*, **14**, 988–995.
- Gremme, G., Brendel, V., Sparks, M.E. and Kurtz, S. (2005) Engineering a software tool for gene structure prediction in higher organisms. *Inform. Software Technol.*, **47**, 965–978.
- Kiryutin, B., Souvorov, A. and Tatusova, T. (2007) *11th Annual International Conference in Research in Computational Molecular Biology*. San Francisco.
- Gotoh, O. (2008) Direct mapping and alignment of protein sequences onto genomic sequence. *Bioinformatics*, **24**, 2438–2444.
- Keller, O., Kollmar, M., Stanke, M. and Waack, S. (2011) A novel hybrid gene prediction method employing protein multiple sequence alignments. *Bioinformatics*, **27**, 757–763.
- Keilwagen, J., Wenk, M., Erickson, J.L., Schattat, M.H., Grau, J. and Hartung, F. (2016) Using intron position conservation for homology-based gene prediction. *Nucleic Acids Res.*, **44**, e89.
- Gotoh, O., Morita, M. and Nelson, D.R. (2014) Assessment and refinement of eukaryotic gene structure prediction with gene-structure-aware multiple protein sequence alignment. *BMC Bioinformatics*, **15**, 189.
- Burge, C. and Karlin, S. (1997) Prediction of complete gene structures in human genomic DNA. *J. Mol. Biol.*, **268**, 78–94.
- Lukashin, A.V. and Borodovsky, M. (1998) GeneMark.hmm: new solutions for gene finding. *Nucleic Acids Res.*, **26**, 1107–1115.
- Stanke, M. and Waack, S. (2003) Gene prediction with a hidden Markov model and a new intron submodel. *Bioinformatics*, **19**(Suppl. 2), ii215–ii225.
- Parra, G., Blanco, E. and Guigo, R. (2000) GeneID in *Drosophila*. *Genome Res.*, **10**, 511–515.
- Souvorov, A., Kapustin, Y., Kiryutin, B., Chetvernin, V., Tatusova, T. and Lipman, D. (2010) Gnomon: NCBI eukaryotic gene prediction tool. *National Center for Biotechnology Information*. Bethesda.

18. Haas,B.J., Salzberg,S.L., Zhu,W., Pertea,M., Allen,J.E., Orvis,J., White,O., Buell,C.R. and Wortman,J.R. (2008) Automated eukaryotic gene structure annotation using EVidenceModeler and the program to assemble spliced alignments. *Genome Biol.*, **9**, R7.
19. Aken,B.L., Ayling,S., Barrell,D., Clarke,L., Curwen,V., Fairley,S., Fernandez Banet,J., Billis,K., Garcia Giron,C., Hourlier,T. *et al.* (2016) The Ensembl gene annotation system. *Database*, **2016**, baw093.
20. Lomsadze,A., Ter-Hovhannisyan,V., Chernoff,Y.O. and Borodovsky,M. (2005) Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res.*, **33**, 6494–6506.
21. Ter-Hovhannisyan,V., Lomsadze,A., Chernoff,Y.O. and Borodovsky,M. (2008) Gene prediction in novel fungal genomes using an *ab initio* algorithm with unsupervised training. *Genome Res.*, **18**, 1979–1990.
22. Parra,G., Bradnam,K. and Korf,I. (2007) CEGMA: a pipeline to accurately annotate core genes in eukaryotic genomes. *Bioinformatics*, **23**, 1061–1067.
23. Kriventseva,E.V., Kuznetsov,D., Tegenfeldt,F., Manni,M., Dias,R., Simao,F.A. and Zdobnov,E.M. (2019) OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic Acids Res.*, **47**, D807–D811.
24. Rodriguez,J.M., Rodriguez-Rivas,J., Di Domenico,T., Vazquez,J., Valencia,A. and Tress,M.L. (2018) APPRIS 2017: principal isoforms for multiple gene sets. *Nucleic Acids Res.*, **46**, D213–D217.
25. Buchfink,B., Xie,C. and Huson,D.H. (2015) Fast and sensitive protein alignment using DIAMOND. *Nat. Methods*, **12**, 59–60.
26. Borodovsky,M. and Mcininch,J. (1993) GenMark: parallel gene recognition for both DNA strands. *Comput. Chem.*, **17**, 123–133.
27. Camacho,C., Coulouris,G., Avagyan,V., Ma,N., Papadopoulos,J., Bealer,K. and Madden,T.L. (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.
28. Marchler-Bauer,A., Bo,Y., Han,L., He,J., Lanczycki,C.J., Lu,S., Chitsaz,F., Derbyshire,M.K., Geer,R.C., Gonzales,N.R. *et al.* (2017) CDD/SPARCLE: functional classification of proteins via subfamily domain architectures. *Nucleic Acids Res.*, **45**, D200–D203.
29. Stanke,M., Bruhn,W., Becker,F. and Hoff,K.J. (2019) VARUS: sampling complementary RNA reads from the Sequence Read Archive. *BMC Bioinformatics*, **20**, 558.
30. Leinonen,R., Sugawara,H., Shumway,M. and International Nucleotide Sequence Database Collaboration (2011) The Sequence Read Archive. *Nucleic Acids Res.*, **39**, D19–D21.
31. Hoff,K.J., Lange,S., Lomsadze,A., Borodovsky,M. and Stanke,M. (2016) BRAKER1: unsupervised RNA-Seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics*, **32**, 767–769.