OXFORD

Genome analysis

# Ranked choice voting for representative transcripts with TRaCE

## Andrew J. Olson ⓘ [1,*] and Doreen Ware[1,2]

[1]Plant Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11768, USA and [2]USDA ARS Robert W. Holley Center for Agriculture and Health Cornell University, Ithaca, NY 14853, USA

*To whom correspondence should be addressed.

Associate Editor: Janet Kelso

## Abstract

**Summary:** Genome sequencing projects annotate protein-coding gene models with multiple transcripts, aiming to represent all of the available transcript evidence. However, downstream analyses often operate on only one representative transcript per gene locus, sometimes known as the canonical transcript. To choose canonical transcripts, Transcript Ranking and Canonical Election (TRaCE) holds an 'election' in which a set of RNA-seq samples rank transcripts by annotation edit distance. These sample-specific votes are tallied along with other criteria such as protein length and InterPro domain coverage. The winner is selected as the canonical transcript, but the election proceeds through multiple rounds of voting to order all the transcripts by relevance. Based on the set of expression data provided, TRaCE can identify the most common isoforms from a broad expression atlas or prioritize alternative transcripts expressed in specific contexts.

**Availability and implementation:** Transcript ranking code can be found on GitHub at {{https://github.com/warelab/TRaCE}}.

**Contact:** : olson@cshl.edu

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

## 1 Introduction

Genome sequencing projects often use complex, automated annotation pipelines to build reference sets of gene models. These pipelines mask repeats in the assembled genome, align protein and transcript evidence, and build gene models by aggregating overlapping alignments that adhere to known or inferred splice site patterns (Campbell *et al.*, 2014; Haas *et al.*, 2003; Hoff *et al.*, 2019). Before a project releases a set of high-confidence gene models, additional filtering steps may remove transcript models that lack homology or are subject to non-sense-mediated degradation.

Alternative splicing contributes to the functional diversity of a genome (Black, 2003); and new sequencing technology such as PacBio IsoSeq can capture splice variants at an unprecedented scale (Bruijnesteijn *et al.*, 2018; Wang *et al.*, 2016; Zhang *et al.*, 2019). However, this heightened sensitivity can lead to the detection of transcriptional noise, which can be misreported by gene builders as biologically relevant splice variants. Furthermore, it is possible for partially processed transcripts containing retained introns that neither disrupt the reading frame nor introduce stop codons to be promoted to canonical transcripts (Fig. 1).

Comparative gene tree analysis platforms such as Ensembl Compara (Herrero *et al.*, 2016) operate on a single canonical transcript for each gene locus. In the absence of a curated canonical

transcript, this is usually defined as the longest transcript with the longest translation, but this definition does not necessarily select the best representative transcript for a gene locus. Subsequently developed techniques have defined canonical isoforms based on expression level, sequence conservation, annotation of functional domains or some combination of these features (Li *et al.*, 2014; Pruitt *et al.*, 2002; Rodriguez *et al.*, 2018; The UniProt Consortium *et al.*, 2016). For example, NCBI's RefSeq Select dataset uses an evidence hierarchy to identify a transcript in each protein-coding human and mouse gene model. The Matched Annotation from NCBI and EMBL-EBI (MANE) project has the goal of providing a unified set of human protein-coding gene annotations, but it is not known if and when such efforts will be applied to other species.

We developed Transcript Ranking and Canonical Election (TRaCE) to choose canonical transcripts based on data typically available at the time of a new genome annotation. In this approach, transcripts are ranked by length, domain coverage, and how well they represent a diverse population of transcriptome RNA-seq data. An 'election' based on ranked-choice voting selects a canonical transcript that is the first- or second-choice transcript for the majority of samples. The election proceeds through multiple rounds, effectively sorting all transcripts by relevance. Here we present the TRaCE algorithm and results obtained by running TRaCE on *Zea mays* and *Homo sapiens* gene annotations. In addition, we describe validation
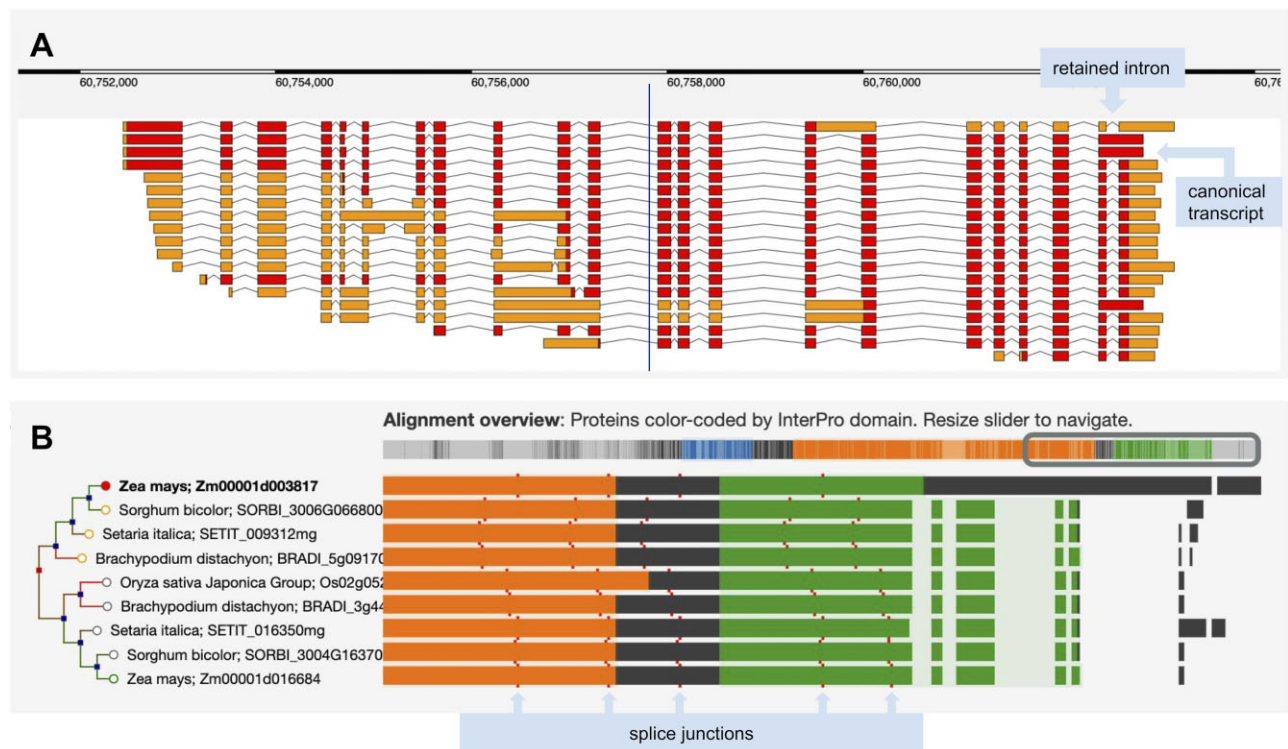
**Fig. 1.** (**A**) The complex set of transcript models for the Zea mays B73 gene sbe4 (starch branching enzyme4). Red blocks show the predicted coding regions, and orange blocks are untranslated regions. The longest translation contains a retained intron and was selected as the canonical transcript for Compara gene tree analysis. (**B**) The left side shows a portion of the gene tree focused on this maize gene and displaying homologs from *Sorghum bicolor*, *Setaria italica*, *Brachypodium distachyon* and *Oryza sativa* Japonica. The right side shows regions of protein sequences participating in the multiple sequence alignment, color coded by InterPro domain. The first row shows a unique region relative to other species that derives from the retained intron

of TRaCE predictions by manual curation (Tello-Ruiz *et al.*, 2019) and compare TRaCE to RefSeq/MANE Select and APPRIS (Rodriguez *et al.*, 2018) human transcript classifications.

## 2 Materials and methods

The first step in preparing to run TRaCE is to gather a diverse set of RNA-seq expression data covering a wide variety of tissues or conditions to act as 'voters' in the upcoming elections. The next step is to align the reads, assemble sample-specific transcripts, and quantify their expression. Each reference gene model with multiple transcripts (candidates) will hold an election to sort the reference transcripts by relevance (Fig. 2).

In each election, samples rank the candidate transcripts based on the annotation edit distance (AED) to the most highly expressed overlapping sample-specific transcripts (Eilbeck *et al.*, 2009). AED scores range from 0 (perfect agreement) to 1 (no overlap) and are calculated from the pairwise similarity of reference transcripts and aligned evidence based on the proportion of exonic overlap. Because there may be insufficient data to assemble full-length transcripts from samples in which the gene is expressed at low levels, the AED score calculation is restricted to overlapping portions of candidate transcripts. A maximum AED score cutoff (default, 0.5) prevents samples from voting for candidate transcripts with very little similarity. There are also cutoff parameters for minimum expression level (default TPM, 0.5) and proportion overlapping (default, 0.5) to filter out some noise in the sample transcriptome data. The election includes additional voters that rank transcripts based on domain coverage, protein length and transcript length. To avoid overwhelming the length-based voters when running TRaCE with many samples, sample votes are weighted to balance the electorate. Default weights were selected to prioritize functional domain coverage over protein length and total transcript length.

Once each sample voter and the length-based voters have ranked the transcripts, the election proceeds in multiple rounds selecting winners until no candidates remain. In each round, TRaCE tallies votes for top-ranked candidates; and so long as there is a tie for first place, votes for the subsequent rankings are added to the tally.

## 3 Results

We ran TRaCE on a pre-release set of *Zea mays* B73 gene models with the set of 10 RNA-seq samples that had already been aligned to the genome as part of the evidence-based gene annotation pipeline (Hufford *et al.*, 2021). The samples were derived from shoot, root, embryo, endosperm, ear, tassel, anther and three leaf sections (base, middle and tip). StringTie version 1.3.5 (with the –rf flag) was used for transcript assembly and quantification (Pertea *et al.*, 2016) and InterProScan version 5.38-76.0 was run to identify Pfam domains (Mulder and Apweiler, 2007). The *Zea mays* B73 V5 annotation set (Zm00001eb) has 15 162 multi-transcript protein-coding gene models; for 5616 of these (37%), the canonical transcript chosen by TRaCE was not the longest isoform. TRaCE selected canonical transcripts for the genome annotations of 25 additional maize accessions, 33–38% of which were not the longest isoform (Supplementary Table S1).

We used two approaches to validate TRaCE's predictions on maize genes. First, we modified an interactive gene tree viewer, designed to flag problematic gene models by visual inspection of the multiple sequence alignment and domain annotations (Tello-Ruiz *et al.*, 2021). We used this interface to compare maize B73 V5 canonical transcripts (Zm00001eb) selected by TRaCE with the prior set of maize V4 canonical transcripts (Zm00001d) selected by length criteria alone. A random selection of 173 pairs of genes for which the TRaCE canonical was not the longest transcript were evaluated in the gene tree viewer and flagged if the alignment was inconsistent with outgroup orthologs. Genes were flagged if there was a relative gain or loss of conserved sequence within the transcript or at either
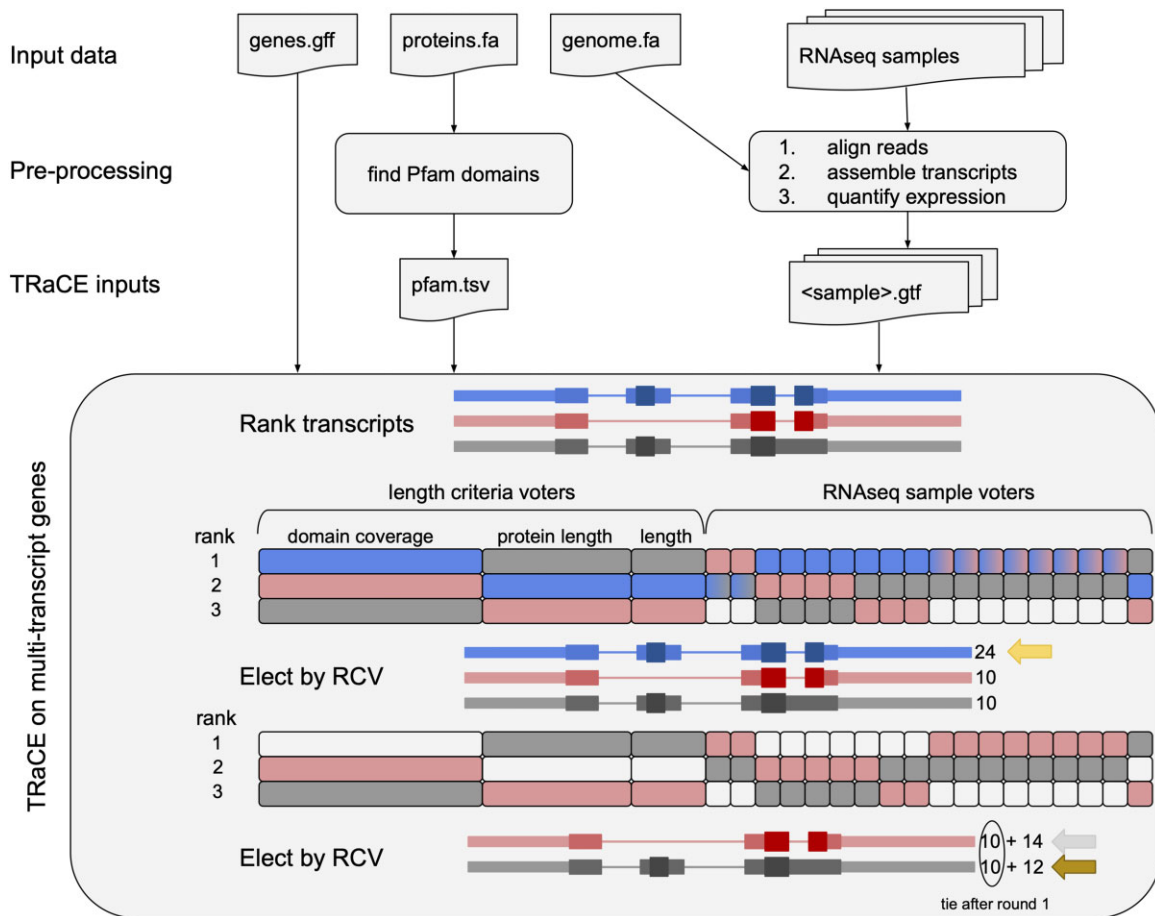
**Fig. 2.** Flowchart of preparation of TRaCE inputs and a schematic of the rank-choice voting (RCV) approach to select transcripts for an example gene with three transcripts (blue, red, gray). Exon thickness corresponds to non-coding, coding and functional regions with Pfam domains. Voters are represented by rectangles, and rank transcripts by length criteria (9, 6 or 3 votes) or AED (1 vote per sample). Eight of the samples rank the red and blue transcripts equally (blue-red gradient), so both get tallied in round 1. RCV selects the blue transcript first with 24 rank 1 votes. After removing the blue votes from consideration, the red and gray transcripts tie with 10 rank 1 votes, but the red transcript is elected with 14 rank 2 votes

end. Of these gene pairs, 32% were flagged as problematic in Zm00001d only, 4% in Zm00001eb only and 5% in both versions (Supplementary Table S2). The most common issue in the flagged Zm00001d gene models was gain of sequence due to an intron retention. Thus, according to this approach, TRaCE was selecting better-conserved isoforms than the prior length-based algorithm.

In the second approach, TRaCE predictions were validated by student curators who were given a subset of 48 gene models with two to five transcripts, for which TRaCE's top-ranked isoform was not the longest isoform. The students, who were not aware of TRaCE's output, were asked to rate transcripts as best, good or poor, based on viewing the gene structure and expression evidence in the Apollo genome browser (Dunn *et al.*, 2019). Each gene model was curated by at least three different students. The transcript ratings were mapped to a score (best 2, good 1, poor -1). Transcript rankings from TRaCE and rankings based on length alone were compared to rankings based on curator scores. For each rank (1–5), we calculated the sum of the curator scores for the associated transcripts. The correlation of these sums between the length-based ranking and the curator-based ranking was 0.917, whereas the TRaCE and curator ranking sums had a higher correlation coefficient of 0.985 (Supplementary Table S3).

We also ran TRaCE on human GRCh38 annotations (Frankish *et al.*, 2019) with a diverse panel of 127 samples of human RNA-seq data covering the development of seven major organs (brain, cerebellum, heart, kidney, liver, ovary and testis) from 4 weeks post-conception to adulthood (https://www.ebi.ac.uk/gxa/experiments/E-MTAB-6814/Results). Reads were aligned with hisat2 version 2.1.0

(–dta –reorder), transcripts were assembled and quantified with stringtie version 2.1.4 (–conservative) and protein-coding reference transcripts were annotated with Pfam domains using InterProScan version 5.38-76.0 (Mulder and Apweiler, 2007; Pertea *et al.*, 2016).

The GRCh38 annotation set has 13 848 multi-transcript protein-coding gene models that were classified by both APPRIS and MANE Select. The TRaCE canonical was not the longest isoform in 3717 (27%) of these gene models. For comparison, the principal isoform according to APPRIS and the MANE Select transcript was not the longest isoform in 3061 (22%) and 4292 (31%) of gene models, respectively. There are 1202 gene models where APPRIS and MANE Select disagree. In these cases, TRaCE agrees with APPRIS on 408 (34%) genes, MANE Select on 597 (50%) genes and neither on 197 (16%) genes. On the 12 646 multi-transcript gene models where APPRIS and MANE Select agree, TRaCE gives 10 677 (84%) transcripts rank 1, 1470 (12%) rank 2, 351 (3%) rank 3 and 148 (1%) rank 4 or higher. To assess TRaCE's performance on gene models with many transcripts, we compared TRaCE to APPRIS and MANE Select on the 90% of genes with 2–10 transcripts and the remaining 10% of human protein-coding gene models with 11–151 transcripts. There are 1399 genes with many transcripts where APPRIS and MANE Select agree. In these cases, TRaCE selects 1021 (73%) of these as the canonical transcript, 215 (15%) have rank 2, 92 (7%) have rank 3 and 71 (5%) have rank 4 or higher. On the 11 247 genes with fewer transcripts where APPRIS and MANE Select agree TRaCE assigns 9656 (86%) rank 1, 1255 (11%) rank 2, 259 (2%) rank 3 and 84 (1%) rank 4 or higher. For the initial release of TRaCE, we manually tuned the weights on TRaCE's length-based

votes, but future versions may benefit from an automated parameter sweep to minimize these differences.

## References

Black,D.L. (2003) Mechanisms of alternative pre-messenger RNA splicing. *Annu. Rev. Biochem.*, **72**, 291–336.

Bruijnesteijn,J. *et al.* (2018) Human and rhesus macaque haplotypes defined by their transcriptomes. *J. Immunol.*, **200**, 1692–1701.

Campbell,M.S. *et al.* (2014) Genome annotation and curation using MAKER and MAKER-P. *Curr. Protoc. Bioinf.*, **48**, 4.11.1–4.11.39.

Dunn,N.A. *et al.* (2019) Apollo: democratizing genome annotation. *PLoS Comput. Biol.*, **15**, e1006790.

Eilbeck,K. *et al.* (2009) Quantitative measures for the management and comparison of annotated genomes. *BMC Bioinformatics*, **10**, 67.

Frankish,A. *et al.* (2019) GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.*, **47**, D766–73.

Haas,B.J. *et al.* (2003) Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res.*, **31**, 5654–5666.

Herrero,J. *et al.* (2016) Ensembl comparative genomics resources. *Database*, **2016**, baw053.

Hoff,K.J. *et al.* (2019) Whole-genome annotation with BRAKER. *Methods Mol. Biol.*, **1962**, 65–95.

Hufford,M.B. *et al.* (2021) De novo assembly, annotation, and comparative analysis of 26 diverse maize genomes. Science, **373**, 655–662.

Li,H.-D. *et al.* (2014) Revisiting the identification of canonical splice isoforms through integration of functional genomics and proteomics evidence. *Proteomics*, **14**, 2709–2718.

Mulder,N. and Apweiler,R. (2007) InterPro and InterProScan: tools for protein sequence classification and comparison. *Comparat. Genomics*, **396**, 59–70.

Pertea,M. *et al.* (2016) Transcript-level expression analysis of RNA-Seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.*, **11**, 1650–1667.

Pruitt,K. *et al.* (2002) The Reference Sequence (RefSeq) database. 2002 Oct 9 [Updated 2012 Apr 6]. In: McEntyre J., Ostell J. (eds.) *The NCBI Handbook [Internet]*, **Chapter 18**. National Center for Biotechnology Information (US), Bethesda, MD. Available from: https://www.ncbi.nlm.nih.gov/books/NBK21091/ (2 April 2021, date last accessed).

Rodriguez,J.M. *et al.* (2018) APPRIS 2017: principal isoforms for multiple gene sets. *Nucleic Acids Res.*, **46**, D213–D217.

Tello-Ruiz,M.K. *et al.* (2019) Double triage to identify poorly annotated genes in maize: the missing link in community curation. *PLoS One*, **14**, e0224086.

Tello-Ruiz,M.K. *et al.* (2021) Gramene 2021: harnessing the power of comparative genomics and pathways for plant research. *Nucleic Acids Res.*, **49**, D1452–D1463.

The UniProt Consortium. *et al.* (2016) UniProt: the universal protein knowledgebase. *Nucleic Acids Res.*, **45**, D158–D169.

Wang,B. *et al.* (2016) Unveiling the complexity of the maize transcriptome by single-molecule long-read sequencing. *Nat. Commun.*, **7**, 11708.

Zhang,G. *et al.* (2019) PacBio full-length cDNA sequencing integrated with RNA-Seq reads drastically improves the discovery of splicing transcripts in rice. *Plant J. Cell Mol. Biol.*, **97**, 296–305.