

SOFTWARE

Open Access



ClusTrast: a short read de novo transcript isoform assembler guided by clustered contigs

Karl Johan Westrin¹, Warren W. Kretzschmar^{1,2} and Olof Emanuelsson^{1*}

*Correspondence:
olofem@kth.se

¹ Science for Life Laboratory,
Department of Gene
Technology, KTH Royal Institute
of Technology, 171 65 Solna,
Sweden

² Department
of Medicine Huddinge,
Center for Hematology
and Regenerative Medicine
(HERM), Karolinska Institute, 141
52 Flemingsberg, Sweden

Abstract

Background: Transcriptome assembly from RNA-sequencing data in species without a reliable reference genome has to be performed de novo, but studies have shown that de novo methods often have inadequate ability to reconstruct transcript isoforms. We address this issue by constructing an assembly pipeline whose main purpose is to produce a comprehensive set of transcript isoforms.

Results: We present the de novo transcript isoform assembler ClusTrast, which takes short read RNA-seq data as input, assembles a primary assembly, clusters a set of guiding contigs, aligns the short reads to the guiding contigs, assembles each clustered set of short reads individually, and merges the primary and clusterwise assemblies into the final assembly. We tested ClusTrast on real datasets from six eukaryotic species, and showed that ClusTrast reconstructed more expressed known isoforms than any of the other tested de novo assemblers, at a moderate reduction in precision. For recall, ClusTrast was on top in the lower end of expression levels (<15% percentile) for all tested datasets, and over the entire range for almost all datasets. Reference transcripts were often (35–69% for the six datasets) reconstructed to at least 95% of their length by ClusTrast, and more than half of reference transcripts (58–81%) were reconstructed with contigs that exhibited polymorphism, measuring on a subset of reliably predicted contigs. ClusTrast recall increased when using a union of assembled transcripts from more than one assembly tool as primary assembly.

Conclusion: We suggest that ClusTrast can be a useful tool for studying isoforms in species without a reliable reference genome, in particular when the goal is to produce a comprehensive transcriptome set with polymorphic variants.

Keywords: De novo transcriptome assembly, Guiding contigs, Isoform assembly, RNA-seq, Recall/sensitivity

Background

In eukaryotes, many genes can produce RNA transcripts of differing base sequences called transcript isoforms. Transcript isoforms are created by alternative transcriptional start sites, splicing, or polyadenylation. Controlling transcript isoform expression is one way for a cell to regulate protein expression and thereby its behavior [1–3]. Changes in transcript isoform expression have been associated with developmental changes and tissue specificity in eukaryotes, and disease in humans [4–7]. Thus, it is often important



to clarify not only what genes are expressed but also which transcript isoforms are expressed.

The expression of genes and transcripts is often studied by RNA-sequencing, where short reads (SRs) derived from massively parallel shotgun sequencing are aligned to an organism's reference genome. With this approach, reconstructing transcripts is possible by using the reference genome as a guide [8]. However, many non-model organisms do not have a high-quality reference genome available. In such cases, a commonly used approach is de novo assembly in which transcripts are assembled from the reads only. The assembled transcripts are sometimes referred to as contigs or reconstructed transcripts. Popular tools to perform de novo transcriptome assembly include Trans-ABYSS [9], Trinity [10], Oases [11], and SOAP-denovo-Trans [12]. An overview of current transcriptome assemblers is available [13]. According to that study, the best performing assemblers were Trans-ABYSS, Trinity and rnaSPAdes [14].

In principle, these tools can also reconstruct transcript isoforms of the expressed genes, but in practice their sensitivity is poor. In *Mus musculus*, Schultz et al. [11] reported that Oases, Trans-ABYSS, and Trinity assembled 1.21, 1.25, and 1.01 transcripts per gene, respectively, whereas a reference-based assembler reconstructed 1.56 transcripts per gene. Bushmanova et al. [14] also observed poor transcript isoform reconstruction performance of transcriptome assembly methods: while their method, rnaSPAdes, outperformed the other compared assemblers in gene reconstruction in *Mus musculus*, it assembled only 1.02 transcripts per gene. In the same comparison, Trinity managed to assemble the most transcripts, with a ratio of 1.11 transcripts per gene. The insufficient ability of current de novo transcriptome assembly approaches to reconstruct all expressed transcript isoforms of a gene was evident to us in our work on the DAL19 gene in spruce, *Picea abies* [7]: Only one out of four confirmed DAL19 transcript isoforms was reconstructed to at least 90% using Oases and two using Trinity. We performed a directed assembly that managed to reconstruct three of the four transcript isoforms, but this method did not scale to whole transcriptome assembly. These examples, and others [15, 16], demonstrate that there is still much room for improvement in de novo transcript isoform assembly.

Another observation concerns the imperfect overlap between the sets of reconstructed transcripts from different de novo assembly tools. Smith-Unna et al. [17] noted that out of Oases, Trinity, and SOAP-denovo-Trans, each assembler reconstructed a large number of *bona fide* transcripts that neither of the other assemblers managed to reconstruct. They concluded that combining assembly methods may be an effective way to improve the detection rate of transcripts.

We report the de novo transcriptome assembler ClusTrast, which builds upon our previous experience of transcript isoform assembly [7]. The main purpose of ClusTrast is to provide a comprehensive set of transcript isoforms, using only sequence reads as input, and with the explicit intent to prioritize recall. The ClusTrast pipeline combines two assembly methods, Trans-ABYSS and Shannon, incorporates a novel approach to clustering guiding contigs, assigns short reads to the clusters, and finally performs a cluster-wise assembly of the clustered short reads. We assessed transcript isoform reconstruction performance of ClusTrast and several de novo transcriptome assemblers in six eukaryotic organisms and found that ClusTrast reconstructed more

known transcript isoforms than any other assembler and reconstructed unknown (including misassembled) transcripts at a rate comparable to other assemblers.

Implementation

ClusTrast method

We developed an approach for transcriptome assembly from short reads called ClusTrast.

Overview

Figure 1 shows a flowchart of the ClusTrast pipeline. The only required input to ClusTrast is a file with short RNA-seq reads, referred to as SRs, short reads or SR RNA-seq. Guiding contigs (GCs) is an optional input. Additional file 1: Fig. S.1 illustrates an example of how the method works.

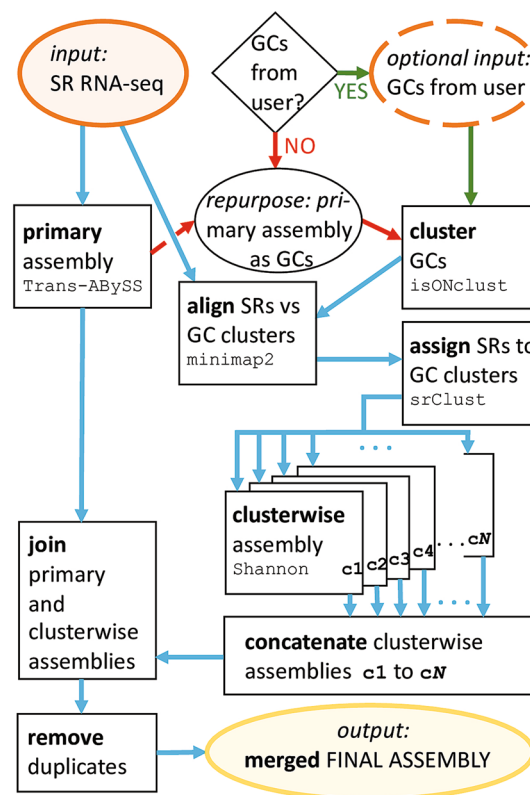


Fig. 1 The ClusTrast pipeline. Orange ovals denote input data, black squares denote actions taken by ClusTrast, blue arrows denote intermediate data file transfers, yellow oval denotes output final assembly, and the black rhomboid and oval handle the guiding contig user input status. The only required input to ClusTrast is a short read RNA-seq data set (orange oval in top left corner). Guiding contigs may optionally be input by the user (dashed orange oval in top right corner, and green arrows followed if provided). They may also be used as primary assembly (not depicted). If guiding contigs are not provided by the user, the primary assembly is repurposed and used also as guiding contigs (red arrows). SR = short reads, GCs = guiding contigs

Primary assembly and guiding contigs (GCs)

In ClusTrast, Trans-ABYSS [9] is employed to create a “primary assembly” from the short reads. The primary assembly will by default be used as the set of guiding contigs (GCs) in ClusTrast. The guiding contigs are used in the next step to assign the short reads to clusters. Guiding contigs may also be provided separately by the user, and could then also serve as primary assembly if desired. The primary assembly is by default merged into the final set of assembled transcripts.

Trans-ABYSS is one of the leading de novo transcriptome assemblers according to Hölzer and Marz [13]. They evaluated the original strategy of Trans-ABYSS, which used several different k -mers and merged the resulting assemblies, in order to get both recall from small values of k and precision from high values of k . Since a single- k run uses much less memory (or is substantially faster) than a multi- k run, we tried both strategies with ClusTrast. In this report, we have appended **-M** to the name of a method if it used a multi- k strategy. We also tried other assemblers as potential primary assemblers for ClusTrast, see *Primary assembly alternatives*.

Clustering of (a) guiding contigs and (b) Assigning short reads to clusters

(a) The clustering of the guiding contigs is performed with isONclust [18], a tool originally developed for clustering long reads (from PacBio or ONT sequencing technologies) into gene families. It uses a greedy algorithm for the clustering and handles variable error rates by the means of the quality values in the FASTQ input files. When the set of guiding contigs is in FASTA-format, ClusTrast will convert it to FASTQ-format with a static quality.

(b) The short reads are aligned to the guiding contigs with minimap2 [19], using the preset option `-x sr`, intended for short read alignment, but included secondary alignments. Secondary alignments can optionally be excluded in ClusTrast. Next, the short reads are assigned to the guiding contig clusters based on the alignment results. If a short read x is aligned to guiding contigs X_1 and X_2 , and X_1 belongs to cluster n_1 and X_2 belongs to cluster n_2 , x will be included in n_1 and n_2 . Reads not aligning to any guiding contig are put in a separate cluster. Thus, the short reads have now been clustered. A read can only occur once per cluster. See also Additional file 1: Section C.2.

Clusterwise assembly

The cluster-wise assembly in ClusTrast is performed by the transcriptome assembler Shannon [20] (also used in refShannon, a genome-guided transcript assembler [21]), and aims to be information theoretically optimal. Kannan et al. claim that Shannon can finish in linear time given (i) sufficient diversity of transcript abundance and (ii) no loops in the graph, but they do not address how it will deal with datasets not meeting these criteria. However, dividing the reads in the short read dataset into clusters before assembly will reduce the complexity of each individual assembly and lower the risk of violating these requirements. Because of this, and its aim to reconstruct as many transcripts as possible, Shannon is used for the cluster-wise assemblies in ClusTrast.

Merging the primary and clusterwise assemblies

The final steps of ClusTrast are to join the concatenated clusterwise assemblies with the primary assembly, and next, to remove duplicate instances of reconstructed transcripts. Similar but non-identical reconstructed transcripts are kept, because of possible polymorphic variants that we want ClusTrast to retain. The output of ClusTrast is the merged final assembly.

Datasets and annotations

Short read RNA-seq datasets for assembly generation

We evaluated the de novo transcriptome assemblies using the NCBI SRA datasets in Table 1. They were all non-stranded paired-end short read RNA-seq datasets. We pre-processed the datasets with fastp [22] with default parameters, which means removal of any remaining adapter sequences, quality pruning (max 40% of the bases were allowed to have base quality < 16, and at most five Ns per read), and exclusion of reads that ended up shorter than 15bp (see the Additional file 1: Sections A.1 and B.1 for details).

Reference datasets for assembly evaluation

We downloaded reference genome sequences as well as reference transcript annotations from Ensembl for each of the six species. We used the GTF file annotations of genes and transcripts, not including the ab initio annotations. We estimated the expression of all reference transcripts in each of the six datasets using RSEM [23] and defined a transcript isoform as expressed if the transcripts per million (TPM) reported by RSEM was greater than zero. Versions and commands for RSEM are listed in Additional file 1: Section A–B. We defined a gene as expressed if at least one of the transcript isoforms associated with that gene was expressed. The versions of the annotations used and the number of genes and isoforms we detected in each dataset are shown in Table 2.

Transcriptome assembly generation

We assembled the transcriptomes for all six datasets (Table 1) using Trans-ABYSS [9], Trinity [10], Oases [11], SOAP-denovo-Trans [12], BinPacker [24], Shannon [20], rnaSPAdes [14], TransLiG [25], RNA-Bloom [26], and ClusTrast. We used each assembler’s own default parameters. Trans-ABYSS and Oases can be run in a “multi- k ” mode where the assembler is first run with a single k -mer (“single- k ” mode; where a k -mer is a substring, with fixed length k , of a read) for several different k -mers and the resulting

Table 1 Short read RNA-seq datasets accessed from the NCBI SRA database

SRA ID	RL	Species	RPs	
SRR5133163.1	2 × 150	Human	29.51	29.05
SRR8632985	2 × 76	Mouse	31	30
SRR11341576	2 × 150	Rice	24	23
SRR11278019	2 × 126	Arabidopsis	11.2	11.1
SRR10728575	2 × 150	Zebrafish	21	20
SRR5986240.1	2 × 150	Poplar	25.1	24.4

RL=read length in bases. Species column, indicated in bold is the name to which the data set is referred to throughout this article. RPs=million read pairs, before pre-processing (on the left) and after pre-processing (on the right)

Table 2 Reference transcriptome sequence, genome sequence and annotation versions accessed from Ensembl, where the Version id suffix shows the Ensembl version

Reference		Total		Expressed	
Species	Version id	Genes	Isoforms	Genes	Isoforms
Human	GRCh38.99	40,491	1,904,32	22,510	1,02,552
Mouse	GRCm38.99	36,711	1,19,353	17,400	52,845
Arabid.	TAIR10.48	27,655	48,359	23,085	38,004
Rice	IRGSP-1.0.48	37,967	44,761	29,703	34,956
Zeb.fish	GRCz11.99	30,628	57,775	23,963	35,189
Poplar	Pop_tri_v3.46	41,335	73,012	29,400	44,652

The number of genes and isoforms are counted from the reference transcriptome. Genes and isoforms are considered expressed if TPM>0 as calculated by RSEM on the datasets in Table 1

assemblies are merged into a single assembly. We used both the single- k and multi- k strategies for these two assemblers. We append **-M** to the name of a method if it uses a multi- k strategy, and **-S** if it uses a single- k strategy. Oases-**M** uses by default all odd k -mers from 19 to 31, but it only finished within less than 58 h on the mouse and arabidopsis datasets. On the rice dataset, it finished after \sim 400 h. Therefore, for human and zebrafish, we used only Oases-**S** with $k = 31$. The program versions and the executed commands are listed in Additional file 1: Sections A.2 and B.2.

We also generated a concatenated assembly from the Trans-ABySS and Shannon transcriptomes, referred to as TrAB+Sh, to examine if the clustering approach of ClusTrust improves the assembly quality.

Transcriptome assembly evaluation

We evaluated the transcriptome assemblies by estimating precision (positive predicted value, PPV) and recall (sensitivity or true positive rate, TPR). For this, we used the reference based transcriptome comparison tools Conditional Reciprocal Best BLAST (CRBB) [27], as implemented in the TransRate package [17], and SQANTI [28]. Versions and commands for these tools can be found in the supplementary material. We only used reference transcripts that were considered expressed (Table 2). All assembled transcripts were considered expressed, since they were reconstructed from actual RNA-seq data.

Using SQANTI in evaluation

We used SQANTI (Structural and Quality Annotation of Novel Transcript Isoforms) [28] to classify assembled transcripts according to their splice junction matches with reference genes and transcript isoforms. When an assembled transcript is anti-sense to an annotated gene, SQANTI will classify that transcript as anti-sense. We extracted all transcripts classified as anti-sense, reverse-complemented them, and then reclassified them with SQANTI.

When an assembled transcript and a reference isoform have the same number of exons and same splice junctions, then SQANTI classifies it as a full splice match (FSM). When the assembled transcript has fewer exons than the reference but the splice junctions in the assembled transcript all exist in the reference, it is classified as an incomplete splice match (ISM). In order for SQANTI to classify an assembled transcript as an ISM, all

junctions in the assembled transcript must match the reference, but the exact start and end can differ. In case there are several possible consistent reference isoforms, SQANTI assigns the assembled transcript to the shortest of the matching references. Assembled transcripts classified by SQANTI as novel in catalog (NIC, when the splice junctions are known but there is a novel combination) and novel not in catalog (NNC, with novel splice junctions) were not classified as true positives.

For recall, we counted all expressed reference isoforms with at least one assembled transcript that SQANTI classified as FSM or ISM (and with a certain fraction, 0.25–1.0, of the exons covered) as a true positive, and divided the total number of true positives by the total number of expressed reference isoforms. For precision, we counted each assembled isoform classified by SQANTI as FSM or ISM and covering at least a certain fraction (from 0.25 to 1.0) of the reference exons as a true positive, and divided the total number of true positives by the total number of assembled isoforms.

Using CRBB in evaluation

We used CRBB [27] to classify assembled transcripts according to their similarity to reference transcripts. To this end, we used TransRate [17], which in turn used BLAST [29] to align each assembled transcript to the set of reference transcripts, and each reference transcript to the set of assembled transcripts. By using all transcripts which are top hits in both BLAST alignments reciprocally, an appropriate E-value cutoff is calculated. Transcripts with lower E-values than this cutoff are then considered CRBB hits. We defined recall as the proportion of reference transcripts that have a CRBB hit covering the reference transcript to at least 25%–100%. We defined precision as the proportion of assembled transcripts that are a CRBB hit covering the reference isoform to at least 25%–100%.

Results

Transcriptome assembly evaluation

Transcriptome assemblies for all compared assemblers, including ClusTrast, were generated as described under *Transcriptome assembly generation*. Basic statistics of all assembled transcriptomes are available in Additional file 1: Table S.2–S.7. We collectively refer to all tested approaches as “assemblers”, although assembly pipeline (e.g., ClusTrast) or concatenation (TrAB+Sh) may be more accurate.

Evaluation with SQANTI

We investigated how recall and precision changed when we varied the proportion of exons that an assembled transcript needs to recover in order to be considered a true positive. As this proportion was relaxed for the ISM classifications from 1.0 to 0.25 (for the FSM category it is by definition 1.0), the recall and precision (Fig. 2) increased. ClusTrast-M had the highest SQANTI recall of any assembler for all of the six datasets over the entire range (except roughly tied with TrAB+Sh for arabidopsis). The assembler with the highest precision varied across datasets; it was RNA-Bloom in human, ClusTrast-M in rice and (with TransLiG) poplar, Oases-M in mouse, and TransLiG in arabidopsis and zebrafish.

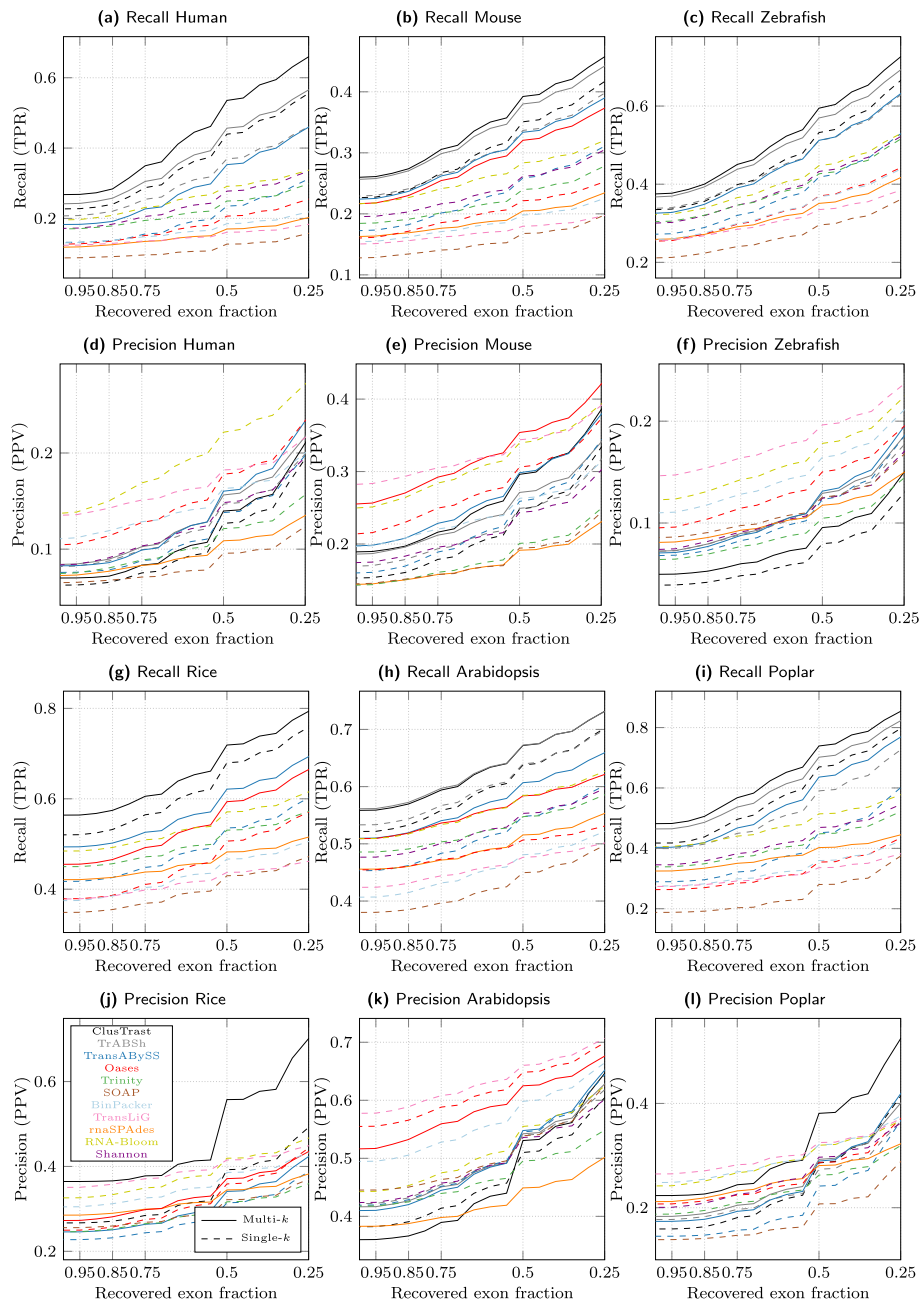


Fig. 2 Recall and Precision as measured by SQANTI. a-c,g-i: Proportion of reference isoforms with at least one SQANTI classification of FSM or ISM vs. the cumulative proportion of exons recovered by the assembly. d-f,j-l: Proportion of reconstructed isoforms classified by SQANTI as FSM or ISM vs. the cumulative proportion of recovered exons from the reference

Fixing the proportion at 0.5 (i.e., at least 50% of exons recovered for ISM), ClusTrast-M detected more transcript isoforms than the established assemblers Trinity (1.23–2.15 fold increase), Oases-S (1.33–2.59 fold increase), and Trans-ABBySS-M (1.1–1.52 fold increase) (Additional file 1: Fig. S.2 and Table S.14). Precision was comparable to Trinity (0.9–1.76 fold change), Oases-S (0.78–1.5 fold change) and Trans-ABBySS-M (0.73–1.63 fold change).

Evaluation with CRBB

We investigated CRBB recall and precision over the same proportion of required recovered exons as for SQANTI and observed an increase in recall and precision as this proportion was decreased from 1.0 to 0.25. We observed some changes in the relative ordering of assemblers as shown in Fig. 3 (CRBB recall and precision). In particular, rnaSPAdes performance levelled off in the lower end.

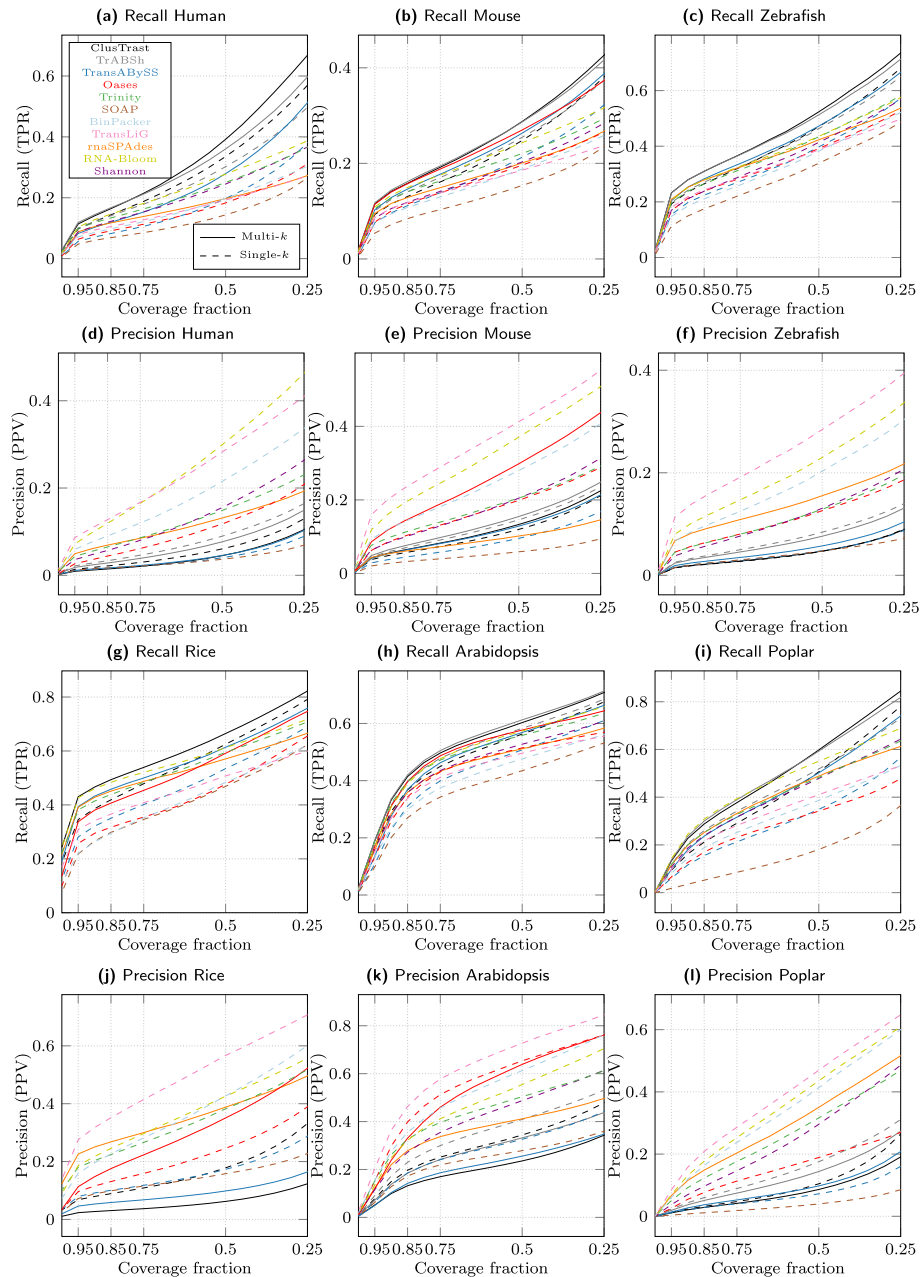


Fig. 3 Recall and Precision as measured by CRBB. a-c-g-i: Proportion of references with a CRBB hit vs. the cumulative proportion of recovered reference length. d-f-j-l: Proportion of reconstructed isoforms with a CRBB hit vs. the cumulative proportion of recovered reference length

Fixing the proportion at 0.5, CRBB recall was higher for ClusTrast-M than for Trinity (1.01–1.34 fold increase) across all datasets, but not compared to all assemblers (Additional file 1: Fig. S.3 and Table S.15). ClusTrast-M performed the best on human, mouse, rice, and zebrafish, it ranked second for arabidopsis, while its ranking varied from first to third as the required reference transcript coverage decreased. ClusTrast-M clearly underperformed compared with Trinity with regard to CRBB precision (0.33–0.68 fold change); the assembler with highest precision in a dataset was always TransLiG or RNA-Bloom.

SQANTI and CRBB evaluation metrics were correlated; true positive sets still differed

The number of transcripts that were considered as true positives by both SQANTI and CRBB or exclusively by only one of them varied between datasets and between assemblers (Additional file 1: Tables S.18 and S.19). ClusTrast, Trans-ABYSS, Oases-M, and (with one exception) Shannon and SOAP-denovo-Trans consistently predicted more transcripts that were considered as true positives exclusively by SQANTI and not by CRBB, while TransLiG was the only assembler that consistently predicted more transcripts that were considered true positives exclusively by CRBB. We used SQANTI categories to classify the ClusTrast true positives that were exclusively detected by either SQANTI or CRBB (Additional file 1: Tables S.23 and S.24, respectively). We observed that the largest category of true positives according to SQANTI but not CRBB was the ISM mono-exon class. The largest category of true positives according to CRBB but not SQANTI was novel not in catalog (NNC) with novel splice sites.

SQANTI and CRBB recall measurements were highly correlated across all assemblies and datasets ($\rho = 0.93$; Additional file 1: Fig. S.4) while SQANTI and CRBB precisions were less correlated ($\rho = 0.75$; Additional file 1: Fig. S.5). We calculated the correlation of precision measurements for each assembler individually: ClusTrast-M obtained $\rho = 0.54$ while for all other assemblers $\rho \geq 0.82$. Next, we excluded the ISM mono-exon class from the set of true positives and recalculated the precision correlation for ClusTrast-M: it increased to $\rho = 0.94$.

Reference transcripts were often covered to at least 95% by FSMs

We investigated the number of expressed reference transcript isoforms that were reconstructed to at least 50% and 95% of their length by a single FSM according to SQANTI, Additional file 1: Table S.16. For all assemblies and both length requirements, either TrAB+Sh or ClusTrast reconstructed the most reference transcript isoforms, with small differences (<5%) except for rice where TrAB+Sh did not produce a result. Between 35.1% (arabidopsis) and 68.8% (rice) of the reference transcript isoforms that had an FSM match were reconstructed by the FSM-classified contig from ClusTrast to at least 95% of their length. The corresponding range for reconstruction to at least 50% of the reference transcript length was between 76.7% (human) and 91.0% (arabidopsis).

An appreciable fraction of reference transcripts were reconstructed with polymorphisms by ClusTrast

We used the subset of reference transcripts with FSM or CRBB hits to estimate how often these reference transcripts were reconstructed as polymorphic variants (SNPs,

indels) or as alternatively spliced contigs. In Additional file 1: Tables S.20 and S.21, the sets labeled *A* contain the FSMs, while the sets labeled *B* contain the CRBB hits. By definition, FSM contigs corresponding to a specific reference transcript are not alternatively spliced, since they contain all splice junctions of their reference transcript. Two (or more) FSM contigs matching one and the same reference transcript are thus polymorphic variants of each other. This is the $A \setminus B$ and $A \cap B$ sets in Additional file 1: Tables S.20 and S.21. On the other hand, two (or more) contigs that are not FSMs but considered as CRBB hits to one and the same reference transcript, are potentially splice variants of that reference transcript. This is the $B \setminus A$ sets. We estimated that 58–81% of the reference transcripts reconstructed by ClusTrast were reconstructed with polymorphic variants, Additional file 1: Table S.22. Conversely, we estimated that 47–78% of ClusTrast assembled contigs contained polymorphic variants, Additional file 1: Table S.22.

Recall varied over expression levels and number of exons in isoforms

To determine if the assemblers differed in how well they recovered isoforms of genes with more than one annotated isoform, we calculated SQANTI recall of isoforms binned by genes according to the number of isoforms these genes expressed (Fig. 4). In most cases the ranking of assemblers by recall did not change with increasing number of expressed isoforms per gene. ClusTrast-M came out on top over almost the entire range for 5 out of 6 datasets, although for mouse it was tied with TrAB+Sh and for arabidopsis it was tied with Oases-M and TrAB+Sh.

We binned reference transcripts by expression quantiles as measured by RSEM. The SQANTI recall increased with increased expression level, for all assemblers and for all data sets except that some assemblers levelled off in the range 80–100%. We observed that recall was higher for ClusTrast-M than all other assemblers in the lower

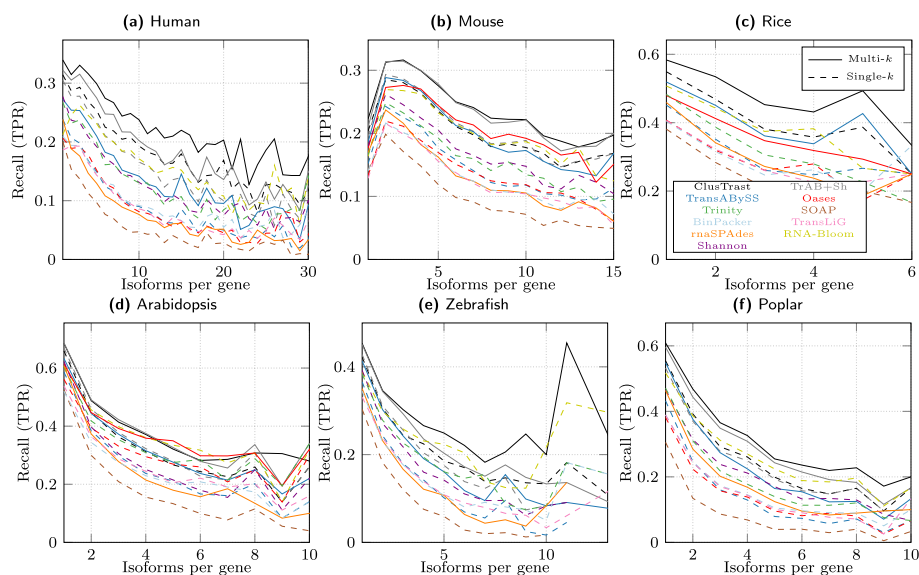


Fig. 4 SQANTI recall of reference isoforms (FSM) binned by number of expressed isoforms per gene

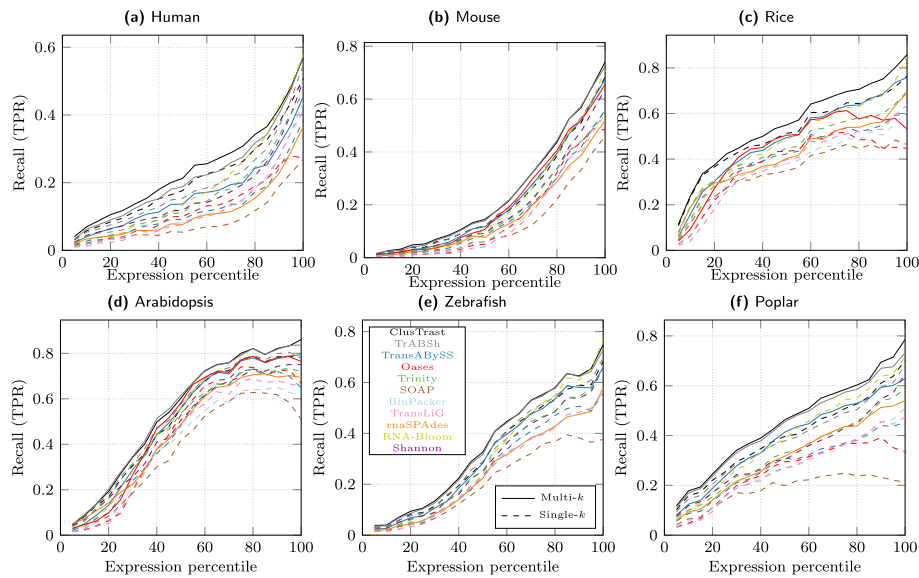


Fig. 5 SQANTI recall of expressed transcript isoforms stratified according to RSEM expression. Recall within each bin (5 percentiles) is defined as the proportion of transcript isoforms that have an FSM match to an assembled contig that is ≥ 200 bp

end of expression levels (expression quantile $<15\%$) and across the entire range of expression levels for all datasets except arabidopsis and zebrafish where ClusTrast-M was tied with TrAB+Sh (Fig. 5).

We observed cases where ClusTrast detected highly expressed isoforms missed by other methods. We illustrate this with examples of genes where the highest expressed isoform (according to RSEM) was reconstructed only by ClusTrast and not by any other method. Additional file 1: Figs. S.26 to S.31 contain Sashimi plots for these example genes, one for each of the six datasets (see Additional file 1: Section C.4 for more details).

Simulated datasets

We evaluated ClusTrast on two simulated datasets: one human dataset from Hölzer and Marz [13] and one mouse dataset from Hayer et al. [15]. The results (Figs. 6 and 7) were analogous to the ones from the real datasets: On the simulated human dataset, ClusTrast was the leading method for recall (according to both SQANTI and CRBB) and the worst for precision (according to SQANTI). On the simulated mouse dataset, it was even between ClusTrast and TrAB+Sh on recall, and a mediocre to low precision for ClusTrast.

Run time and memory usage

Across all real datasets, all assemblers except Oases-M completed within 48 h and required less than 300 GB of memory (Additional file 1: Table S.25). ClusTrast-M took between 660 and 2145 min to complete, the longest time of all assemblers for mouse and arabidopsis. Oases-M took the longest time to complete for one dataset, and did not complete for the remaining three. SOAP-denovo-Trans was the fastest assembler for all datasets. ClusTrast-M peak memory use was between 57.15 and 267.2 GB for the six datasets, highest of all assemblers for one dataset, while Trinity and Oases-M had the

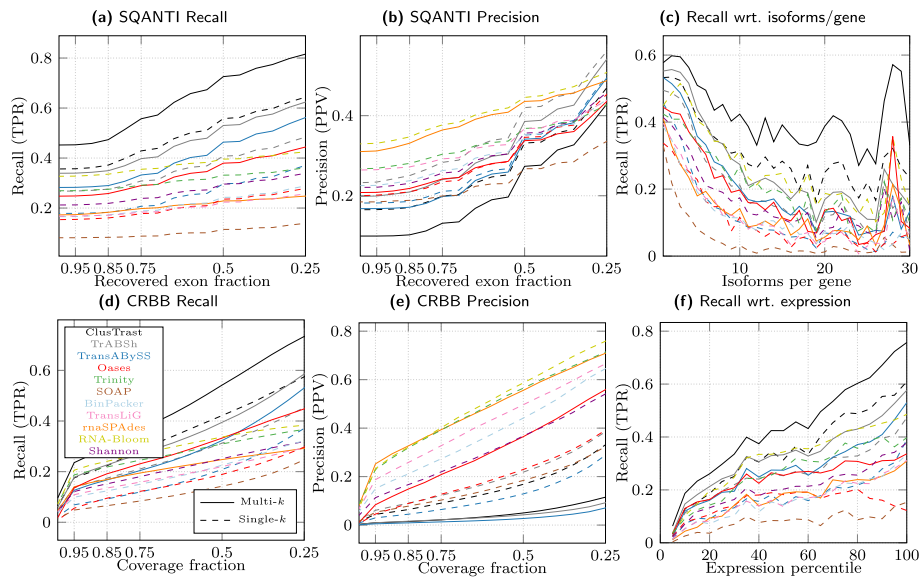


Fig. 6 Results for a simulated human dataset from Hölzer and Marz [13]

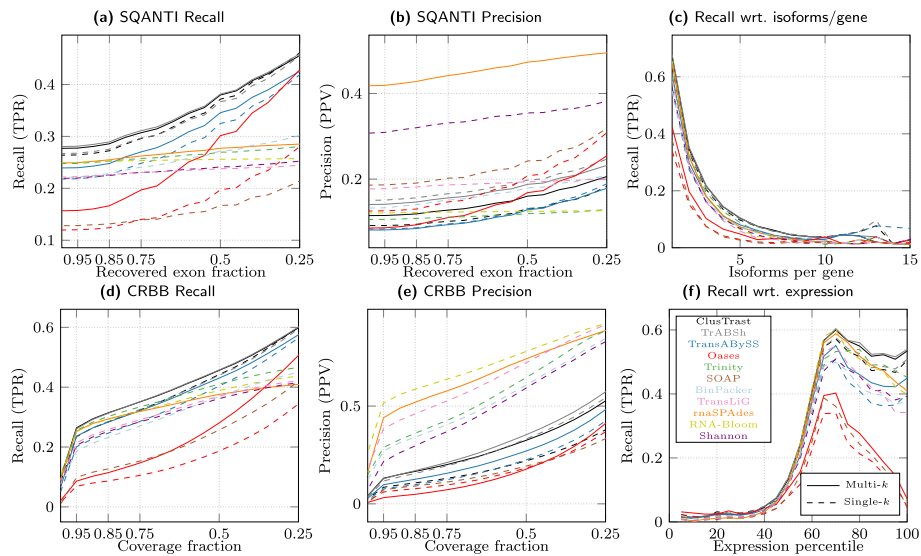


Fig. 7 Results for a simulated mouse dataset from Hayer et al. [15]

highest peak memory use for two datasets each. RNA-Bloom had the lowest peak memory usage. Our computational setup is described in Additional file 1: Table S.26.

Primary assembly alternatives

We assessed alternatives to Trans-ABySS-M for primary assembly in ClusTrast: the widely used Trinity, and the recent RNA-Bloom (both of which faring well individually in our own evaluation), and a “meta assembler”, that we call META, constructed by taking the union of the individual assemblies from Trans-ABySS-M, Trinity, and RNA-Bloom (Figs. 8, 9). ClusTrast was tested in four versions, each with a different primary assembly, Trans-ABySS-M, Trinity, RNA-Bloom, and META (solid lines). Recall was

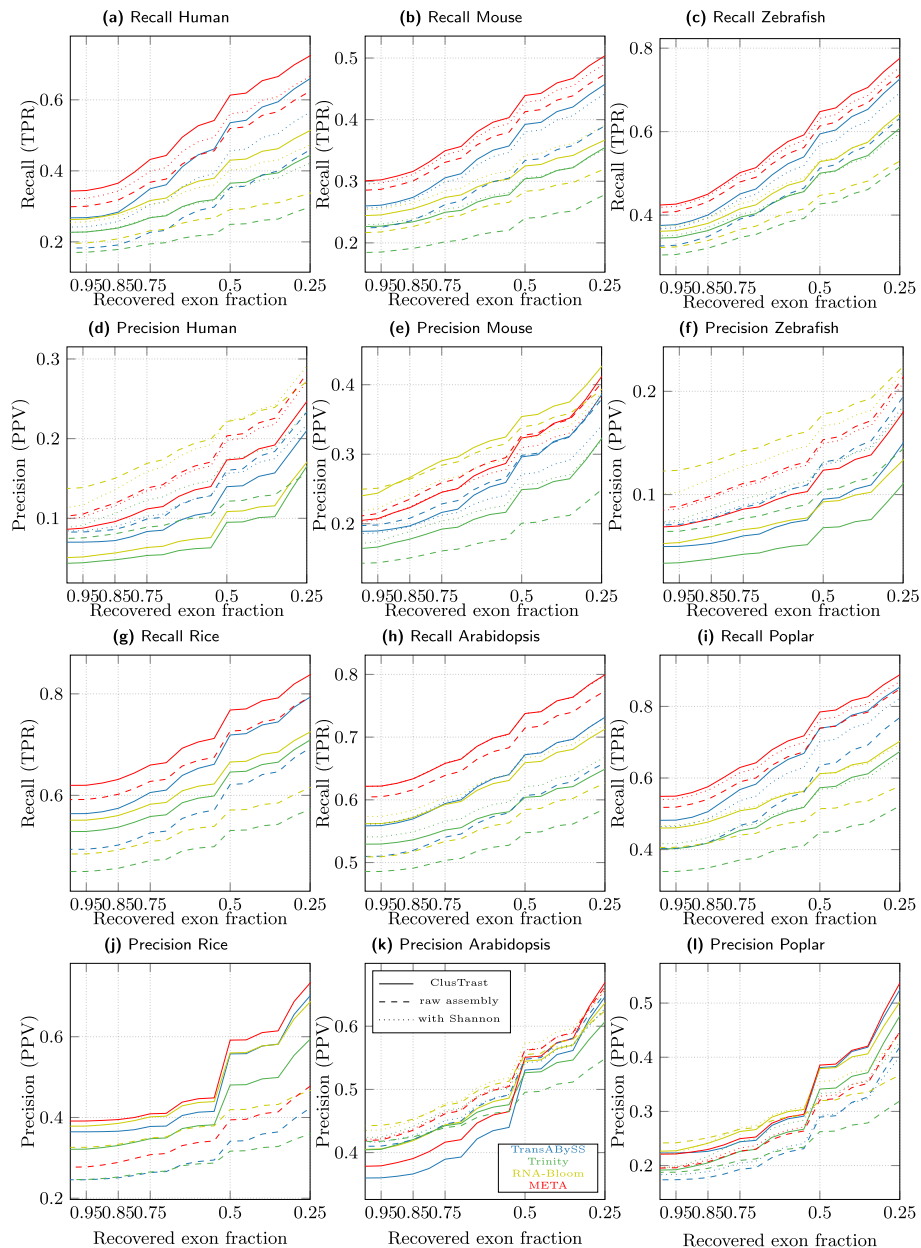


Fig. 8 Recall and Precision as measured by SQANTI for ClusTrast when ran with different tools for primary assembly (solid lines), compared with these tools on their own (dashed lines) and concatenated with Shannon (dotted lines)

improved as compared to the four individual primary assemblies (dashed lines), according to both SQANTI, Fig. 8, and CRBB, Fig. 9, while precision was lower for 4/6 datasets according to SQANTI and overall according to CRBB. We also concatenated the individual assembly from Trans-ABySS-M, Trinity, and RNA-Bloom, respectively, to the individual assembly from Shannon (dotted lines). These concatenated assemblies (where the concatenation of Trans-ABySS-M and Shannon, TrAB+Sh, is included also in Figs. 2, 3, 4, 5, 6 and 7) yielded higher recall for RNA-Bloom and Trinity as compared to what

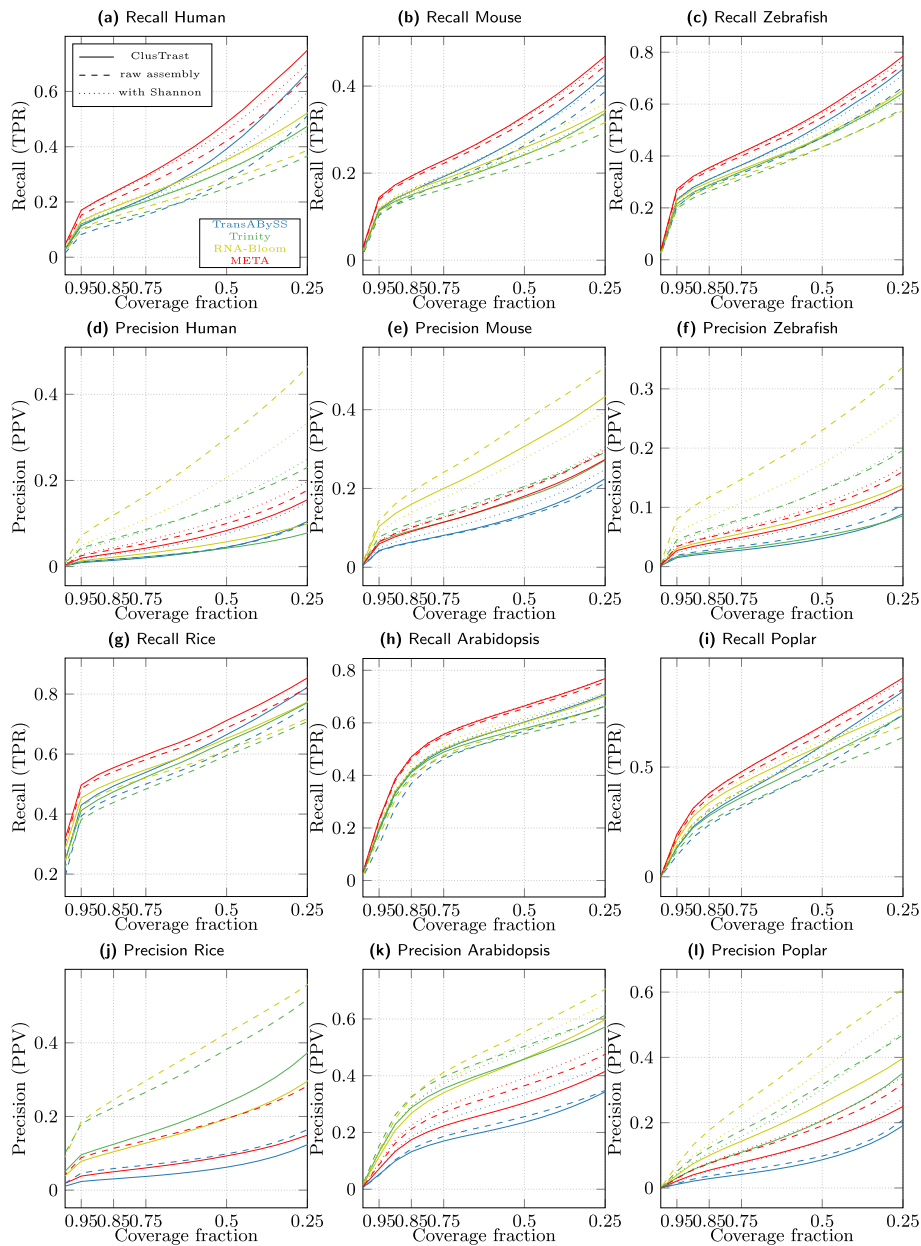


Fig. 9 Recall and Precision as measured by CRBB for ClusTrast when ran with different tools for primary assembly (solid lines), compared with these tools on their own (dashed lines) and concatenated with Shannon (dotted lines)

ClusTrast reached, while ClusTrast demonstrated higher recall for Trans-ABYSS-M and META.

The META assembly in itself yielded higher recall than ClusTrast when ran with any other primary assembly (Trans-ABYSS-M, Trinity, or RNA-Bloom). But ClusTrast with META as primary assembly yielded higher recall than both META on its own and META with Shannon concatenated.

In Fig. 10, we show the number of annotated and expressed isoforms that, of all the methods tested in our study, only ClusTrast managed to reconstruct, and with which

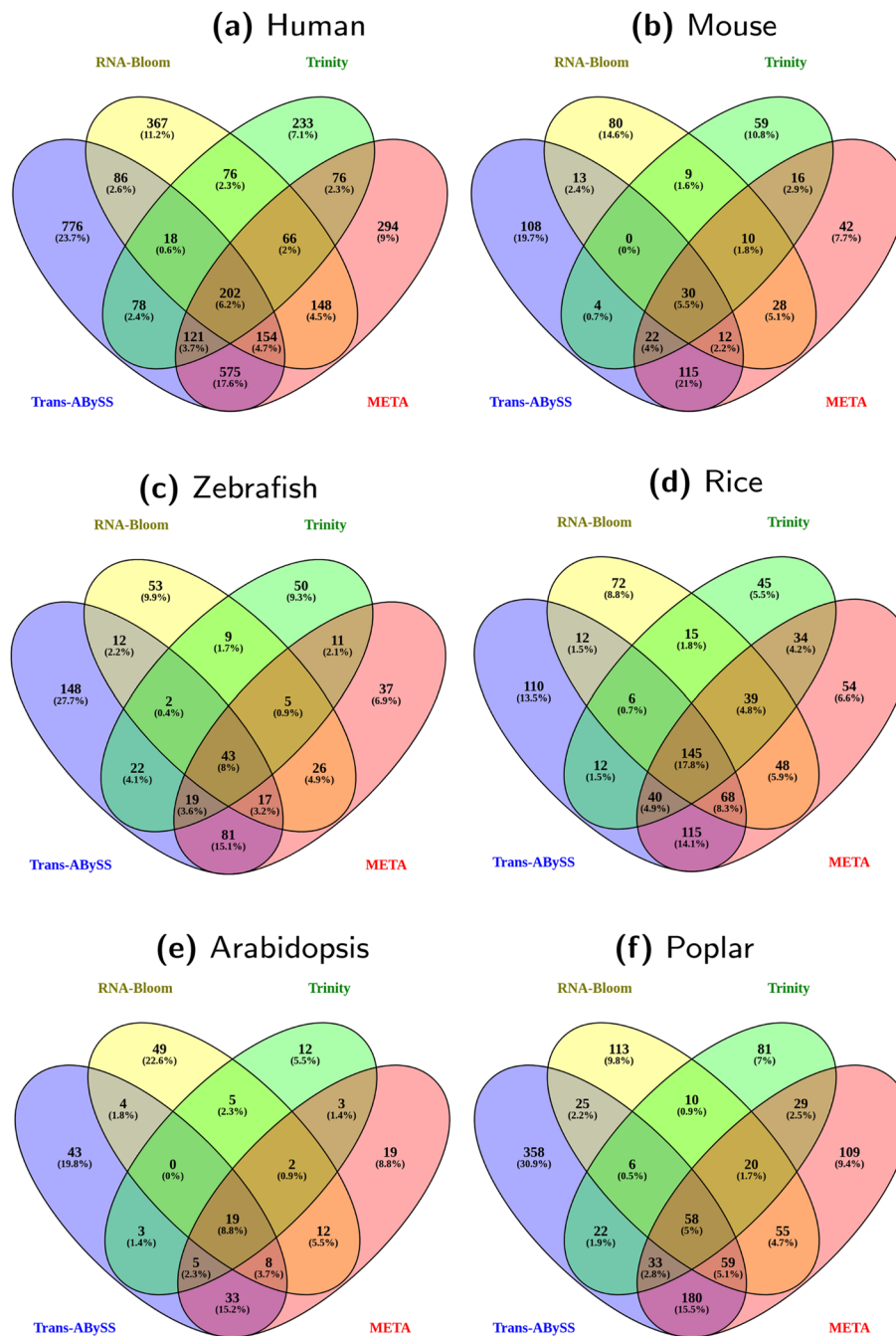


Fig. 10 ClusTrast results using four different tools - RNA-Bloom, Trinity, Trans-ABYSS-M, and META - for the primary assembly. The Venn diagrams show the number of annotated and expressed transcript isoforms reconstructed solely by ClusTrast and no other tested method. ClusTrast with primary assembly from Trans-ABYSS-M generates the highest number of unique transcript isoforms for all datasets (a) - (f)

tool for primary assembly (RNA-Bloom, Trinity, Trans-ABYSS-M, and META). ClusTrast with Trans-ABYSS-M for primary assembly (the default version of ClusTrast) managed to reconstruct most such isoforms in all instances except rice, where it placed second after ClusTrast with META as primary assembly.

The runtimes for ClusTrast with Trinity as primary assembly were comparable to using Trans-ABYSS-M for primary assembly. With RNA-Bloom (and thus also META) as primary assembly, however, ClusTrast required more computational resources: the clustering step took ~ 68 h to complete, and required memory of over 4TB.

Thus, all taken together, we believe that Trans-ABYSS-M is a reasonable choice for primary assembly in ClusTrast.

Discussion

We described and assessed the de novo transcriptome assembler ClusTrast. Across all six tested datasets, ClusTrast-M and ClusTrast-S created assemblies that had the highest or among the highest recall as measured by SQANTI and CRBB.

ClusTrast consists of several steps (Fig. 1). In the primary assembly step, we have chosen Trans-ABYSS as the primary assembler included in the distributed version of ClusTrast. We have shown that it is possible to use other assemblers as well, but with lower performance. Using the union of assembled transcripts from more than one assembly tool as the primary assembly improves the recall of ClusTrast. This might be suitable for some users, but requires larger computational resources. The clustering of the guiding contigs is a way to conceptually focus the assembly to genes, and their corresponding isoforms. The most common case is that a reference gene is represented in (i.e., is aligned to) only one cluster, although it is quite common that a reference gene, and its corresponding isoforms, is represented in more than one cluster (Additional file 1: Fig. S.19). In fact, ClusTrast was used in our recent study [30] to detect isoforms in individual gene families of specific biological interest. The core step of ClusTrast is the clusterwise assembly, which is followed by concatenation of the clusterwise assemblies, and finally by merging with the primary assembly (and removal of duplicates). A comparison of the recall and precision using (i) only the primary assembly, (ii) only the concatenated clusterwise assemblies, and (iii) the final, merged assembly is in Additional file 1: Tables S.28-S.31. The recall increases substantially using the merged assembly as opposed to only primary or clusterwise assemblies. The fact that the recall increases also when using a union of assembled transcripts from more than one assembly tool as primary assembly (discussed above) further supports that the clusterwise assembly is a helpful approach. We suggest that a wise choice of method for primary assembly coupled with the power of the clusterwise assembly are the key components of the high recall performance of ClusTrast.

In our evaluation of ClusTrast and the other transcriptome assemblers, we have emphasized metrics that do not penalize the reconstruction of different isoforms of a gene. A compact transcriptome assembly is, for instance, preferable when the reconstructed transcripts are intended to be used for aligning reads to a “reference” for, e.g., differential gene expression analysis. In such a situation, our approach would not be helpful. However, when the goal is to find as many supported transcript isoforms as possible, compactness is not in itself desirable and could in fact be counter productive. A recent review [16] points out that there is a need for metrics that better capture the performance with regards to transcript isoforms. Our use of SQANTI’s approach to evaluation is an attempt to address this. SQANTI [28] was originally designed to evaluate long reads (e.g., PacBio CCS) but the categories it defined and its aim to classify each

non-redundant transcript individually are useful also in evaluation of assembled transcripts. We have been conservative in that we included only the FSM and ISM categories as true positives. Also the NIC category, which encompasses reconstructed transcripts that contain annotated reference exons only, but in novel combinations, could have been included. We compared the results from SQANTI and CRBB (which has been used for transcriptome assembly benchmarking before), and detected a correlation between SQANTI and CRBB scores, particularly strong for recall. We noted that most of the contigs where SQANTI and CRBB classifications of true positives disagreed, belong to the SQANTI class ISM mono-exon. Excluding these from the set of true positives increased the correlation between SQANTI and CRBB precision measurements for ClusTrast. We believe this supports the notion that SQANTI is possible to use for transcriptome assembly evaluation.

Comparing ClusTrast-**M** to one of the most popular transcriptome assemblers, Trinity, revealed that ClusTrast-**M** detected more transcript isoforms than Trinity, and also had a higher precision for isoforms as measured by SQANTI, but clearly underperformed according to CRBB precision. The difference in relative performance of ClusTrast and Trinity according to CRBB and SQANTI precision may be explained by how CRBB handles assembled transcripts with high similarity: If two or more highly similar transcripts exist in the assembly, and some of them have a lower E-value than others, then only the transcripts with E-values below the limit will be considered CRBB hits and thus true positives for precision. SQANTI, in contrast, annotates each transcript independently and therefore calls all assembled transcripts that are similar to a reference transcript as a true positive. We assessed this by recalculating SQANTI precision while only counting one transcript match for every reference transcript (Additional file 1: Fig. S.6), and we observed a marked reduction in precision of Trans-ABySS-**M** and ClusTrast across all assemblies (compare Additional file 1: Figs. S.2 and S.6). We also tested ClusTrast with secondary alignments switched off, and observed a slight improvement for CRBB precision, but at the cost of a reduction in both CRBB and SQANTI recall for most datasets (Additional file 1: Table S.27).

We observed that ClusTrast generally recovered as many or more known isoforms as TrAB+Sh as measured by SQANTI (Additional file 1: Fig. S.2) while suffering only a small reduction in CRBB precision (Additional file 1: Fig. S.3) and that ClusTrast finished successfully on the rice dataset, where Shannon (and thus TrAB+Sh) failed to create an assembly. A possible explanation for both observations is that the clustering performed in ClusTrast may simplify sub-graphs enough to allow better handling by the Shannon heuristic and thereby increase sensitivity. ClusTrast-**M** and TrAB-**M**+Sh were also the best in reconstructing isoforms to their full length according to SQANTI.

In general the performance of the assemblers was rather consistent over species, regardless of evaluation approach (SQANTI or CRBB). We tested two additional human datasets (Additional file 1: Table S.1), for a total of four datasets, one of which simulated, and observed that ClusTrast showed the highest SQANTI and CRBB recall over the range of transcript coverage as well as expression levels for all four human datasets (Fig. 6; Additional file 1: Figs. S.7–S.9). However, the precision performance of ClusTrast was mixed. ClusTrast performance was not correlated to the size of the datasets

(Additional file 1: Fig. S.17). For all results on all additional datasets, see Additional file 1: Section C.3.

We used RSEM to estimate the number of expressed isoforms, and in our recall calculations we used only transcripts with $TPM > 0$. Using all reference transcripts in our evaluation, instead of only those that have a $TPM > 0$, would mean a larger denominator when calculating recall, which would lower the recall of all compared assemblers alike. If RSEM makes a mistake and assigns $TPM = 0$ to a reference transcript that in fact is transcribed, then the recall will be underestimated. If RSEM assigns a $TPM > 0$ to a reference transcript which in fact is not transcribed, recall will be overestimated. Similarly, there might be assembled contigs that correspond to real isoforms that are not present in the reference. These contigs are counted as false negatives while they should be true positives, thus precision performance is likely underestimated. For the SQANTI evaluation, it is possible that many of these would be considered as true positives if we had included the NIC category among the true positives.

Within the scope of this paper, we have shown that the current de novo transcriptome assembly strategy of ClusTrast is successful in finding the most comprehensive set of isoforms from short read RNA-seq datasets, as it outperformed all other tested methods. We do not claim, however, that we have found the optimal version. Future improvements might come, e.g., from testing other assemblers for the clusterwise assembly, or, provided suitable datasets are available, from including long reads as guiding contigs. The use of long reads in transcriptome analysis is not a focus of this study, but the accuracy of long reads has indeed reached a level that enables the use of long read sequencing to reliably address questions not easily amenable with short read sequencing. For instance, long reads can provide direct information about transcript isoforms present in a sample or the methylation status of DNA or RNA molecules [31]. The most telling example is perhaps the telomere-to-telomere (T2T) sequencing of the human genome [32]: long reads revealed approximately 200 Mbases of genomic DNA sequence hitherto undetermined, containing 99 protein coding genes not present in the human reference genome GRCh38. There are emerging methods for analysis of long read transcriptome data, e.g. StringTie2 which can accommodate both short and long reads for a reference-based transcriptome assembly [33] and IsoQuant that enables transcript discovery in long read data sets [34]. Finally, we note that a short read de novo transcriptome assembly approach, such as ClusTrast or any of the ones included in our study, would potentially be able to find the likes of the missing 99 human genes in other species which not yet have had their full T2T genomes sequenced. This attests to the lingering usefulness of short read sequencing and also to the advantages of de novo transcriptome assembly, as a reference-based method by definition would miss any genes or transcripts not present in the reference.

Conclusion

In our tests of model organisms, ClusTrast consistently detected the most transcript isoforms, but at a cost of lower precision. This agrees with our intention of ClusTrast – to provide a comprehensive but non-redundant list of contigs. Therefore, we believe researchers interested in a more complete representation of transcript isoforms from

eukaryotic organisms may wish to use ClusTrast. The resulting list of contigs is amenable for further processing and analysis tailored according to the research question at hand.

Availability and requirements

Project name:	ClusTrast
Project home page:	https://github.com/karljohanw/clustrast
Operating systems:	Linux and MacOS
Programming language:	Bashscript
Requirements:	transabyss, shannon_cpp, isONclust, mini-map2, awk.
License:	GPLv3
Restrictions to use by non-academics:	None

Abbreviations

CRBB	Conditional reciprocal best BLAST
FDR	False discovery rate
FSM	Full splice match
GC	Guiding contig
ISM	Incomplete splice match
NIC	Novel in catalog
NNC	Novel not in catalog
PPV	Positive predictive value
SNP	Single nucleotide polymorphism
SR	Short read
TPM	Transcript per million
TPR	True positive rate

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12859-024-05663-3>.

Additional file 1. Lists tools and versions used, results from additional datasets, additional analyses and illustrations on how the pipeline works.

Acknowledgements

We wish to thank Pelin Akan Sahlén at KTH Royal Institute of Technology for sharing access to the server SAGA.

Author contributions

KJW implemented the method, ran all the experiments and wrote the original version of the manuscript. KJW and OE wrote the final manuscript, with contributions from WWK. KJW and OE analyzed the results. OE and WWK supervised the project. All authors read and approved the final version of the manuscript.

Funding

Open access funding provided by Royal Institute of Technology. This work was supported by FORMAS [2013-650] and the Swedish Research Council [2018-05973]. Computations were enabled by resources, ParallellDatorCentrum (PDC) at KTH Royal Institute of Technology, provided by the Swedish National Infrastructure for Computing (SNIC), partially funded by the Swedish Research Council [2018-05973].

Availability of data and materials

The datasets used in this study are available on NCBI SRA (<https://www.ncbi.nlm.nih.gov/sra>) according to their SRA IDs listed in Table 1 and Additional file 1: Table S.1. The simulated datasets are available as supplementary material of their respective publication [13, 15]. The source code implemented in this study can be found at <https://github.com/karljohanw/clustrast>.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare no competing interests.

Received: 21 November 2022 Accepted: 18 January 2024

Published online: 01 February 2024

References

- Wang ET, Sandberg R, Luo S, Khrebttukova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP, Burge CB. Alternative isoform regulation in human tissue transcriptomes. *Nature*. 2008;456(7221):470–6. <https://doi.org/10.1038/nature07509>.
- Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Gueroussov S, Lee LJ, Slobodeniuc V, Kutter C, Watt S, Colak R, Kim T, Misquitta-Ali CM, Wilson MD, Kim PM, Odom DT, Frey BJ, Blencowe BJ. The evolutionary landscape of alternative splicing in vertebrate species. *Science*. 2012;338(6114):1587–93. <https://doi.org/10.1126/science.1230612>.
- Floor SN, Doudna JA. Tunable protein synthesis by transcript isoforms in human cells. *eLife*. 2016;5:10921. <https://doi.org/10.7554/eLife.10921>.
- Fackenthal JD, Godley LA. Aberrant RNA splicing and its functional consequences in cancer cells. *Dis Models Mech*. 2008;1(1):37–42. <https://doi.org/10.1242/dmm.000331>.
- Sterne-Weiler T, Sanford JR. Exon identity crisis: disease-causing mutations that disrupt the splicing code. *Genome Biol*. 2014;15(1):201. <https://doi.org/10.1186/gb4150>.
- Xiong HY, Alipanahi B, Lee LJ, Bretschneider H, Merico D, Yuen RKC, Hua Y, Gueroussov S, Najafabadi HS, Hughes TR, Morris Q, Barash Y, Krainer AR, Jovic N, Scherer SW, Blencowe BJ, Frey BJ. The human splicing code reveals new insights into the genetic determinants of disease. *Science*. 2015;347(6218):1254806. <https://doi.org/10.1126/science.1254806>.
- Akhter S, Kretzschmar WW, Nordal V, Delhomme N, Street N, Nilsson O, Emanuelsson O, Sundström JF. Integrative analysis of three RNA sequencing methods identifies mutually exclusive exons of MADS-box isoforms during early bud development in *Picea abies*. *Front Plant Sci*. 2018;9:1625. <https://doi.org/10.3389/fpls.2018.01625>.
- Garber M, Grabherr MG, Guttman M, Trapnell C. Computational methods for transcriptome annotation and quantification using RNA-seq. *Nat Methods*. 2011;8(6):469–77. <https://doi.org/10.1038/nmeth.1613>.
- Robertson G, Schein J, Chiu R, Corbett R, Field M, Jackman SD, Mungall K, Lee S, Okada HM, Qian JQ, Griffith M, Raymond A, Thiessen N, Cezard T, Butterfield YS, Newsome R, Chan SK, She R, Varhol R, Kamoh B, Prabhu A-L, Tam A, Zhao Y, Moore RA, Hirst M, Marra MA, Jones SJM, Hoodless PA, Birol I. De novo assembly and analysis of RNA-seq data. *Nat Methods*. 2010;7:909–12.
- Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, di Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N, Regev A. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol*. 2011;29(7):644–52. <https://doi.org/10.1038/nbt.1883>. [arXiv:1512.00567](https://arxiv.org/abs/1512.00567).
- Schulz MH, Zerbino DR, Vingron M, Birney E. Oases: robust de novo RNA-seq assembly across the dynamic range of expression levels. *Bioinformatics*. 2012;28(8):1086–92. <https://doi.org/10.1093/bioinformatics/bts094>.
- Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Li S, Zhou X, Lam T-W, Li Y, Xu X, Wong GK-S, Wang J. SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads. *Bioinformatics*. 2014;30(12):1660–6. <https://doi.org/10.1093/bioinformatics/btu077>.
- Hölzer M, Marz M. De novo transcriptome assembly: a comprehensive cross-species comparison of short-read RNA-Seq assemblers. *GigaScience*. 2019. <https://doi.org/10.1093/gigascience/giz039>.
- Bushmanova E, Antipov D, Lapidus A, Pribelski AD. rnaSPAdes: a de novo transcriptome assembler and its application to RNA-Seq data. *GigaScience*. 2019. <https://doi.org/10.1093/gigascience/giz100>.
- Hayer KE, Pizarro A, Lahens NF, Hogenesch JB, Grant GR. Benchmark analysis of algorithms for determining and quantifying full-length mRNA splice forms from RNA-seq data. *Bioinformatics*. 2015. <https://doi.org/10.1093/bioinformatics/btv488>.
- Thind AS, Monga I, Thakur PK, Kumari P, Dindhoria K, Krzak M, Ranson M, Ashford B. Demystifying emerging bulk RNA-Seq applications: the application and utility of bioinformatic methodology. *Brief Bioinform*. 2021. <https://doi.org/10.1093/bib/bbab259>.
- Smith-Unna R, Boursnell C, Patro R, Hibberd JM, Kelly S. TransRate: reference-free quality assessment of de novo transcriptome assemblies. *Genome Res*. 2016;26(8):1134–44. <https://doi.org/10.1101/gr.196469.115>.
- Sahlin K, Medvedev P. De novo clustering of long-read transcriptome data using a greedy, quality value-based algorithm. *J Comput Biol*. 2020;27(4):472–84. <https://doi.org/10.1089/cmb.2019.0299>.
- Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;34(18):3094–100. <https://doi.org/10.1093/bioinformatics/bty191>.
- Kannan S, Hui J, Mazooji K, Pachter L, Tse D. Shannon: An Information-Optimal de Novo RNA-Seq Assembler. Preprint at bioRxiv (2016). <https://www.biorxiv.org/content/early/2016/02/09/039230.full.pdf>. <https://doi.org/10.1101/039230>. <https://www.biorxiv.org/content/early/2016/02/09/039230>.
- Mao S, Pachter L, Tse D, Kannan S. RefShannon: a genome-guided transcriptome assembler using sparse flow decomposition. *PLoS ONE*. 2020;15(6):1–14. <https://doi.org/10.1371/journal.pone.0232946>.
- Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*. 2018;34(17):884–90. <https://doi.org/10.1093/bioinformatics/bty560>.
- Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinform*. 2011;12(1):323. <https://doi.org/10.1186/1471-2105-12-323>.

24. Liu J, Li G, Chang Z, Yu T, Liu B, McMullen R, Chen P, Huang X. BinPacker: packing-based de novo transcriptome assembly from RNA-seq data. *PLoS Comput Biol*. 2016;12(2):1004772–1004772. <https://doi.org/10.1371/journal.pcbi.1004772>.
25. Liu J, Yu T, Mu Z, Li G. TransLiG: a de novo transcriptome assembler that uses line graph iteration. *Genome Biol*. 2019;20(1):81. <https://doi.org/10.1186/s13059-019-1690-7>.
26. Nip KM, Chiu R, Yang C, Chu J, Mohamadi H, Warren RL, Birol I. RNA-Bloom enables reference-free and reference-guided sequence assembly for single-cell transcriptomes. *Genome Res*. 2020;30(8):1191–200. <https://doi.org/10.1101/gr.260174.119>.
27. Aubry S, Kelly S, Kumpers BMC, Smith-Unna RD, Hibberd JM. Deep evolutionary comparison of gene expression identifies parallel recruitment of trans-factors in two independent origins of c4 photosynthesis. *PLOS Genet*. 2014;10(6):1–16. <https://doi.org/10.1371/journal.pgen.1004365>.
28. Tardaguila M, de la Fuente L, Marti C, Pereira C, Pardo-Palacios FJ, del Risco H, Ferrell M, Mellado M, Macchietto M, Verheggen K, Edelmann M, Ezkurdia I, Vazquez J, Tress M, Mortazavi A, Martens L, Rodriguez-Navarro S, Moreno V, Conesa A. SQANTI: extensive characterization of long read transcript sequences for quality control in full-length transcriptome identification and quantification. *Genome Res*. 2018;28(1):1–16. <https://doi.org/10.1101/gr.222976.117>.
29. 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).
30. Akhter S, Westrin KJ, Zivi N, Nordal V, Kretzschmar WW, Delhomme N, Street NR, Nilsson O, Emanuelsson O, Sundström JF. Cone-setting in spruce is regulated by conserved elements of the age-dependent flowering pathway. *New Phytol*. 2022;236(5):1951–63. <https://doi.org/10.1111/nph.18449>.
31. Kovaka S, Ou S, Jenike KM, Schatz MC. Approaching complete genomes, transcriptomes and epi-omes with accurate long-read sequencing. *Nat Methods*. 2023;20(1):12–6. <https://doi.org/10.1038/s41592-022-01716-8>.
32. ...Nurk S, Koren S, Rhie A, Rautiainen M, Bizkadze AV, Mikheenko A, Vollger MR, Altemose N, Uralsky L, Gershman A, Aganezov S, Hoyt SJ, Diekhans M, Logsdon GA, Alonge M, Antonarakis SE, Borchers M, Bouffard GG, Brooks SY, Caldas GV, Chen N-C, Cheng H, Chin C-S, Chow W, de Lima LG, Dishuck PC, Durbin R, Dvorkina T, Fiddes IT, Formenti G, Fulton RS, Fungtammasan A, Garrison E, Grady PGS, Graves-Lindsay TA, Hall IM, Hansen NF, Hartley GA, Haukness M, Howe K, Hunkapiller MW, Jain C, Jain M, Jarvis ED, Kerpedjiev P, Kirsche M, Kolmogorov M, Korlach J, Kremitzki M, Li H, Maduro VV, Marschall T, McCartney AM, McDaniel J, Miller DE, Mullikin JC, Myers EW, Olson ND, Paten B, Peluso P, Pevzner PA, Porubsky D, Potapova T, Rogaev EI, Rosenfeld JA, Salzberg SL, Schneider VA, Sedlazeck FJ, Shafin K, Shew CJ, Shumate A, Sims Y, Smit AFA, Soto DC, Sovići I, Storer JM, Streets A, Sullivan BA, Thibaud-Nissen F, Torrance J, Wagner J, Walenz BP, Wenger A, Wood JMD, Xiao C, Yan SM, Young AC, Zarate S, Surti U, McCoy RC, Dennis MY, Alexandrov IA, Gerton JL, O'Neill RJ, Timp W, Zook JM, Schatz MC, Eichler EE, Miga KH, Phillippy AM. The complete sequence of a human genome. *Science*. 2022;376(6588):44–53. <https://doi.org/10.1126/science.abj6987>.
33. Kovaka S, Zimin AV, Pertea GM, Razaghi R, Salzberg SL, Pertea M. Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome Biol*. 2019;20(1):278. <https://doi.org/10.1186/s13059-019-1910-1>.
34. Prjibelski AD, Mikheenko A, Joglekar A, Smetanin A, Jarroux J, Lapidus AL, Tilgner HU. Accurate isoform discovery with isoquant using long reads. *Nat Biotechnol*. 2023. <https://doi.org/10.1038/s41587-022-01565-y>.

Publisher's Note

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