

METHOD

Open Access

# Bookend: precise transcript reconstruction with end-guided assembly



Michael A. Schon<sup>1,2\*</sup>, Stefan Lutzmayer<sup>2</sup>, Falko Hofmann<sup>2</sup> and Michael D. Nodine<sup>1,2\*</sup> 

\* Correspondence: [michael.schon@wur.nl](mailto:michael.schon@wur.nl); [michael.nodine@wur.nl](mailto:michael.nodine@wur.nl)

<sup>1</sup>Cluster of Plant Developmental Biology, Laboratory of Molecular Biology, Wageningen University & Research, Wageningen 6708, PB, The Netherlands

<sup>2</sup>Gregor Mendel Institute (GMI), Austrian Academy of Sciences, Vienna Biocenter (VBC), Dr. Bohr-Gasse 3, 1030 Vienna, Austria

## Abstract

We developed Bookend, a package for transcript assembly that incorporates data from different RNA-seq techniques, with a focus on identifying and utilizing RNA 5' and 3' ends. We demonstrate that correct identification of transcript start and end sites is essential for precise full-length transcript assembly. Utilization of end-labeled reads present in full-length single-cell RNA-seq datasets dramatically improves the precision of transcript assembly in single cells. Finally, we show that hybrid assembly across short-read, long-read, and end-capture RNA-seq datasets from *Arabidopsis thaliana*, as well as meta-assembly of RNA-seq from single mouse embryonic stem cells, can produce reference-quality end-to-end transcript annotations.

**Keywords:** RNA-seq, Transcriptome, Single-cell, TSS, PAS, Capping, Polyadenylation, 5' and 3' ends, Long-read, Iso-Seq

## Summary statement

Bookend is a generalized framework that utilizes RNA 5' and 3' end information in RNA-seq datasets to accurately reconstruct transcriptomes including those from single cells.

## Background

The functions of genes depend on the amount and types of RNA molecules that they produce. Variation in transcript initiation, splicing, and polyadenylation can generate an array of RNA isoforms, and cataloging how these RNA variants change across development and disease provides insights into corresponding gene functions [1–3]. Large-scale projects dedicated to the manual curation of gene annotations are extremely valuable, but are labor-intensive and thus limited in scope to the most well-studied organisms [4–7]. Moreover, multicellular organisms have difficult-to-access cell types that will inevitably be overlooked by even the most comprehensive annotation projects [8]. The completeness and accuracy of a reference annotation can considerably impact all downstream data analyses, from gene expression to predictions of gene function [9–11]. To understand how transcriptome architecture varies during development and in response to disease, it is therefore valuable to have an automated



© The Author(s). 2022 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

method that accurately identifies transcript isoforms. Accordingly, many computational tools have been developed for genome annotation including software that utilizes the massive and growing diversity of RNA sequencing (RNA-seq) technologies [12].

A wide array of RNA-seq protocols have been developed to profile different aspects of the transcriptome, from strand-specific coverage of gene bodies [13] to selective amplification of RNA 5' ends [14–17], 3' ends [18, 19] or simultaneous capture of both ends [20, 21]. Major recent advances have enabled the amplification of full-length transcripts from single cells [22, 23] or 3' end capture from millions of cells [24–26]. In parallel, advances have been made for profiling RNA on “third-generation” long-read sequencing platforms such as PacBio and Oxford Nanopore single-molecule sequencers that can read continuous DNA and/or RNA molecules and yield end-to-end complete transcript sequences [27, 28].

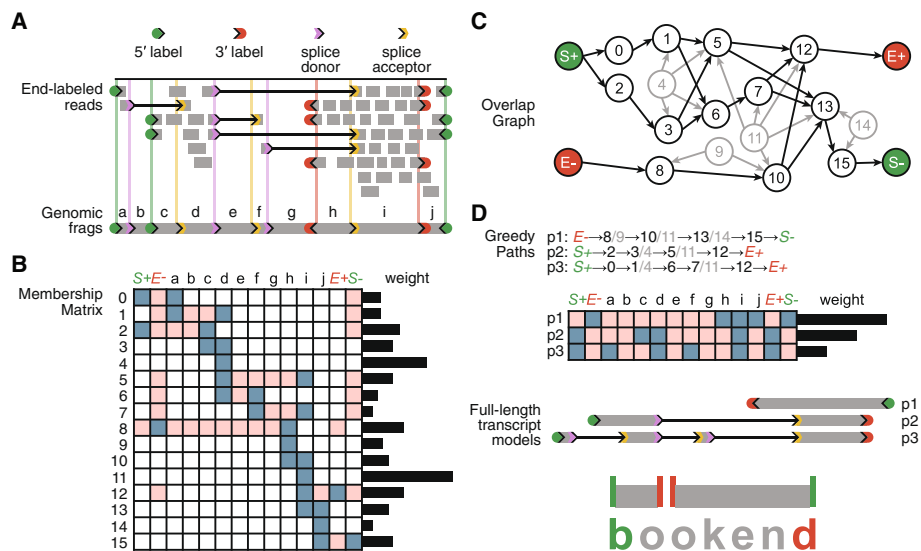
Transcript assembly is the effort to distill information from RNA-seq experiments into a comprehensive annotation of the transcript isoforms present in the corresponding samples. Depending on the method, RNA-seq reads contain a broad spectrum of information content. At one extreme, single-end reads from non-stranded RNA-seq protocols can be 50 nucleotides (nt) or shorter and sequenced from one end of a double-stranded cDNA fragment such that the resulting sequence is a random substring of an RNA molecule or its reverse complement. Paired-end reads contain two ends of a cDNA molecule, and typically, there is a gap of unknown length between the mate pairs. When aligned to a reference genome, paired reads may span more than one splice junction, indicating that these splicing events occurred in the same molecule. Some strand-specific RNA-seq protocols selectively sequence only first-strand or second-strand cDNA to preserve knowledge of the original mRNA molecule's orientation [13]. Other protocols selectively capture and sequence a fragment immediately downstream of the RNA 5' end or upstream of the 3' end, demarcating precisely where that molecule begins or ends, respectively [14, 16–18, 29, 30]. Finally, the most information-rich reads come from long-read sequencing, in which the RNA or cDNA is read in its entirety without fragmentation. Long-read methods are a promising tool for transcript annotation, but current protocols are more error-prone per base sequenced, less sensitive, and more costly than comparable short-read experiments. Because the vast majority of existing RNA-seq data is in short-read format, nearly all assemblers have aimed to reconstruct transcripts from paired-end short reads. A long-recognized problem of assemblers is the inaccurate annotation of transcript start sites (TSS) and polyadenylation sites (PAS) [31, 32]. Existing short-read assemblers infer TSSs and PASs through heuristics such as changes in read coverage, but such changes can also be due to alignment errors, biased RNA fragmentation, sample degradation, or spurious intron retention. Long-read sequencing methods are designed to read RNA from TSS to PAS, but they remain susceptible to a variety of experimental artifacts [32]. The increasing adoption of long reads for transcript annotation has led to a separate suite of tools that summarize, collapse, or “polish” long reads to remove erroneous structures and present a set of representative isoforms from these reads [33, 34]. For example, the recently developed transcript assembler StringTie2 reports the use of long reads in assembly by removing aligned segments with a high error rate and assembling the resulting gapped reads [35]. Transcript annotation would ideally integrate information from a variety of RNA-seq methods to determine the best evidence for transcript

starts, ends, and splicing patterns in a tissue-of-interest. However, current transcriptome assembly methods do not employ information about where RNA molecules begin and end. Here, we describe a method utilizing RNA 5' and 3' end information produced by a variety of RNA-seq protocols to accurately reconstruct transcriptomes including those from single cells.

## Results

### A framework for end-guided transcript assembly

To determine whether RNA 5' and 3' end information can improve transcript assembly algorithms, we developed a generalized framework for identifying RNA ends in sequencing data and using this information to assemble transcript isoforms as paths through a network accounting for splice sites, transcription start sites (TSS), and polyadenylation sites (PAS). Because this software uses end information to guide transcript assembly, we named it Bookend. Importantly, Bookend takes RNA-seq reads from any method as input and after alignment to a reference genome, reads are stored in a lightweight end-labeled read (ELR) file format that records all RNA boundary features (5' labels, splice donors, splice acceptors, gaps, 3' labels), as well as the sample of origin for that read (see Additional file 1: Supporting notes). Assembly is then resolved at each locus with aligned reads through a four-step procedure (Fig. 1; see Methods and Additional file 1: Supporting notes). First, boundary labels from all aligned RNA-seq reads are clustered and filtered to demarcate a unique set of locus TSSs, PASs, and splice junctions. Each locus is partitioned into a set of nonoverlapping "frags" defined



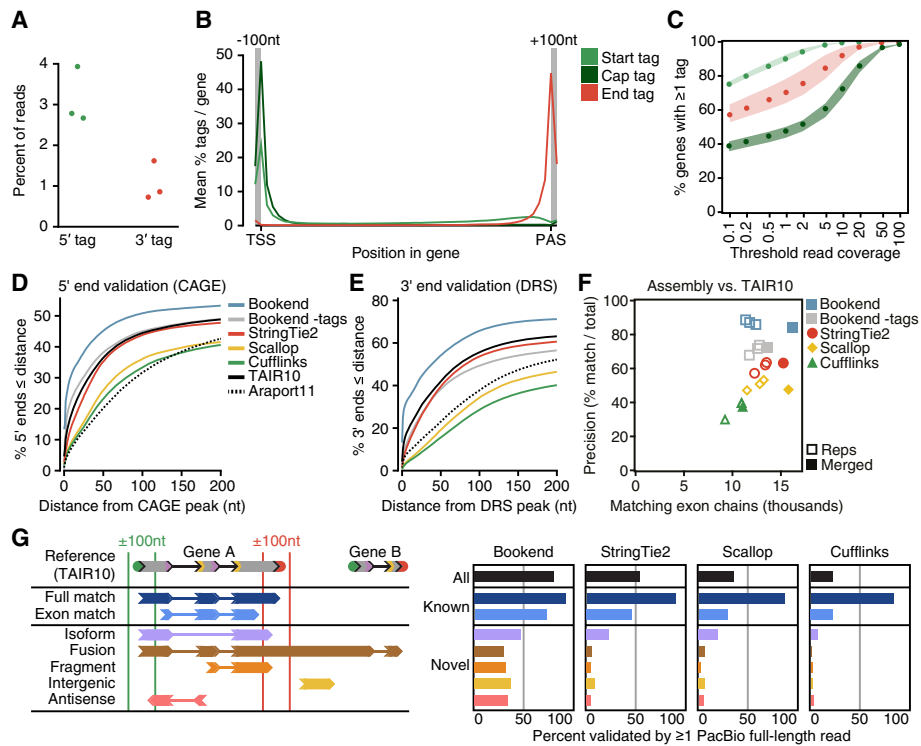
**Fig. 1** End-guided assembly with Bookend. **A** Individual RNA-seq reads are mapped to a genome, recording which reads mark a transcript 5' or 3' end, and which reads span one or more splice junctions. Ranges between adjacent features are recorded as frags. **B** Each unique read structure is recorded in a condensed representation as one element in a Membership Matrix; blue—included, pink—excluded. The weight of each element is the coverage depth of matching reads (sequenced bases/length) across the element. **C** A directed graph is constructed between overlapping elements of the Membership matrix. Weights of contained elements (gray) are distributed proportionally to their containers. **D** A set of optimal paths through the graph is iteratively constructed from the heaviest unassigned elements. Complete Paths are output as full-length transcript annotations

as the spans between adjacent boundary labels. Four additional frags (S+, E+, S-, E-) denote the presence of a Start or End Tag on the forward or reverse strand. Second, a Membership Matrix is generated to redefine all aligned reads with respect to the locus frags. A read's Membership includes each frag it overlaps and excludes each incompatible frag (e.g., a spanned intron, a region upstream of a TSS or downstream of a PAS). Reads with identical patterns of Membership are condensed to a single element (row) of the Membership Matrix, whose weight is the total coverage depth across the element by all reads of that pattern. Third, an Overlap Graph is constructed from the Membership Matrix elements and this directed graph is simplified by collapsing shorter elements into the elements that contain them. Finally, the Overlap Graph is iteratively traversed to resolve an optimal set of Greedy Paths from TSSs to PASs. These Paths describe a set of full-length transcript models best supported by the input reads. The Membership Matrix definition is flexible enough to utilize reads regardless of their length, alignment gaps, strand, or end information (Additional file 1: Fig. S1B).

### End-labeled reads improve the quality of transcript assembly

*Arabidopsis thaliana* (*Arabidopsis*) is an ideal model to benchmark transcript assembly in higher eukaryotes. The *Arabidopsis* genome is compact (~119 megabases), contains few repetitive elements, and the TAIR10 reference annotation was extensively curated from expressed sequence tag (EST) data [7]. To determine whether end-labeled reads improve assembly, we examined libraries generated with the low-input sequencing method Smart-seq2 from *Arabidopsis* floral buds [16]. Two crucial steps in the Smart-seq2 protocol, template switching and preamplification, enrich for full-length cDNA with an oligo label at both the 5' (template switching oligo, TSO) and 3' (oligo-dT) ends [22]. These oligos were trimmed from all reads and a record was kept of which end label was found (5', 3', or no label) before mapping to the genome. As anticipated, a small percentage of reads were found with either label (Fig. 2A; Additional file 1: Table S1). All reads were aligned to the *Arabidopsis* genome, and the terminal positions of 5'- and 3'-labeled reads were retained as "Start Tags" and "End Tags," respectively. Of End Tags mapping to annotated genes, 88% mapped near PASs, defined as the last decile of the gene or up to 100 nt downstream (Fig. 2B). Start Tags had lower specificity for TSSs, with only 48% of Start Tags in the first decile of genes or up to 100 nt upstream. Template switching is known to readily occur at RNA 5' ends derived from in vivo or in vitro RNA decay. However, a subset of reads contain an intervening G between the TSO and the genome-aligned sequence, indicating a 7-methylguanosine cap on the template RNA [16, 29, 37]. The upstream untemplated G (uuG)-containing Start Tags were classified as Cap Tags. Cap Tags were rare relative to all Start Tags (9%), but were much more specific to TSSs with an average of 88% of Cap Tags within each gene mapping near the 5' end (Fig. 2B). To optimize detection of true transcript 5' and 3' ends, the Tag Clustering algorithm designed for Bookend defines Tag weight as a function of total read depth and applies a bonus to Cap Tags over non-uuG Start Tags (See Additional file 1: Supporting notes: "Tag Clustering").

Despite end-labeled reads being relatively rare, the preamplification process should ensure that a TSO or oligo-dT sequence is at each end of every cDNA molecule prior to tagmentation. Therefore, we expected end-labeled reads to be distributed widely



**Fig. 2** End-labeled Smart-seq2 reads accurately detect transcript 5' and 3' ends. **A** Percent of reads in three Smart-seq2 libraries that contained a 5'-labeled or 3'-labeled junction, respectively. **B** Average signal strength per gene of Start, End, and Cap Tags along gene bodies in 50 bins with an additional 100 nt flanking each gene boundary. Start Tag, any 5' label; Cap Tag, 5' label with upstream untemplated G (uuG); End Tag, 3' label. **C** Likelihood of a gene to possess  $\geq 1$  Start, Cap, or End Tag as a function of aligned read coverage (average read depth/base). **D** Cumulative frequency of annotated 5' ends as a function of distance from the closest CAGE peak [36]. **E** Distance of 3' ends from the nearest DRS peak [30] as in (D). **F** Performance of three transcript assemblers, measured by total number of reference-matching exon chains (x-axis) vs. percent of assembled transcripts that match the reference (y-axis). **G** (Left) Schematic depicting classifications of assembled transcripts against the closest TAIR10 reference isoform. (Right) Rate of validation by PacBio full-length non-chimeric (FLNC) reads for different assemblies, grouped by classification

across the genome wherever reads exist. As predicted, the majority of genes with  $> 0$  read coverage contained  $\geq 1$  Start Tag and  $\geq 1$  End Tag, and the likelihood of finding a Start or End Tag increased as a function of total read coverage (Fig. 2C). Of all genes with at least 1 $\times$ , 10 $\times$ , and 100 $\times$  read coverage, 73.3, 94.4, and 99.2% possessed both a Start and End Tag, respectively.

To assess whether end-labeled reads mark real TSSs and PASs at nucleotide precision, Bookend was used to assemble all floral bud Smart-seq2 reads either with or without utilizing Start and End Tags. Additionally, three leading short-read transcript assemblers were used with comparable settings (see “Methods”): StringTie2 [35, 38], Scallop [39], and Cufflinks [40]. Publicly available Arabidopsis CAGE [36] and Direct RNA-seq (DRS [30]) datasets were used to validate 5' and 3' ends, respectively. All three of these widely used assemblers output thousands of single-exon unstranded fragments, which were ambiguous with regard to which end is 5' or 3' and thus were discarded from further analyses (Additional file 1: Table S2). Bookend-defined TSSs based on Start/Cap Tags were more likely to have a CAGE peak within 200 nt than 5' ends reported either by Bookend without the use of Start Tags, the three leading assemblers,

or even the current *Arabidopsis* reference annotations (Fig. 2D). Likewise, a higher proportion of Bookend-identified PASs were supported by DRS reads than PASs reported by the other transcript assemblers and *Arabidopsis* reference annotations (Fig. 2E). At the nucleotide level, Bookend-defined transcript boundaries were more than twice as likely to agree with the exact experimentally determined TSS and PAS peak positions than the most accurate reference annotation (TAIR10), while the other three assemblers reported transcript boundaries less accurate than TAIR10 (Additional file 1: Fig. S2A,B). Strikingly, even the Bookend 5' and 3' ends > 100 nt from any reference still possessed known sequence motifs associated with TSS and PAS, respectively, whereas sequence content around novel ends from Cufflinks, Scallop, and StringTie2 is largely incoherent (Additional file 1: Fig. S2C,D). In addition to a dramatic increase in transcript boundary accuracy, 16,158 exon chains predicted by Bookend fully matched a TAIR10 reference transcript, which was higher than when end-labeled reads were ignored (13,660) and exceeded the totals from Scallop (15,785), StringTie2 (15,253), or Cufflinks (11,051) (Fig. 2F). Therefore, Bookend correctly builds more known transcripts than other assemblers and Bookend-annotated 5' and 3' ends were more precise than even the most accurate *Arabidopsis* reference annotation.

In addition to known transcripts, Bookend constructed 2979 isoforms not present in TAIR10, which was 66% fewer than StringTie2 (8,886), 83% fewer than Scallop (17,400), and 84% fewer than Cufflinks (18,934). An assembled transcript may fail to match TAIR10 either because the assembly is incorrect or because the reference is incomplete. To distinguish between these possibilities, two long-read SMRT cells of floral bud RNA were sequenced with the PacBio platform to yield 547,910 full-length non-chimeric (FLNC) reads. All short-read assemblies were partitioned into 7 different classifications based on their relationship to the most similar TAIR10 model (Fig. 2G). A transcript model was considered experimentally validated if at least one aligned PacBio read fully matched the model (entire exon chain,  $\pm 100$  nt ends). Of all Bookend transcripts, 81.2% were supported by PacBio data, which surpassed the validation of transcripts predicted by StringTie2 (54.7%), Scallop (35.9%), or Cufflinks (22.3%) (Fig. 2G; Additional file 1: Table S2). Reference-matching transcripts have a higher average estimated abundance than non-reference transcripts, making the latter more difficult to validate with the limited throughput of long-read sequencing (Additional file 1: Fig. S2E). Despite this limitation, 42.3% of non-reference Bookend assemblies were fully supported by at least one PacBio read, which was substantially higher than the validation rate of non-reference transcript assemblies generated by StringTie2 (15.9%), Scallop (11.6%), and Cufflinks (4.3%) (Fig. 2G). Taken together, these results demonstrate that end-guided assembly using latent RNA end information enables precise transcript reconstruction from end-labeled short-read datasets.

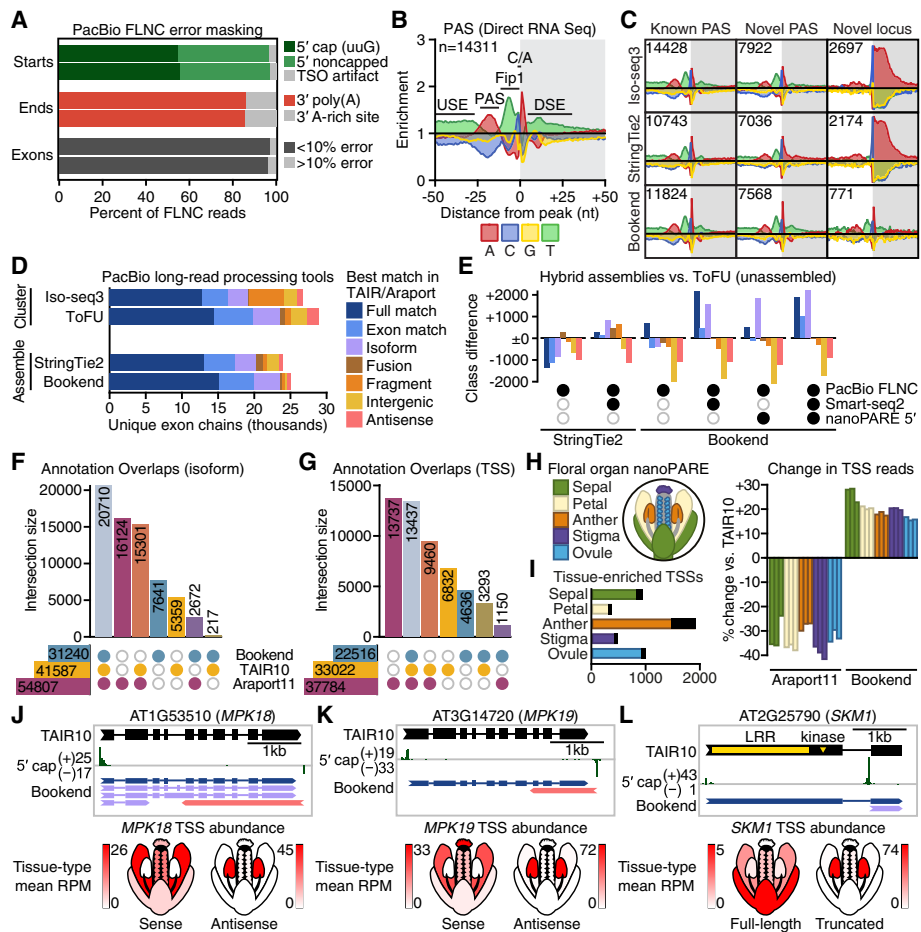
### Hybrid assembly refines and complements long-read RNA-seq

Long-read sequencing technologies do not obviate the need for transcript reconstruction. Various sources of technical and biological noise result in fragmented or improperly spliced long reads [32, 41]. Long-read approaches also suffer from a higher base-level error rate compared to short-read platforms [42]. Error correcting methods such as Circular Consensus Sequencing (CCS) require reverse transcription and cDNA amplification, which are susceptible to mispriming and template-switching artifacts [43, 44]. This has driven the

ongoing development of tools to refine transcript models derived from long reads [33, 34]. Additionally, StringTie2 was recently repurposed to assemble long reads [35].

To quantify potential sources of error, PacBio FLNC reads were aligned to the genome and processed by the Bookend pipeline to identify and remove template-switching artifacts, oligo-d(T) mispriming events at A-rich regions, and exons with a high alignment error (Fig. 3A). Across both SMRT cells, 95.4% of reads aligned successfully, and 97.0% of alignments did not contain any high-error exons, defined as the total length of mismatches, inserts, and deletions exceeding 10% of the exon length. However, 14.1% of all FLNC 3' end labels were removed due to alignment failure or the presence of an A-rich region immediately downstream of the oligo-d(T) junction. If treated as genuine 3' ends, these reads can cause false annotation of 3'-UTRs or putative transcripts antisense or intergenic to known genes [43] (Additional file 1: Fig. S3A). Direct RNA sequencing bypasses oligo-d(T) priming and was used to produce a map of genuine Arabidopsis PAS [30]. These sites show a distinct pattern of nucleotide enrichment, including a C/A dinucleotide motif at the cleavage and polyadenylation site itself, and a U-rich upstream element (USE) and downstream element (DSE) (Fig. 3B). Three tools were used to reduce the PacBio FLNC data into a unique set of transcripts: the Iso-seq3 clustering algorithm from PacBio, assembly by StringTie2, and end-guided assembly by Bookend. All 3 methods could recapitulate known PAS motifs at the set of 3' ends within 100 nt of a TAIR10-annotated PAS. StringTie2-annotated 3' ends showed a slight A-richness at novel 3' ends, and both Iso-seq3 and StringTie2 annotations contain thousands of putative novel antisense or intergenic RNAs whose 3' ends are extremely A-rich (Fig. 3C), which is a hallmark of mispriming artifacts [32]. In contrast, Bookend-defined 3' ends at both known and novel locations showed canonical PAS motifs. Similarly, known and novel Bookend 5' ends showed features distinct to transcription start sites, including the TATA-box and Y-patch (Additional file 1: Fig. S3B). Therefore, Bookend retains genuine novelty from long-read datasets by filtering against known sources of error.

Another major source of transcript assembly error is truncated 5' ends due to premature template switching during reverse transcription or amplification of degraded RNA. Although 79% of FLNC alignments fully matched the exon chain of a TAIR10 or Araport11 transcript, most were copies of a few highly expressed genes. After collapsing alignments into sets of unique exon chains, full-length reference transcripts accounted for only 31.4% of all unique chains, and 24.8% of unique chains were fragments of known transcript models, missing one or more exons (Additional file 1: Table S3). Clustering by Iso-seq3 removes some fragments, and they can be further reduced after alignment by collapsing 5' truncations with Transcript isoforms: Full-length and Unassembled (ToFU) [45] (Fig. 3D). To determine whether precision could be improved through assembly, the FLNC data was processed by StringTie2 or Bookend. StringTie2 yielded 12% fewer full-length reference matches than ToFU, but also reported 27% fewer transcripts that failed to match a reference (Fig. 3D,E, Additional file 1: Table S3). Bookend reported a 1% increase in reference matches over ToFU with a 45% reduction in non-matches. Because the Arabidopsis genome is compact with an average of only 1.5 kilobases (kb) between adjacent genes, assembly algorithms agnostic to 5' and 3' end information risk creating "fusions" of adjacent genes due to spurious read-through transcripts (Additional file 1: Fig. S3A). StringTie2 reported 838 fusion



**Fig. 3** Long-read sequencing is augmented by hybrid assembly. **A** Artifacts identified in PacBio FLNC reads from two SMRT cells by alignment to the Arabidopsis reference genome. **B** Nucleotide frequency enrichment in a  $\pm 50$  nt window around poly(A) sites (PAS) identified by Direct RNA-seq [30]. **C** Nucleotide enrichment around 3' ends of transcripts constructed from PacBio reads by Iso-seq3 (top), StringTie2 (middle), and Bookend (bottom) at sites overlapping a TAIR10 PAS (left), novel PAS at a known gene (middle), and novel antisense or intergenic loci (right); colors and scales as in (B). **D** Classification against the closest match in TAIR10 or Araport11 of transcripts constructed by four long-read processing strategies: Iso-seq3 clustering, cluster collapse by ToFU, and FLNC assembly by StringTie2 or Bookend. **E** Effect of long-read assembly on the number of transcripts by class (colored as in (D)) by StringTie2 (left) or Bookend (right) using hybrid assembly with one or more tissue-matched sequencing libraries. Bars show difference vs. ToFU-collapsed Iso-seq3 clusters. **F** UpSet plot depicting the number of overlapping transcript isoforms present in TAIR10, Araport11, and the Bookend Floral Bud assembly. **G** UpSet plot for the union of transcription start site peaks, allowing a  $\pm 50$  nt overlap window. **H** Diagram of floral organs analyzed by nanoPARE (left) and percent change in the number of TSS-overlapping reads for alternative annotations to TAIR10 in 15 tissue-specific nanoPARE libraries (right). **I** Bar chart of the number of enriched TSSs in each tissue type ( $\geq 4$ -fold mean RPM ingroup vs. outgroup, ANOVA  $p < 0.01$ , Benjamini-Hochberg correction). Shaded portions are TSSs exclusive to Bookend Floral Bud. **J** IGV browser image depicting Bookend Floral Bud assemblies for *MPK18* (top) and nanoPARE abundance heatmaps for *MPK18* sense and antisense TSSs (bottom). **K** Architecture and abundance of *MPK19* isoforms as in (K). **L** Full-length and truncated isoforms of *SKM1* assembled by Bookend (top). Heatmaps as in (J–K) for *SKM1* TSSs (bottom). TPM, transcripts per million; TSS, transcription start site

transcripts, 41% more than ToFU. By contrast, end-guided assembly of PacBio FLNCs with Bookend yielded 32% fewer fusions than ToFU and 52% fewer than StringTie2 while reporting more full-length matches than either (Fig. 3D, Additional file 1: Table S3).



Bookend's assembly model is general enough to combine reads from different sequencing strategies to produce a single "hybrid assembly." We used Bookend to assemble combinations of long reads (PacBio FLNCs), short reads (Smart-seq2), and transcript start site reads [16], and compared all transcript models against their closest match in either TAIR10 or Araport11 (Additional file 1: Table S3). Assembly was most sensitive when all three read types were combined, and concordance with reference annotations was higher for all Bookend hybrid assemblies than for all other methods. When short reads, long reads, and transcript start site reads were combined, Bookend could identify 2841 more reference-matching transcripts than ToFU (Fig. 3E). A recent update to StringTie2 implemented hybrid assembly of short and long reads and reported an improvement over long-read assembly alone [46]. Consistent with this report, StringTie2 hybrid assembly on the floral bud libraries was more sensitive than with only long reads, assembling 374 more matches than ToFU (Fig. 3E). However, StringTie2 also yielded more than twice as many transcript fragments as ToFU (Fig. 3E). Unlike StringTie2, Bookend can also integrate information from RNA 5' ends and requiring Cap Tags at transcript 5' ends during Bookend hybrid assembly yielded 30,219 transcript models with a 74.6% global concordance with the reference annotations (Additional file 1: Table S3). We report this hybrid assembly of long, short, and 5' end reads as the Bookend Floral Bud annotation (Fig. 3F, Additional file 2: Dataset 1-2) [16].

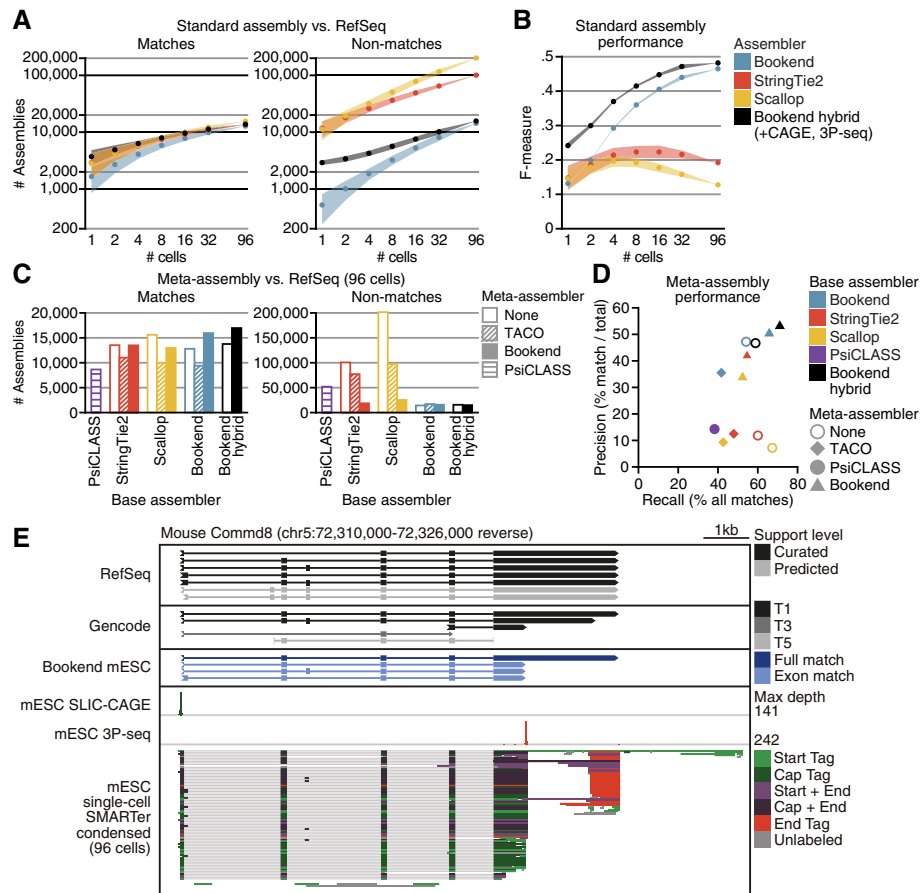
The Bookend Floral Bud annotation was assembled from RNA-seq of floral buds, which contain petal, sepal, anther, stigma, and ovule organs. To examine novel transcript models assembled by Bookend in more detail, we first quantified 15 previously published tissue-specific nanoPARE libraries (3 biological replicates from petals, sepals, anthers, stigma, or ovules that comprise floral buds) against TSSs from TAIR10, Araport11, or Bookend Floral Bud annotations. Although Bookend Floral Bud full-length isoforms had a greater overlap with Araport11, Bookend Floral Bud TSSs overlap more closely with TAIR10 (Fig. 3F,G). Araport11 had the largest set of TSSs, but they could only account for an average of 33.5% fewer nanoPARE reads than the TSS set from TAIR10 (Fig. 3H). This is consistent with reports that Araport11 TSSs are systematically placed too far upstream [16, 36]. The Bookend Floral Bud annotation possessed the smallest TSS set, but accounted for 20.7% more nanoPARE reads than TAIR10 on average.

Tissue-enriched TSSs were then calculated across the union of TAIR10, Araport11, and Bookend Floral Bud annotations. A tissue-enriched TSS was at least 4-fold more abundant with an ANOVA Benjamini-Hochberg-adjusted  $p$ -value  $< 0.01$  in either sepals, petals, anthers, stigma, or ovules relative to the other tissues (Additional file 2: Dataset 3). Previously unannotated TSSs account for 8.8% of all sites but 16.5% of tissue-enriched TSSs suggesting that TSSs that vary across tissue types tend to be missing from reference annotations. Novel TSSs are especially overrepresented in anthers, where 459 of 1932 anther-enriched TSSs were exclusive to the Bookend annotation (2.7-fold enrichment,  $p = 2.54e^{-90}$ , hypergeometric test, Fig. 3I). Sixty-three of these anther-enriched TSSs belong to unannotated antisense RNAs, including two transcripts running antisense to the Mitogen Activated Protein Kinase genes *MPK18* and *MPK19*, which are both inversely correlated with the abundance of their sense transcript (Fig. 3J,K). Bookend also uncovered a number of striking novel isoforms in anthers,

including transcripts with a TSS near the 3' end of the gene. One example is *STERILITY-REGULATING KINASE MEMBER1 (SKM1)*, a leucine-rich repeat receptor-like kinase involved in signaling between the pistil and pollen tube during pollination at high temperatures [47]. While the full-length isoform is detectable at a low level in sepals, petals, and ovules, a truncated isoform missing all leucine-rich repeats and the kinase domain is an order of magnitude more abundant in anthers (Fig. 3L). Therefore, end-guided assembly with Bookend enables the identification and initial characterization of tissue-specific transcript isoforms.

### Transcript discovery from single-cell sequencing

Bookend achieved comparable precision assembling *Arabidopsis* transcriptomes from either long reads or short reads generated by Smart-seq2, which is a protocol routinely used for single-cell RNA sequencing (scRNA-seq) (Additional file 1: Fig. S3C). However, scRNA-seq poses multiple hurdles to accurate assembly. Amplifying the few picograms of RNA in a single cell exacerbates biases and artifacts during reverse transcription [22], and dropouts from inefficient RNA capture place limits on accurate isoform quantification from scRNA-seq [48]. Additionally, scRNA-seq has been most widely adopted in the study of mammalian systems. The mouse genome (and likewise the human genome) is roughly 30 times larger than the *Arabidopsis* genome with an average of twice as many introns per gene and nearly three times the number of annotated isoforms. Additionally, mouse introns can exceed 100 kb and are on average 36 times longer than in *Arabidopsis*. Many isoforms per gene and large spans of non-genic sequence make it considerably more challenging both to assemble transcripts and to validate which assemblies are correct. To evaluate Bookend's utility on mammalian scRNA-seq data, we tested it on a dataset designed for single-cell benchmarking [49] which contains a set of synthetic Spike-In RNA Variants (SIRVs) added prior to cell lysis. SIRVs were designed to present a challenge to isoform quantification tools by mimicking complex mammalian genes [50]. The 69 synthetic transcripts map to 7 regions on a hypothetical genome in a way that recapitulates canonical and non-canonical splicing variation, antisense transcription, and alternative 5' and 3' ends with up to 18 isoforms per gene (Additional file 1: Fig. S4A). SIRV Mix E2 contains molecules in four discrete concentrations so that each locus has major and minor isoforms that vary in relative abundance by up to 128-fold. SMARTer library preparations from 96 single mouse embryonic stem cells (mESCs) were deeply sequenced, with an average of 7 million aligned paired-end 100 bp reads per cell (Additional file 1: Table S4) including an average of just over 500,000 SIRV-mapping reads per cell. Bookend correctly reconstructed (full splice match and  $\leq 100$  nt error on both ends) an average of 22.6 transcripts per cell, which was higher than either Scallop (16.3) or StringTie2 (13) (Fig. 4A,B). Moreover, Bookend assembled fewer false SIRVs than StringTie2 and especially Scallop (Fig. 4B). To test for a relationship between performance and sequencing depth, cells were progressively combined into pairs, then sets of 4, 16, 32, and all reads from the 96 cells. The relative performance of the three assemblers was stable over two orders of magnitude of input with the F-measure (harmonic mean of precision and recall) slightly rising for Bookend as the sequencing depth increased and slightly decreasing for the others (Fig. 4B). Importantly, Bookend consistently assigned a higher estimated



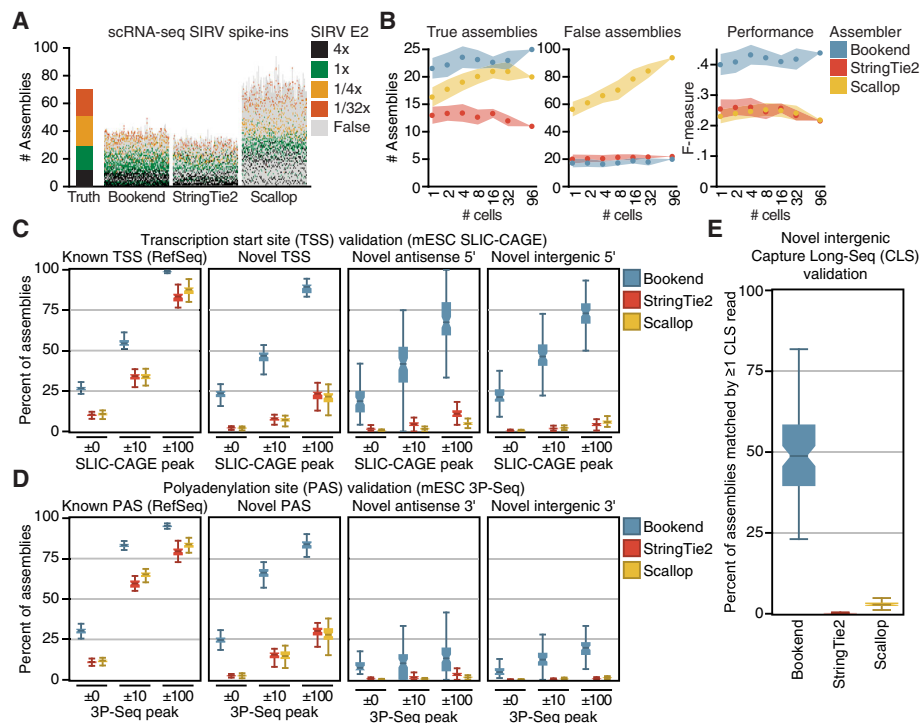
**Fig. 4** Bookend performance on single mouse cells. **A** Reconstruction of Spike-In RNA Variants (SIRVs) from 96 paired-end 100 bp SMARTer libraries of single mESCs. Each vertical bar depicts the assemblies from one cell, ordered from highest (bottom) to lowest (top) estimated abundance. Colored boxes match a true isoform of the given input concentration; gray boxes are false assemblies. **B** SIRV assembly performance as a function of increasing sequencing depth. F-measure (right) is the harmonic mean of sensitivity and precision. **C** Boxplots showing percent validation of 5' ends with SLIC-CAGE support within the given windows for 96 single mESC assemblies. **D** Boxplots as in (C) showing 3' end validation by 3P-Seq peaks. **E** Percent of intergenic assemblies (no overlap with RefSeq) in single cells which have  $\geq 1$  matching Capture Long-Seq read from the mouse CLS atlas

abundance to true transcripts, and false assemblies were more concentrated in the low abundance regime than for other assemblers (Fig. 4A). Overall precision on SIRVs averaged 55.9% for Bookend (vs. 39.6% StringTie2, 22.5% Scallop), and precision on the most abundant half of assemblies was 74.2% (vs. 48.2% StringTie2, 28.4% Scallop).

End-labeled reads mapping to the mouse genome were assembled for each cell using five different assemblers including the paired-end assembler TransComb [51]. Of the methods examined, StringTie2 was the fastest and most memory efficient, whereas Bookend had the second-lowest memory footprint and was comparable in processing time to TransComb (Additional file 1: Fig. S4B). To test assembly quality, transcript models were compared to RefSeq mm39. All matching exon chains were considered matches, and precision was measured as the percent of all assemblies that match RefSeq. Recall was defined by tallying all transcripts correctly assembled at least once and counting the proportion of this transcript set found per cell. Although recall was considerably lower for Bookend (average 7.8%) than other methods (StringTie2 13.9%,

Scallop 13.6%, TransComb 11.3%, Cufflinks 12.8%), precision was multiple times higher (Bookend 71.7%, StringTie2 21.2%, Scallop 18.1%, TransComb 14.6%, Cufflinks 14.5%) (Additional file 1: Fig. S4C). Most of Bookend's differences in sensitivity and precision can be attributed to discarding incomplete transcript models. If Bookend is instructed not to discard models that lack 5' and/or 3' end tags, then precision is reduced but sensitivity becomes comparable to other assemblers (Additional file 1: Fig. S4C).

As with TAIR10, RefSeq is almost certainly incomplete, and non-reference-matching assemblies could still be valid. To experimentally validate non-RefSeq mESC assemblies, three validation datasets were used: uuG-containing SLIC-CAGE [17] reads from mESCs for 5' end validation, mESC 3P-Seq [52] reads for 3' end validation, and a database of long noncoding RNAs identified by intergenic capture long-read sequencing (CLS [53]) for full-length validation of novel intergenic loci. An assembly was considered validated by a method if at least one read directly supported an assembled transcript's respective structure(s). Assemblies with 5' ends  $\leq 100$  nt away from a RefSeq TSS contained "known" TSSs, and all others possessed "novel" TSSs. Likewise, assemblies with 3' ends  $\leq 100$  nt from their matching reference polyadenylation sites were considered "known" PASs and all others were "novel". An average of 99.7% of Bookend, 83.9% of Scallop, and 79.0% of Stringtie2 single-cell assemblies with a known TSS had at least one SLIC-CAGE read within 100 nt (Fig. 5C). Moreover, the majority of novel, antisense, and intergenic TSSs



**Fig. 5** End-guided meta-assembly accurately integrates single-cell data. **A** Performance of assemblers with input from increasing numbers of single mESC cells. Assemblies with a matching exon chain to a RefSeq transcript (left) or no match to a RefSeq transcript (right). **B** F-measure of assemblies, where recall is the proportion of all transcripts assembled by  $\geq 1$  strategy and precision is matches/total assemblies. **C** Comparison of Bookend meta-assembly to standard assembly and other meta-assemblers. Number of RefSeq-matching transcripts assembled (left) or the number of non-matches (right). **D** Precision/recall plot of the 12 assemblies from **C**; recall and precision calculated as in **(B)**. **E** IGV browser image of the *Commd8* gene. From top to bottom: RefSeq, Gencode, and Bookend mESC annotations, 5' ends from mESC SLIC-CAGE, 3' ends from mESC 3P-seq, Bookend-condensed partial assemblies from 96 single mESCs

from Bookend transcripts were supported by at least 1 capped SLIC-CAGE read, whereas no novel group from StringTie2 or Scallop surpassed a 25% validation rate. The 3P-Seq dataset had fewer total reads and was less sensitive overall, but it still supported 19.9% of intergenic Bookend assembly 3' ends, compared to 1.4% for Scallop and 0.8% for StringTie2 (Fig. 5D). By comparing against the CLS atlas, we could validate the full structure of intergenic mESC assemblies. Bookend assembled a very small number of novel intergenic transcripts per cell (average 33 vs. 1209 by StringTie2 and 1073 by Scallop), but 49% of these were supported by one or more reads from the CLS atlas, compared to just 3% for Scallop intergenic assemblies and 0.3% for StringTie2 (Fig. 5E). Finally, because Cap and End Tags were extremely sparse in Droevendaalsesteeg 1, 6708 PB Wageningen, Netherlands cell (Additional file 1: Table S4), we hypothesized that the lower sensitivity could be explained by dropout of end labels. Supplying the mESC SLIC-CAGE (5' end) and 3P-seq (3' end) datasets to a Bookend hybrid assembly raised recall from 7.8 to 18.2% and retained a precision of 67.2% (Additional file 1: Fig. S4D). Assemblies were repeated for two replicates of Smart-seq2 data from the same experiment with comparable results, which demonstrates that end-guided assembly is consistent between two different full-length single-cell sequencing protocols (Additional file 1: Fig. S4D). Therefore, end-guided assembly of single-cell RNA-seq data can be used to identify genuine transcriptional novelty that is otherwise masked by noise.

### Condensed assembly and meta-assembly

A defining feature of single-cell experiments is that many individual cells are profiled in parallel. While sensitivity in an individual cell is low, information across multiple cells can be combined to achieve a more complete view of the experiment. Tools have been developed for transcript "meta-assembly" of reads from multiple sources. By modeling for variation across samples, meta-assemblers achieve higher precision than standard assembly on the same set of reads [54, 55]. To measure the impact of meta-assembly, a series of assemblies on subsamples of all 706 million aligned single-cell mESC reads was first performed with StringTie2 and Scallop, as well as Bookend with and without the addition of mESC SLIC-CAGE and 3P-seq libraries (Fig. 4A). The mean number of reference-matching transcripts varied greatly across assemblers on single cells (1656 Bookend, 3711 Bookend hybrid, 2904 StringTie2, 2831 Scallop), but the magnitude of difference decreased with progressive doublings, up to the full set of 96 cells (12,794 Bookend, 13,762 Bookend hybrid, 13,524 StringTie2, 15,611 Scallop). By contrast, non-matches grew linearly with input. Bookend consistently assembled roughly an order of magnitude fewer non-matching transcripts than other assemblers across all input levels. Scallop identified the most matches from the full 96-cell dataset, but this was dwarfed by nearly 13 times the number of assemblies that failed to match RefSeq (201,631 Scallop, 100,646 StringTie2, 14,301 Bookend, 15,711 Bookend hybrid). By assuming non-matches to be mostly false, we calculated recall and precision as before and combined them to track the relationship between overall performance (F-measure) and input. F-measure of Bookend and Bookend hybrid assembly continued to improve with increasing input, but Scallop and StringTie2 began to decline above 4 and 16 cells, respectively, due to the growth of non-matches outpacing matches (Fig. 4B). Consistent with previous reports, we see that standard assemblers suffer from an input-dependent decay in precision [54, 55].

As an alternative approach, two published meta-assemblers were used to process the 96-cell dataset. TACO builds a consensus annotation by re-defining transcript boundaries through “change-point detection” on a set of files from any standard assembler [54], whereas PsiCLASS generates the individual assemblies and performs meta-assembly through a consensus voting system [55]. The flexibility of Bookend’s framework allows its assembly algorithm to be run on assemblies, including its own output. To test the efficacy of meta-assembly with Bookend, each of the 96 single mESC datasets were “condensed” by a first pass through Bookend Assemble in which no incomplete transcripts were discarded (Additional file 1: Fig. S5A; “Path Filtering” section of Additional file 1: Supporting notes). Assembly was run again on the 96 condensed files, only retaining complete transcript models during the second pass. Bookend was also used to meta-assemble the 96 single-cell assemblies by StringTie2 and Scallop. Compared to standard assembly by StringTie2 or Scallop, all meta-assemblies produced substantially fewer non-matching transcripts (Fig. 4C). However, single-cell meta-assemblies surprisingly also recalled fewer RefSeq matches than standard assembly, with the exception of Bookend-to-Bookend and hybrid Bookend-to-Bookend meta-assemblies. PsiCLASS and TACO both showed somewhat higher precision than standard assembly, but at the expense of a severe drop in recall (Fig. 4D). PsiCLASS had the lowest recall of any method, but higher precision than StringTie2-to-TACO or Scallop-to-TACO meta-assembly. Bookend-to-Bookend meta-assembly considerably outperformed PsiCLASS in both recall (relative increase of 72%) and precision (relative increase of 253%). PsiCLASS produced an unusually large number of partial transcript fragments, likely due to the fact that scRNA-seq often has substantial 3’ bias that is not adequately accounted for (Additional file 1: Fig. S5A,B). Notably, when TACO was applied to single-cell Bookend assemblies, it showed both a 23% relative reduction in recall and a 25% relative reduction in precision compared to standard Bookend assembly. In contrast, Bookend-to-Bookend meta-assembly increased recall by 22% and precision by 7% (+ 58% recall and + 42% precision vs. Bookend-to-TACO). Across all three base assemblers, TACO reported fewer full reference matches than the standard assembly, while Bookend reported the same number or more full matches with a greater reduction in all non-matching classes than TACO (Additional file 1: Fig. S5C). Of all combinations tested, both sensitivity and precision were highest at the intron chain and full transcript level in a Bookend-to-Bookend hybrid meta-assembly in which SLIC-CAGE and 3P-seq data were supplied alongside the single-cell condensed assemblies [56] (Additional file 1: Table S5). We report this assembly as the “Bookend mESC” annotation (Additional file 2: Datasets 4,5). Requiring that both transcript ends are replicable across at least two different samples raised the transcript-level concordance with RefSeq to 54.1%, a relative increase of 271% over the most precise non-Bookend method (PsiCLASS), and a substantially higher agreement than even Gencode, an alternative mouse reference annotation that only shares 31.7% of its transcripts at assembled loci with RefSeq (Additional file 1: Fig. S5D). While Gencode isoforms contain a broader set of alternative TSS and PAS than RefSeq, we noticed that they can be contained in low-confidence or fragmented transcript models, as in the gene *Commd8* (Fig. 4E). By combining multiple unique advantages of end-guided assembly, Bookend could assemble more reference matches than any other strategy while maintaining a majority concordance with known annotations.

## Discussion

Computational gene annotation pipelines have long struggled to produce a reliable picture of plant and animal transcriptomes at the isoform level [11, 31, 57]. Studying the details of gene regulation and isoform usage remains restricted to a small number of model organisms in which manually curated accurate transcript models are available. Even with specialized methods for sequencing RNA ends, connecting those ends to a gene model can be computationally challenging, especially for noncoding RNAs [37]. By generating accurate end-to-end transcript assemblies from a range of widely accessible sequencing methods, Bookend enables the automated annotation of promoter architecture, alternative polyadenylation, and splicing dynamics in tissues in response to developmental, environmental, and disease state cues.

The utility of Bookend is limited by the availability of end-labeled RNA-seq data, which is only produced by a subset of all RNA-seq protocols. Fortunately, large-scale projects have been undertaken to catalog RNA ends, including thousands of human and mouse tissue-specific CAGE datasets from the FANTOM5 consortium [58]. Bookend assembly of unlabeled RNA-seq can be augmented by providing tissue-matched datasets of RNA ends. If 5' end data for a tissue of interest is missing, template switching protocols are straightforward alternatives to standard RNA-seq, and Smart-seq3 was developed to yield a far higher ratio of 5' labeled reads than Smart-seq2 [23].

## Conclusion

Despite rapid advancements in scale and sensitivity of single-cell RNA sequencing, the accurate detection of transcript isoforms is still an outstanding challenge [48]. Multiple approaches to apply long-read sequencing to single cells have been developed, but limits on throughput, error rate, and cost restrict their use [59–61]. Notably, large-scale Smart-seq2 experiments across multiple organisms have already been sequenced, including tens of thousands of cells from 20 mouse tissues and 24 human tissues by the Tabula Muris and Tabula Sapiens Consortia, respectively [62, 63]. Through meta-assembly of existing and future scRNA-seq datasets, we envision that Bookend will enable the comprehensive reannotation of transcriptomes at single-cell resolution.

## Methods

### PacBio sequencing

Two PacBio Iso-seq libraries were generated each using 10 µg of total RNA from *Arabidopsis* inflorescences containing unopened floral buds. Total RNA was extracted with TRIzol following the method described in Schon et al. [16] to yield two biological replicates with an RNA integrity number (RIN) of 9.0 and 9.2, respectively. SMRTbell libraries were constructed by the Vienna BioCenter Core Facilities (VBCF) and sequenced on a Sequel SMRT Cell 1M.

### Published RNA sequencing data

Smart-seq2 datasets from 5 ng *Arabidopsis* floral bud RNA and tissue-matched nano-PARE libraries from 10 µg total RNA were downloaded from the NCBI Gene Expression Omnibus (GEO), series accession GSE112869. Single-cell RNA-seq datasets of mouse embryonic stem cells and SIRVs from Natarajan et al. [49] were downloaded





*BED12*]. For cluster collapse, ToFU was installed with the *cDNA\_Cupcake* package ([https://github.com/Magdoll/cDNA\\_Cupcake](https://github.com/Magdoll/cDNA_Cupcake)), and the function *collapse\_isoforms\_by\_sam.py* was run on a sorted BAM file of all FLNC reads with the arguments `-fnc_coverage 2 -max_5_diff 200 -max_3_diff 200`. Additionally, aligned unsorted SAM files were converted to ELR with the command `bookend elr --stranded -s -e --start_seq ATGGG --genome [TAIR10.fa]`.

### Assembly

To make assembly setting maximally uniform across Bookend, StringTie2, Scallop, and Cufflinks, the following arguments were used. For Arabidopsis assemblies: `bookend assemble --max_gap 50 --min_cov 2 --min_len 60 --min_proportion 0.02 --min_overhang 3 --cap_bonus 5 --cap_filter 0.02; stringtie -g 50 -c 2 -m 60 -f 0.02 -a 3 -M 1 -s 5; scallop --min_bundle_gap 50 --min_transcript_coverage 2 --min_transcript_length_base 60 --min_flank_length 3 --min_single_exon_coverage 5 --min_transcript_length increase 50; cufflinks -F 0.02 --overhang-tolerance 3 --min-frags-per-transfrag 10 -j 0.15 -A 0.06`. For mouse assemblies, the same settings were used with the following exceptions: `--min_proportion` was set to 0.01, `--min_len` to 200, and `--require_cap` was enforced on mouse assemblies except when assembling spike-in transcripts, which do not possess caps. For meta-assembly, Bookend was run with the same settings as above for mouse. TACO was run with the arguments `--filter-min-expr 2 --filter-min-length 200 --isoform-frac 0.01`, and PsiCLASS was run with default settings

### Assembly algorithms

A brief overview of the end-guided assembly process implemented in Bookend is below. For a full breakdown of the algorithms used, see the “Bookend Algorithms” section of Additional file 1: Supporting notes. For detailed instructions on using the Bookend software package, see Additional file 3: Bookend user guide.

*(Generate Chunks)* First, reads are streamed in from an ELR file in sorted order and separated into overlapping chunks. *(Tag Clustering)* In each chunk, Start Tags and End Tags are clustered on each strand by grouping tags by genomic position and assigning each position a signal score of counts  $\times$  proportion of total coverage. A signal threshold is set and positions below the threshold are discarded. Remaining positions are grouped within a user-specified distance to yield Start and End clusters on each strand. *(Calculate Membership Matrix)* Start/End clusters are added to a catalog of boundaries, which include splice donor/acceptor sites that are also filtered by a threshold of total overlapping coverage. Adjacent boundary pairs define a “frag”, and each read is assigned a Membership array that describes whether the read overlaps or excludes each frag. Redundant membership arrays are combined, and the unique set of elements is stored as the Membership Matrix. *(Calculate Overlap Matrix)* A matrix describing the relationship between each element pair *1* and *2* is generated by asking (from left to right in genomic coordinates): can *1* extend into *2*? Can *2* extend into *1*? More formally, *2* extends *1* if they disagree in no column of the Membership Matrix and there exists at least one instance of a run of non-zero membership in *2* that ends at a higher index than the overlapping run in *1*. Each comparison returns a pair of Overlaps,  $O_{1,2}$  and  $O_{2,1}$ , respectively: 1 = extends, -1 = excludes, 2 = is contained by, 0 = does not overlap.

The values  $-1$  and  $0$  are symmetric, but  $1$  and  $2$  are directed relationships that can be used as edges in a directed graph. (*Collapse Linear Chains*) It is possible to identify and collapse non-branching sets of elements (“linear chains”) prior to assembly. Two graphs are constructed with elements as nodes: a directed graph with extensions as edges, and an undirected graph with exclusions as edges. A depth-first search is conducted by visiting each element in increasing order of information content (number of non-zero memberships). During a visit, the element’s edges are traversed recursively to record all traversed nodes’ exclusions. An element with no edges is assigned to a new chain. Otherwise, when an element’s edges are all traversed, the element is compared against its outgroup, the set of all elements reached. If all outgroup elements belong to one chain and the element and outgroup have the same set of exclusions, then the element is added to the same chain. If the element’s outgroup is assigned to multiple chains, the element begins a new chain. After completion of the search, each chain is combined to form a single reduced element. (*Generate Overlap Graph*) From the set of reduced elements, a second directed graph is constructed with a global source (Start+/End-) and sink (Start-/End+), where each node records the element weight (sequenced bases / genomic length), outgroup (extends to), ingroup (extends from), containments and exclusions. (*Resolve Containment*) All elements contained by one or more longer elements have their weight redistributed proportionally to their containers as long as not all containers exclude any single node the element does not already exclude. (*Greedy Paths*) All elements begin unassigned. Starting with the heaviest unassigned element, choose an extension (ingroup/outgroup pair) that maximizes a score that equally combines the following: maximal weight of the extension, maximal similarity of coverage distribution across samples between element and extension, minimal coverage variance across covered frags, and does not cause the source or sink to become unreachable. The highest-scoring extension is iteratively added to a path until both source and sink are reached or no further extensions are possible. Paths are generated in this manner until the total weight of unassigned elements falls below a given signal threshold.

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-022-02700-3>.

**Additional file 1: Figure S1.** The Bookend workflow. **Figure S2.** Nucleotide-level precision of Arabidopsis assembly 5’ and 3’ ends. **Figure S3.** Artifacts in long-read data. **Figure S4.** Single mESC assembly details. **Figure S5.** Meta-assembly details. **Table S1.** Floral bud Smart-seq2 end-labeled read mapping statistics. **Table S2.** Long-read validation of floral bud assemblies by class. **Table S3.** Floral bud hybrid assembly details. **Table S4.** End-labeling and alignment of single mESCs. **Table S5.** GffCompare performance statistics for mESC meta-assemblies. **Supporting notes.** Detailed information about End Labeled Read file format and assembly algorithms.

**Additional file 2: Dataset 1.** Bookend Floral Bud, hybrid assembly of Arabidopsis stage 12 inflorescence. **Dataset 2.** Classification of Bookend Floral Bud transcripts against TAIR10 and Araport11. **Dataset 3.** Tissue-specific quantification of TSSs from TAIR10, Araport11 and Bookend Floral Bud annotations. **Dataset 4.** Bookend mESC, hybrid assembly of mouse embryonic stem cells. **Dataset 5.** Classification of Bookend mESC transcripts against RefSeq and Gencode.

**Additional file 3.** Bookend user guide: Detailed guidelines and instructions for using all utilities in the Bookend software package.

**Additional file 4.** Review history.

## Acknowledgements

We thank the Next Generation Sequencing Facility at Vienna BioCenter Core Facilities GmbH (VBCF) for their outstanding services and technical support.

## Peer review information

Andrew Cosgrove was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

**Authors' contributions**

M.A.S. and M.D.N. conceived the project; M.A.S. developed the methodology; M.A.S. and S.L. performed the experiments; M.A.S. and F.H. analyzed data; M.A.S. prepared figures; M.A.S. wrote the article; M.A.S. and M.D.N. edited the article; M.D.N. acquired funding and supervised the project. All authors read and approved the final manuscript.

**Authors' information**

Twitter handles: @MSchon0 (Michael A. Schon); @mnode1 (Michael D. Nodine).

**Funding**

This work was supported by funding from the European Research Council under the European Union's Horizon 2020 research and innovation program (Grant 637888 to M.D.N.) and the DK Graduate Program in RNA Biology (DK-RNA) sponsored by the Austrian Science Fund (FWF, DK W 1207-B09).

**Availability of data and materials**

Bookend software is available on the Python Package Index and can be installed with the command `pip install bookend-ma`. Source code is available as a repository on GitHub at <https://github.com/Gregor-Mendel-Institute/bookend> under the MIT open-source license and is also available at Zenodo <https://doi.org/https://doi.org/10.5281/zenodo.6486387> [66]. All sequencing data generated in this study have been submitted to the National Center for Biotechnology Information Gene Expression Omnibus (NCBI GEO, <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE189482 [67].

**Declarations****Ethics approval and consent to participate**

Not applicable.

**Competing interests**

The authors declare that they have no competing interests.

Received: 6 December 2021 Accepted: 5 June 2022

Published online: 29 June 2022

**References**

- Liu F, Marquardt S, Lister C, Swiezewski S, Dean C. Targeted 3' processing of antisense transcripts triggers Arabidopsis FLC chromatin silencing. *Science*. 2010;327:94–7.
- Rhinn H, Qiang L, Yamashita T, Rhee D, Zolin A, Vanti W, et al. Alternative  $\alpha$ -synuclein transcript usage as a convergent mechanism in Parkinson's disease pathology. *Nat Commun*. 2012;3:1084.
- Solana J, Irimia M, Ayoub S, Orejuela MR, Zywitza V, Jens M, et al. Conserved functional antagonism of CELF and MBNL proteins controls stem cell-specific alternative splicing in planarians. *Elife*. 2016;5:e16797. <https://doi.org/10.7554/eLife.16797>.
- Mudge JM, Harrow J. The state of play in higher eukaryote gene annotation. *Nat Rev Genet*. 2016;17:758–72.
- Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res*. 2019;47:D766–73.
- McGarvey KM, Goldfarb T, Cox E, Farrell CM, Gupta T, Joardar VS, et al. Mouse genome annotation by the RefSeq project. *Mamm Genome*. 2015;26:379–90.
- Berardini TZ, Reiser L, Li D, Mezheritsky Y, Muller R, Strait E, et al. The Arabidopsis information resource: making and mining the "gold standard" annotated reference plant genome. *Genesis*. 2015;53:474–85.
- FANTOM Consortium and the RIKEN PMI and CLST (DGT), Forrest ARR, Kawaji H, Rehli M, Baillie JK, de Hoon MJL, et al. A promoter-level mammalian expression atlas. *Nature*. 2014;507:462–70.
- Wu P-Y, Phan JH, Wang MD. Assessing the impact of human genome annotation choice on RNA-seq expression estimates. *BMC Bioinformatics*. 2013;14(Suppl 11):S8.
- SEQC/MAQC-III Consortium. A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nat Biotechnol*. 2014;32:903–14.
- Guigó R, Flicek P, Abril JF, Reymond A, Lagarde J, Denoeud F, et al. EGASP: the human ENCODE Genome Annotation Assessment Project. *Genome Biol*. 2006;7(Suppl 1):S2.1–31.
- Stark R, Grzelak M, Hadfield J. RNA sequencing: the teenage years. *Nat Rev Genet*. 2019;20:631–56.
- Levin JZ, Yassour M, Adiconis X, Nusbaum C, Thompson DA, Friedman N, et al. Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nat Methods*. 2010;7:709–15.
- Murata M, Nishiyori-Sueki H, Kojima-Ishiyama M, Carninci P, Hayashizaki Y, Itoh M. Detecting expressed genes using CAGE. *Methods Mol Biol*. 2014;1164:67–85.
- Adiconis X, Haber AL, Simmons SK, Levy Moonshine A, Ji Z, Busby MA, et al. Comprehensive comparative analysis of 5'-end RNA-sequencing methods. *Nat Methods*. 2018;15:505–11.
- Schon MA, Kellner MJ, Plotnikova A. NanoPARE: parallel analysis of RNA 5' ends from low-input RNA. *Genome Res*. 2018;28:1931–42 Available from: <https://genome.cshlp.org/content/28/12/1931.short>.
- Cvetesic N, Leitch HG, Borkowska M, Müller F, Carninci P, Hajkova P, et al. SLIC-CAGE: high-resolution transcription start site mapping using nanogram-levels of total RNA. *Genome Res*. 2018;28:1943–56.
- Jan CH, Friedman RC, Ruby JG, Bartel DP. Formation, regulation and evolution of *Caenorhabditis elegans* 3'UTRs. *Nature*. 2011;469:97–101.
- Moll P, Ante M, Seitz A, Reda T. QuantSeq 3' mRNA sequencing for RNA quantification. *Nat Methods*. 2014;11:i–iii.
- Pelechano V, Wei W, Steinmetz LM. Extensive transcriptional heterogeneity revealed by isoform profiling. *Nature*. 2013;497:127–31.
- Wang J, Li B, Marques S, Steinmetz LM, Wei W, Pelechano V. TIF-Seq2 disentangles overlapping isoforms in complex human transcriptomes. *Nucleic Acids Res*. 2020;48:e104.

22. Picelli S, Björklund ÅK, Faridani OR, Sagasser S, Winberg G, Sandberg R. Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods*. 2013;10:1096–8.
23. Hagemann-Jensen M, Ziegenhain C, Chen P, Ramsköld D, Hendriks G-J, Larsson AJM, et al. Single-cell RNA counting at allele and isoform resolution using Smart-seq3. *Nat Biotechnol*. 2020;38:708–14.
24. Cao J, Spielmann M, Qiu X, Huang X, Ibrahim DM, Hill AJ, et al. The single-cell transcriptional landscape of mammalian organogenesis. *Nature*. 2019;566:496–502.
25. Cao J, Packer JS, Ramani V, Cusanovich DA, Huynh C, Daza R, et al. Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science*. 2017;357:661–7.
26. Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. *Nat Commun*. 2017;8:14049.
27. Garalde DR, Snell EA, Jachimowicz D, Sipos B, Lloyd JH, Bruce M, et al. Highly parallel direct RNA sequencing on an array of nanopores. *Nat Methods*. 2018;15:201–6.
28. Wan Y, Liu X, Zheng D, Wang Y, Chen H, Zhao X, et al. Systematic identification of intergenic long-noncoding RNAs in mouse retinas using full-length isoform sequencing. *BMC Genomics*. 2019;20:559.
29. Cumbie JS, Ivanchenko MG, Megraw M. NanoCAGE-XL and CapFilter: an approach to genome wide identification of high confidence transcription start sites. *BMC Genomics*. 2015;16:597.
30. Sherstnev A, Duc C, Cole C, Zacharaki V, Hornyik C, Oszolak F, et al. Direct sequencing of *Arabidopsis thaliana* RNA reveals patterns of cleavage and polyadenylation. *Nat Struct Mol Biol*. 2012;19:845–52.
31. Steijger T, Abril JF, Engström PG, Kokocinski F, RGASP Consortium, Hubbard TJ, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods*. 2013;10:1177–84.
32. Tardaguila M, de la Fuente L, Marti C, Pereira C, Pardo-Palacios FJ, Del Risco H, et al. SQANTI: extensive characterization of long-read transcript sequences for quality control in full-length transcriptome identification and quantification. *Genome Res*. 2018;28(3):396–411. <https://doi.org/10.1101/gr.222976.117>.
33. Kuo RI, Tseng E, Eory L, Paton IR, Archibald AL, Burt DW. Normalized long read RNA sequencing in chicken reveals transcriptome complexity similar to human. *BMC Genomics*. 2017;18:323.
34. Tang AD, Soulette CM, van Baren MJ, Hart K, Hrabeta-Robinson E, Wu CJ, et al. Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals downregulation of retained introns. *Nat Commun*. 2020;11:1438.
35. Kovaka S, Zimin AV, Perteau GM, Razaghi R, Salzberg SL, Perteau M. Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome Biol*. 2019;20:278.
36. Thieffry A, Vigh ML, Bornholdt J, Ivanov M, Brodersen P, Sandelin A. Characterization of *Arabidopsis thaliana* promoter bidirectionality and antisense RNAs by inactivation of nuclear RNA decay pathways. *Plant Cell*. 2020;32:1845–67.
37. de Rie D, Abugessaisa I, Alam T, Arner E, Arner P, Ashoor H, et al. An integrated expression atlas of miRNAs and their promoters in human and mouse. *Nat Biotechnol*. 2017;35:872–8.
38. Perteau M, Perteau GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*. 2015;33:290–5.
39. Shao M, Kingsford C. Accurate assembly of transcripts through phase-preserving graph decomposition. *Nat Biotechnol*. 2017;35:1167–9.
40. Trapnell C, Roberts A, Goff L, Perteau G, Kim D, Kelley DR, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc*. 2012;7:562–78.
41. Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. Opportunities and challenges in long-read sequencing data analysis. *Genome Biol*. 2020;21:30.
42. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-generation sequencing technologies. *Nat Rev Genet*. 2016;17:333–51.
43. Balázs Z, Tombác D, Csabai Z, Moldován N, Snyder M, Boldogkői Z. Template-switching artifacts resemble alternative polyadenylation. *BMC Genomics*. 2019;20:824.
44. Tang DTP, Plessy C, Salimullah M, Suzuki AM, Calligaris R, Gustincich S, et al. Suppression of artifacts and barcode bias in high-throughput transcriptome analyses utilizing template switching. *Nucleic Acids Res*. 2013;41:e44.
45. Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, et al. Widespread polycistronic transcripts in fungi revealed by single-molecule mRNA sequencing. *PLoS One*. 2015;10:e0132628.
46. Shumate A, Wong B, Perteau G, Perteau M. Improved transcriptome assembly using a hybrid of long and short reads with StringTie. *PLoS Comput Biol*. 2022;18(6):e1009730. <https://doi.org/10.1101/2021.12.08.471868>.
47. Endo S, Shinohara H, Matsubayashi Y, Fukuda H. A novel pollen-pistil interaction conferring high-temperature tolerance during reproduction via CLE45 signaling. *Curr Biol*. 2013;23:1670–6.
48. Westoby J, Artemov P, Hemberg M, Ferguson-Smith A. Obstacles to detecting isoforms using full-length scRNA-seq data. *Genome Biol*. 2020;21:74.
49. Natarajan KN, Miao Z, Jiang M, Huang X, Zhou H, Xie J, et al. Comparative analysis of sequencing technologies for single-cell transcriptomics. *Genome Biol*. 2019;20:70.
50. Paul L, Kubala P, Horner G, Ante M, Holländer I, Alexander S, et al. SIRVs: Spike-In RNA Variants as external isoform controls in RNA-sequencing. *bioRxiv*. 2016:080747 Available from: <https://www.biorxiv.org/content/10.1101/080747v1.abstract>.
51. Liu J, Yu T, Jiang T, Li G. TransComb: genome-guided transcriptome assembly via combing junctions in splicing graphs. *Genome Biol*. 2016;17:213.
52. Nam J-W, Rissland OS, Koppstein D, Abreu-Goodger C, Jan CH, Agarwal V, et al. Global analyses of the effect of different cellular contexts on microRNA targeting. *Mol Cell*. 2014;53:1031–43.
53. Lagarde J, Uszczyńska-Ratajczak B, Carbonell S, Pérez-Lluch S, Abad A, Davis C, et al. High-throughput annotation of full-length long noncoding RNAs with capture long-read sequencing. *Nat Genet*. 2017;49:1731–40.
54. Niknafs YS, Pandian B, Iyer HK, Chinnaiyan AM, Iyer MK. TACO produces robust multisample transcriptome assemblies from RNA-seq. *Nat Methods*. 2017;14:68–70.
55. Song L, Sabuncian S, Yang G, Florea L. A multi-sample approach increases the accuracy of transcript assembly. *Nat Commun*. 2019;10:5000.

56. Pertea G, Pertea M. GFF Utilities: GffRead and GffCompare. *F1000Res*. 2020;9:304.
57. Wang B, Tseng E, Regulski M, Clark TA, Hon T, Jiao Y, et al. Unveiling the complexity of the maize transcriptome by single-molecule long-read sequencing. *Nat Commun*. 2016;7:11708.
58. Noguchi S, Arakawa T, Fukuda S, Furuno M, Hasegawa A, Hori F, et al. FANTOM5 CAGE profiles of human and mouse samples. *Sci Data*. 2017;4:170112.
59. Gupta I, Collier PG, Haase B, Mahfouz A, Joglekar A, Floyd T, et al. Single-cell isoform RNA sequencing characterizes isoforms in thousands of cerebellar cells. *Nat Biotechnol*. 2018; Available from: <https://doi.org/10.1038/nbt.4259>.
60. Philpott M, Watson J, Thakurta A, Brown T Jr, Brown T Sr, Oppermann U, et al. Nanopore sequencing of single-cell transcriptomes with scCOLOR-seq. *Nat Biotechnol*. 2021;39(12):1517–20. <https://doi.org/10.1038/s41587-021-00965-w>.
61. Zheng YF, Chen ZC, Shi ZX, Hu KH, Zhong JY. HIT-sclSOseq: High-throughput and high-accuracy single-cell full-length isoform sequencing for corneal epithelium. *bioRxiv*. 2020; Available from: <https://www.biorxiv.org/content/10.1101/2020.07.27.222349v1.abstract>.
62. Tabula Muris Consortium, Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature*. 2018;562:367–72.
63. Quake SR, Sapiens Consortium T. The Tabula Sapiens: a single cell transcriptomic atlas of multiple organs from individual human donors. *bioRxiv*. 2021; Available from: <https://www.biorxiv.org/content/10.1101/2021.07.19.452956.abstract>.
64. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnetjournal*. 2011;17:10–2.
65. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013;29:15–21.
66. Schon MA, Lutzmayr S, Hofmann F, Nodine MD. Bookend [Internet]. 2022. Available from: <https://zenodo.org/record/6486387>
67. Nodine MD, Schon MA. Precise transcript reconstruction with end-guided assembly. *GSE189482*. *Gene Expression Omnibus*. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE189482> (2022).

## Publisher's Note

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

**Ready to submit your research? Choose BMC and benefit from:**

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

**At BMC, research is always in progress.**

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

