1    **Title:** Precise Transcript Reconstruction with End-Guided Assembly

2    **Running title:** End-Guided Assembly with Bookend

3

4    **Authors:** Michael A. Schon[1,2,*], Stefan Lutzmayer[1], Falko Hofmann[1] & Michael D. Nodine[1,2,*]

5

6    **Affiliations:** [1]Gregor Mendel Institute (GMI), Austrian Academy of Sciences, Vienna Biocenter

7    (VBC), Dr. Bohr-Gasse 3, 1030 Vienna, Austria; [2]Laboratory of Molecular Biology, Wageningen

8    University, Wageningen, 6708 PB, the Netherlands; *Correspondence

9

10    **Correspondence:** michael.nodine@wur.nl and michael.schon@wur.nl

11

12    **Key words:** RNA-seq, transcriptome, single-cell, TSS, PAS, capping, polyadenylation, 5′ and 3′

13    ends, long-read, Iso-Seq

14

15 **Summary statement:**

16 Bookend is a generalized framework that utilizes RNA 5′ and 3′ end information hidden in RNA-

17 seq datasets to accurately reconstruct transcriptomes including those from single cells.

18

19

20 **ABSTRACT**

21 Accurate annotation of transcript isoforms is crucial to understand gene functions, but automated

22 methods for reconstructing full-length transcripts from RNA sequencing (RNA-seq) data remain

23 imprecise. We developed Bookend, a software package for transcript assembly that incorporates

24 data from different RNA-seq techniques, with a focus on identifying and utilizing RNA 5′ and 3′

25 ends. Through end-guided assembly with Bookend we demonstrate that correct modeling of

26 transcript start and end sites is essential for precise transcript assembly. Furthermore, we

27 discovered that utilization of end-labeled reads present in full-length single-cell RNA-seq (scRNA-

28 seq) datasets dramatically improves the precision of transcript assembly in single cells. Finally,

29 we show that hybrid assembly across short-read, long-read, and end-capture RNA-seq datasets

30 from Arabidopsis, as well as meta-assembly of RNA-seq from single mouse embryonic stem cells

31 (mESCs) can produce end-to-end transcript annotations of comparable quality to reference

32 annotations in these model organisms.

**INTRODUCTION**

33

34 The functions of genes depend on the amount and types of RNA molecules that they produce.

35 Variation in transcript initiation, splicing and polyadenylation can generate an array of RNA

36 isoforms, and cataloging how these RNA variants change across development and disease

37 provides insights into corresponding gene functions[1–3]. Large-scale projects dedicated to the

38 manual curation of gene annotations are extremely valuable, but are labor-intensive and thus

39 limited in scope to the most well-studied organisms[4–7]. Moreover, multicellular organisms have

40 difficult-to-access cell types that will inevitably be overlooked by even the most comprehensive

41 annotation projects[8]. The completeness and accuracy of a reference annotation can considerably

42 impact all downstream data analyses, from gene expression to predictions of gene function[9–11].

43 To understand how transcriptome architecture varies during development and in response to

44 disease, it is therefore valuable to have an automated method that accurately identifies transcript

45 isoforms. Accordingly, many computational tools have been developed for genome annotation

46 including software that utilizes the massive and growing diversity of RNA sequencing (RNA-seq)

47 technologies[12].

48 A wide array of RNA-seq protocols have been developed to profile different aspects of the

49 transcriptome, from strand-specific coverage of gene bodies[13] to selective amplification of RNA

50 5′ ends[14–17], 3′ ends[18,19] or simultaneous capture of both ends[20,21]. Major recent advances have

51 enabled the amplification of full-length transcripts from single cells[22,23] or 3′ end capture from

52 millions of cells[24–26]. In parallel, advances have been made for profiling RNA on "third-generation"

53 long-read sequencing platforms such as PacBio and Oxford Nanopore single-molecule

54 sequencers that can read a continuous DNA and/or RNA molecule many times the length of a

55 typical transcript and yield end-to-end complete sequences of RNA molecules[27,28].

56 Transcript assembly is the effort to distill information from RNA-seq experiments into a

57 comprehensive annotation of the transcript isoforms present in the corresponding samples.

58 Depending on the method, RNA-seq reads contain a broad spectrum of information content. At

3

59 one extreme, single-end reads from a non-stranded RNA-seq protocol can be 50 nucleotides (nt)

60 or shorter and sequenced from one end of a double-stranded cDNA fragment such that the

61 resulting sequence is a random substring of an RNA molecule or its reverse complement. Paired-

62 end reads contain two ends of a cDNA molecule and typically there is a gap of unknown length

63 between the mate pairs. When aligned to a reference genome, paired reads may span more than

64 one splice junction, indicating that these splicing events occurred in the same molecule. Some

65 strand-specific RNA-seq protocols selectively sequence only first-strand or second-strand cDNA

66 to preserve knowledge of the original mRNA molecule's orientation[13]. Other protocols selectively

67 capture and sequence a fragment immediately downstream of the RNA 5′ end or upstream of the

68 3′ end, demarcating precisely where that molecule begins or ends, respectively. Finally, the most

69 information-rich reads come from long-read sequencing, in which the RNA or cDNA is read in its

70 entirety without fragmentation. Long-read methods are a promising tool for transcript annotation,

71 but current protocols are more error-prone per base sequenced, less sensitive, and more costly

72 than comparable short-read experiments. Because the vast majority of existing RNA-seq data is

73 in short-read format, nearly all assemblers have aimed to reconstruct transcripts from paired-end

74 short reads. A long-recognized problem of assemblers is the inaccurate annotation of transcript

75 start sites (TSS) and polyadenylation sites (PAS)[29,30]. Existing short-read assemblers infer TSSs

76 and PASs at sharp changes in read coverage, but such changes can also be due to alignment

77 errors, biased RNA fragmentation, sample degradation, or spurious intron retention. Long-read

78 sequencing methods are designed to read RNA from TSS to PAS, but they remain susceptible to

79 a variety of experimental artifacts[30]. The increasing adoption of long reads for transcript

80 annotation has led to a separate suite of tools that summarize, collapse, or "polish" long reads to

81 remove erroneous structures and present a set of representative isoforms from these reads[31,32].

82 For example, a recently developed transcript assembler reports the use of long reads in assembly

83 by removing aligned segments with a high error rate and assembling the resulting gapped reads[33].

84 Transcript annotation would ideally integrate information from a variety of RNA-seq methods to

85    determine the best evidence for transcript starts, ends and splicing patterns in a tissue-of-interest.

86    However, current transcriptome assembly methods do not employ information about where RNAs

87    begin and end. Here, we describe a method utilizing RNA 5′ and 3′ end information contained in

88    RNA-seq datasets to accurately reconstruct transcriptomes including those from single cells.
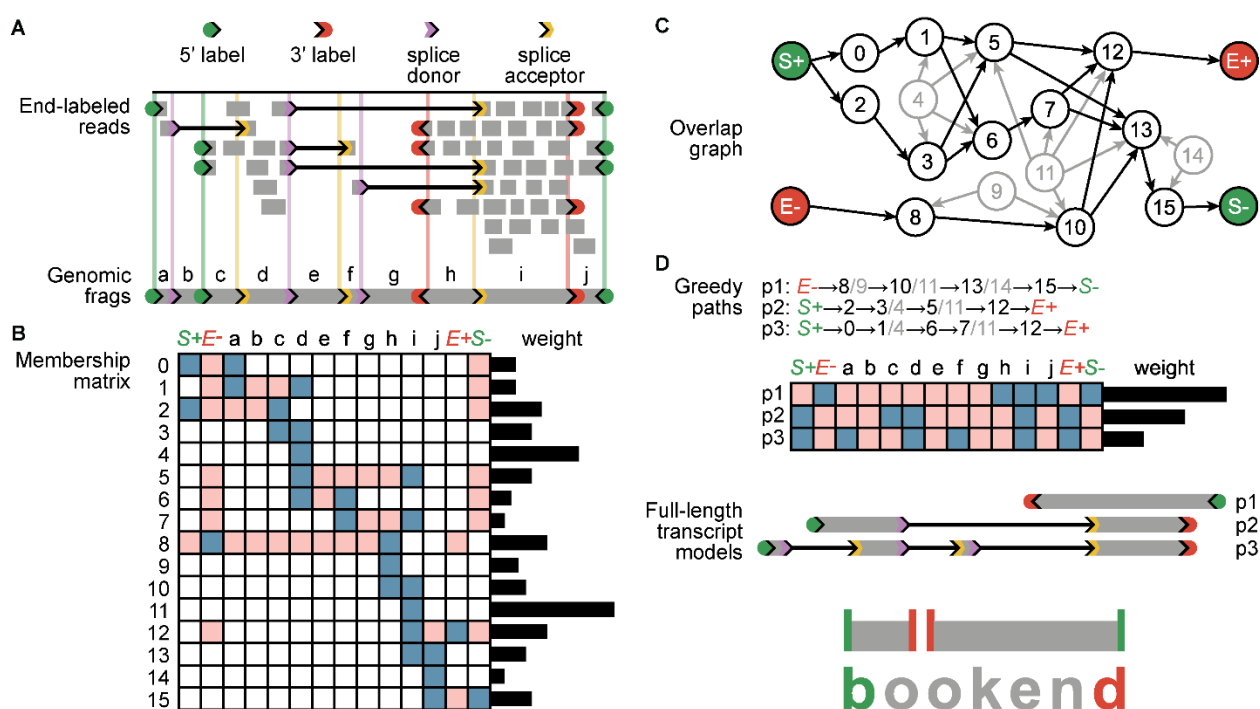
89

90    **RESULTS**

91    *A framework for end-guided transcript assembly*

92    To determine whether RNA 5′ and 3′ end information can improve transcript assembly algorithms,

93    we developed a generalized framework for identifying RNA ends in sequencing data and using

94    this information to assemble transcript isoforms as paths through a network accounting for splice

95    sites, transcription start sites (TSS) and polyadenylation sites (PAS). Because this software uses

96    end information to guide transcript assembly, we named it Bookend. Importantly, Bookend takes

97    RNA-seq reads from any method as input and after alignment to a reference genome, reads are

98    stored in a lightweight end-labeled read (ELR) file format that records all RNA boundary features

99    (5′ labels, splice donors, splice acceptors, gaps, 3′ labels), as well as the sample of origin for that

100   read (see Supporting Notes). Assembly is then resolved at each locus with aligned reads through

101   a four-step procedure (Fig1; see Methods and Supporting Notes). First, boundary labels from all

102   aligned RNA-seq reads are clustered and filtered to demarcate a unique set of locus TSSs, PASs

103   and splice junctions. Each locus is partitioned into a set of nonoverlapping "frags" defined as the

104   spans between adjacent boundary labels. Four additional frags (S+, E+, S-, E-) denote the

105   presence of a Start or End Tag on the forward or reverse strand. Second, a Membership Matrix

106   is generated to redefine all aligned reads with respect to the locus frags. A read's Membership

107   includes each frag it overlaps and excludes each incompatible frag (e.g. a spanned intron, a

108   region upstream of a TSS or downstream of a PAS). Reads with identical patterns of Membership

109   are condensed to a single element (row) of the Membership Matrix, whose weight is the total

110   coverage depth across the element by all reads of that pattern. Third, an Overlap Graph is
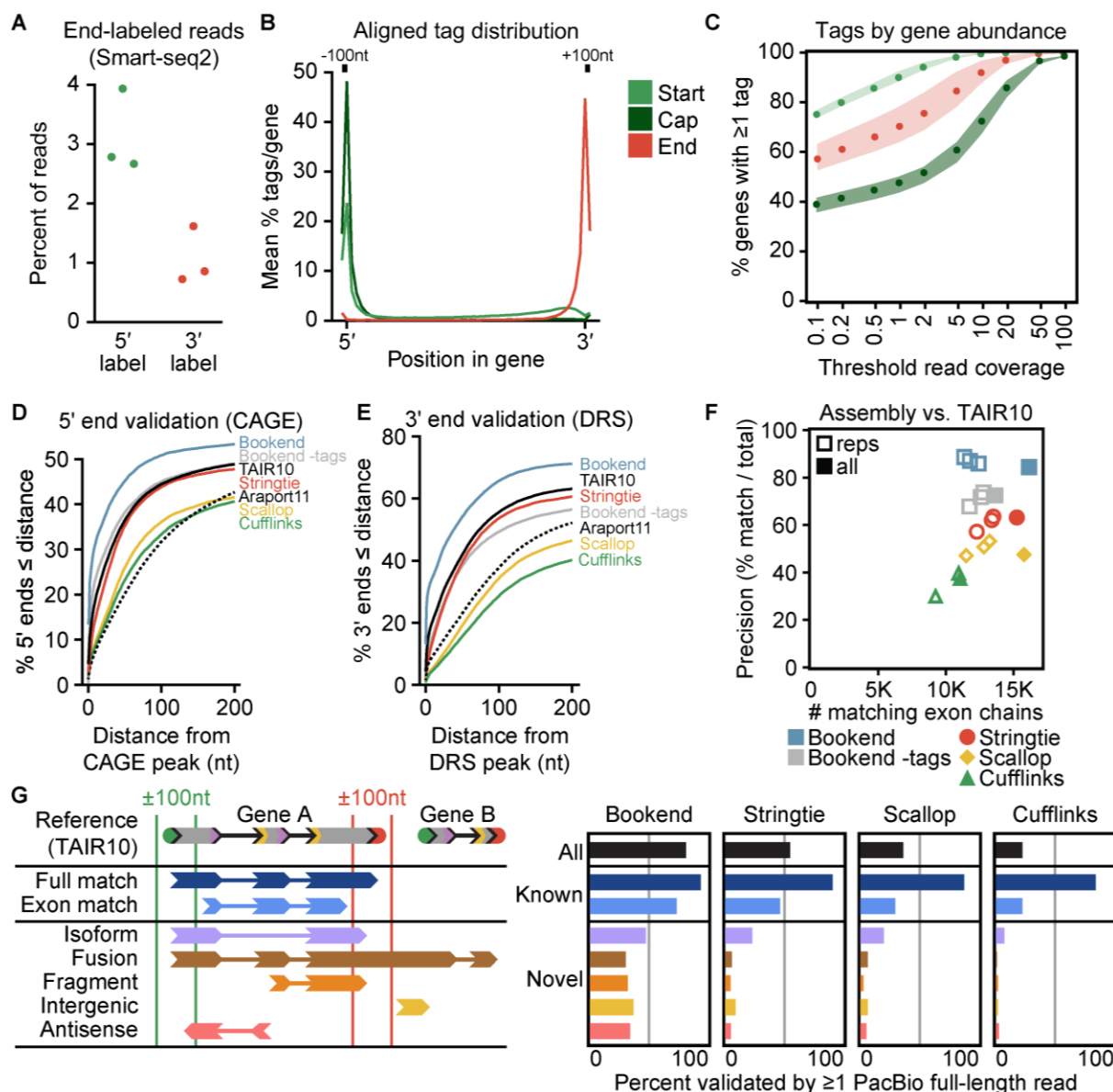
111     constructed from the Membership Matrix elements and this directed graph is simplified by

112     collapsing shorter elements into the elements that contain them. Finally, the Overlap Graph is

113     iteratively traversed to resolve an optimal set of Greedy Paths from TSSs to PASs. These Paths

114     describe a set of full-length transcript models best supported by the input reads. The Membership

115     Matrix definition is flexible enough to utilize reads regardless of their length, alignment gaps,

116     strand, or end information (FigS1B).

117



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

129

6

130     *End-labeled reads improve the quality of transcript assembly*

131     Arabidopsis is an ideal model to benchmark transcript assembly in higher eukaryotes. The

132     Arabidopsis genome is compact (~119 megabases), contains few repetitive elements, and the

133     TAIR10 reference annotation was extensively curated from expressed sequence tag (EST) data[7].

134     To determine whether assembly benefits from end-labeled reads, we examined libraries

135     generated with the low-input sequencing method Smart-seq2 from Arabidopsis floral buds[16]. Two

136     crucial steps in the Smart-seq2 protocol, template switching and preamplification, enrich for full-

137     length cDNA with an oligo label at both the 5′ (template switching oligo, TSO) and 3′ (oligo-dT)

138     end[22]. These oligos were trimmed from all reads and a record was kept of which end label was

139     found (5′, 3′, or no label) before mapping to the genome. As anticipated, a small percentage of

140     reads were found with either label (Fig 2A; Supplemental Table 1). All reads were aligned to the

141     Arabidopsis genome, and the terminal positions of 5′- and 3′-labeled reads were retained as "Start

142     Tags" and "End Tags", respectively. Of End Tags mapping to annotated genes, 88% mapped

143     near PASs, defined as the last decile of the gene or up to 100nt downstream (Fig2B). Start Tags

144     had lower specificity for TSSs, with only 48% of Start Tags in the first decile of genes or up to

145     100nt upstream. Template switching is known to readily occur at RNA 5′ ends derived from in vivo

146     or in vitro RNA decay. However, a subset of reads contain an intervening G between the TSO

147     and the genome-aligned sequence, indicating a 7-methylguanosine cap on the template

148     RNA[16,34,35]. The upstream untemplated G (uuG)-containing Start Tags were classified as Cap

149     Tags. Cap Tags were rare relative to all Start Tags (9%), but were much more specific to TSSs

150     with an average of 88% of Cap Tags within each gene mapping near the 5' end (Fig2B). To

151     optimize detection of true transcript 5' and 3' ends, the Tag Clustering algorithm designed for

152     Bookend defines Tag weight as a function of total read depth and applies a bonus to Cap Tags

153     over non-uuG Start Tags (See Supplemental Note: "Tag Clustering").

**154**

**Figure 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 100nt 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 ≥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[37] 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.

8

168    Despite end-labeled reads being relatively rare, the preamplification process should ensure

169    that a TSO or oligo-dT sequence is at each end of every cDNA molecule prior to tagmentation.

170    Therefore, we expected end-labeled reads to be distributed widely across the genome wherever

171    reads exist. As predicted, the majority of genes with >0 read coverage contained ≥1 Start Tag

172    and ≥1 End Tag, and the likelihood of finding a Start or End Tag increased as a function of total

173    read coverage (Fig2C). Of all genes with at least 1x, 10x and 100x read coverage, 73.3%, 94.4%

174    and 99.2% possessed both a Start and End Tag, respectively.

175    To assess whether end-labeled reads mark real TSSs and PASs at nucleotide precision,

176    Bookend was used to assemble all floral bud Smart-seq2 reads either with or without utilizing

177    Start and End Tags. Additionally, three leading short-read transcript assemblers were used with

178    comparable settings (see Methods): StringTie2[33,38], Scallop[39], and Cufflinks[40]. Publicly available

179    Arabidopsis CAGE[36] and Direct RNA-seq (DRS[37]) datasets were used to validate 5′ and 3′ ends,

180    respectively. All three of these widely-used assemblers output thousands of single-exon

181    unstranded fragments, which were ambiguous with regard to which end is 5' or 3' and thus were

182    discarded from further analyses (Supplemental Table 2). Bookend-defined TSSs based on

183    Start/Cap Tags were more likely to have a CAGE peak within 200nt than 5′ ends reported either

184    by Bookend without the use of Start Tags, the three leading assemblers, or even the current

185    Arabidopsis reference annotations (Fig2D). Likewise, a higher proportion of Bookend-identified

186    PASs were supported by DRS reads than PASs reported by the other transcript assemblers and

187    Arabidopsis reference annotations (Fig2E). At the nucleotide level, Bookend-defined transcript

188    boundaries were more than twice as likely to agree with the exact experimentally-determined TSS

189    and PAS peak positions than the most accurate reference annotation (TAIR10), while the other

190    three assemblers reported transcript boundaries less accurate than TAIR10 (FigS2A-B).

191    Strikingly, even the Bookend 5′ and 3′ ends >100nt from any reference still possessed known

192    sequence motifs associated with TSS and PAS, respectively, whereas sequence content around

193    novel ends from Cufflinks, Scallop, and StringTie2 is largely incoherent (FigS2C-D). In addition to

194  a dramatic increase in transcript boundary accuracy, 16,158 exon chains predicted by Bookend

195  fully matched a TAIR10 reference transcript, which was higher than when end-labeled reads were

196  ignored (13,660) and exceeded the totals from Scallop (15,785), StringTie2 (15,253) or Cufflinks

197  (11,051) (Fig2F). Therefore, Bookend correctly builds more known transcripts than other

198  assemblers and Bookend-annotated 5′ and 3′ ends were more precise than even the most

199  accurate Arabidopsis reference annotation.

200      In addition to known transcripts, Bookend constructed 2,979 isoforms not present in TAIR10,

201  which was 66% fewer than StringTie2 (8,886), 83% fewer than Scallop (17,400), and 84% fewer

202  than Cufflinks (18,934). An assembled transcript may fail to match TAIR10 either because the

203  assembly is incorrect or because the reference is incomplete. To distinguish between these

204  possibilities, two long-read SMRT cells of floral bud RNA were sequenced with the PacBio

205  platform to yield 547,910 full-length non-chimeric (FLNC) reads. All short-read assemblies were

206  partitioned into 7 different classifications based on their relationship to the most similar TAIR10

207  model (Fig2G). A transcript model was considered experimentally validated if at least one aligned

208  PacBio read fully matched the model (entire exon chain, ±100nt ends). Of all Bookend transcripts,

209  81.2% were supported by PacBio data, which surpassed the validation of transcripts predicted by

210  StringTie2 (54.7%), Scallop (35.9%) or Cufflinks (22.3%) (Fig2G; Supplemental Table 2).

211  Reference-matching transcripts have a higher average estimated abundance than non-reference

212  transcripts, making the latter more difficult to validate with the limited throughput of long-read

213  sequencing (FigS2E). Despite this limitation, 42.3% of non-reference Bookend assemblies were

214  fully supported by at least one PacBio read, which was substantially higher than the validation

215  rate of non-reference transcript assemblies generated by StringTie2 (15.9%), Scallop (11.6%),

216  and Cufflinks (4.3%) (Fig. 2G). Taken together, these results demonstrate that end-guided

217  assembly using latent RNA end information enables precise transcript reconstruction from short-
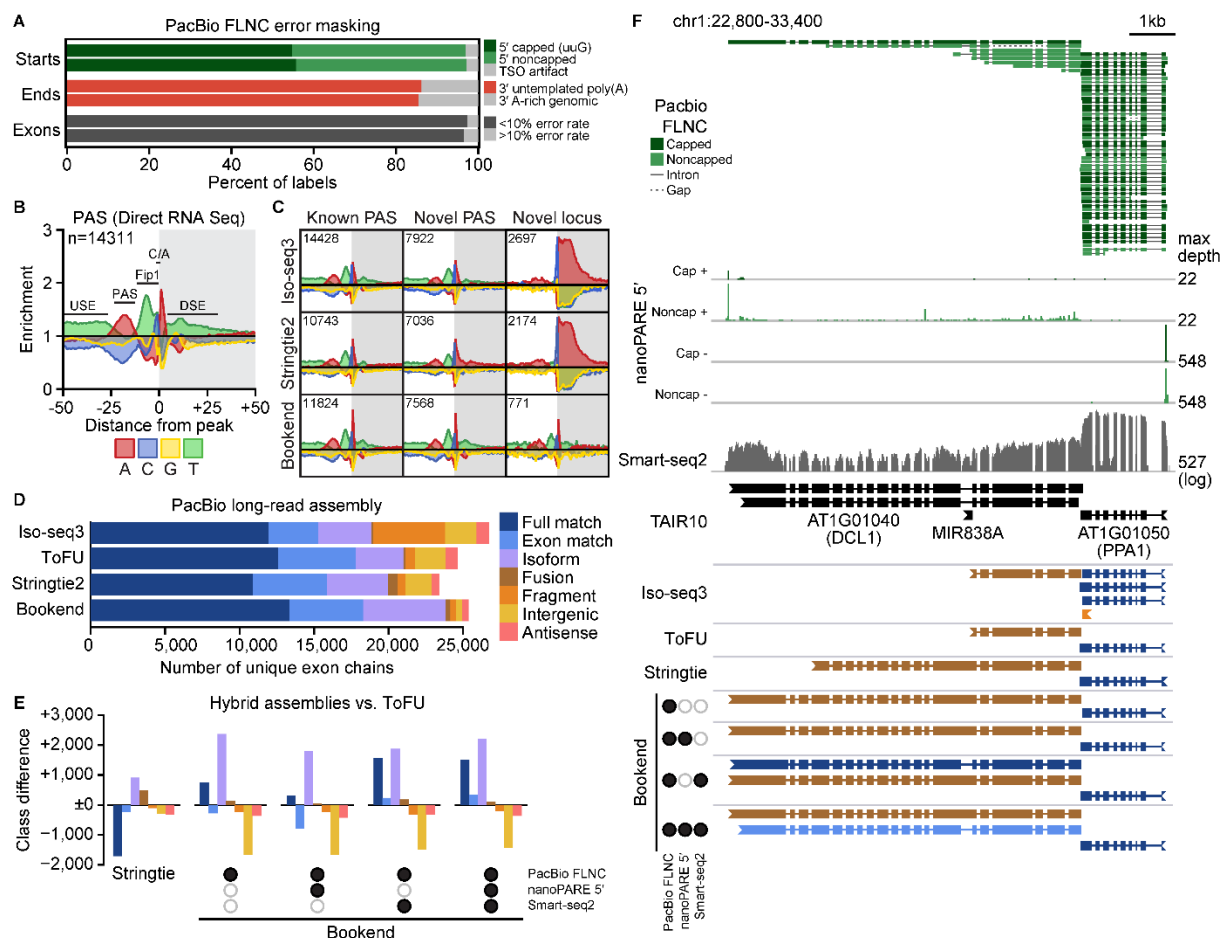
218  read datasets.

219    *Hybrid assembly refines and complements long-read RNA-seq*

220    Long-read sequencing technologies do not obviate the need for transcript reconstruction. Various

221    sources of technical and biological noise result in fragmented or improperly spliced long reads[30,41].

222    Long-read approaches also suffer from a higher base-level error rate compared to short-read

223    platforms[42]. Error correcting methods such as Circular Consensus Sequencing (CCS) require

224    reverse transcription and cDNA amplification, which are susceptible to mispriming and template-

225    switching artifacts[43,44]. This has driven the ongoing development of tools to refine transcript

226    models derived from long reads[31,32]. Additionally, StringTie2 was recently repurposed to assemble

227    long reads[33].

228        To quantify potential sources of error, PacBio FLNC reads were aligned to the genome and

229    processed by the Bookend pipeline to identify and remove template-switching artifacts, oligo-d(T)

230    mispriming events at A-rich regions, and exons with a high alignment error (Fig3A). Across both

231    SMRT cells, 95.4% of reads aligned successfully, and 97.0% of alignments did not contain any

232    high-error exons, defined as the total length of mismatches, inserts, and deletions exceeding 10%

233    of the exon length. However, 14.1% of all FLNC 3′ end labels were removed due to alignment

234    failure or the presence of an A-rich region immediately downstream of the oligo-d(T) junction. If

235    treated as genuine 3′ ends, these reads can cause false annotation of 3′-UTRs or putative

236    transcripts antisense or intergenic to known genes[43] (FigS3A). Direct RNA sequencing bypasses

237    oligo-d(T) priming and was used to produce a map of genuine Arabidopsis PAS[37]. These sites

238    show a distinct pattern of nucleotide enrichment, including a C/A dinucleotide motif at the

239    cleavage and polyadenylation site itself, and a U-rich upstream element (USE) and downstream

240    element (DSE) (Fig3B). Three tools were used to reduce the PacBio FLNC data into a unique set

241    of transcripts: the Iso-seq3 clustering algorithm from PacBio, assembly by StringTie2, and end-

242    guided assembly by Bookend. All 3 methods could recapitulate known PAS motifs at the set of 3′

243    ends within 100nt of a TAIR10-annotated PAS. In contrast to Bookend, StringTie2-annotated 3′

244    ends showed a slight A-richness at novel 3′ ends, and both Iso-seq3 and StringTie2 annotations

245    contain thousands of putative novel antisense or intergenic RNAs whose 3′ ends are extremely

246    A-rich (Fig3C). Therefore, Bookend retains genuine novel PAS by filtering against known 3′

247    artifacts.

248



250    **Figure 3.** Long-read sequencing is augmented by hybrid assembly
251    **(A)** Artifacts identified in PacBio FLNC reads from two SMRT cells by alignment to the Arabidopsis
252    reference genome. **(B)** Nucleotide frequency enrichment in a ±50nt window around poly(A) sites
253    (PAS) identified by Direct RNA Seq[37]. **(C)** Nucleotide enrichment around 3′ ends of transcripts
254    constructed from PacBio reads by Iso-seq3 (top), StringTie2 (middle), and Bookend (bottom) at
255    sites overlapping a TAIR10 PAS (left), novel PAS at a known gene (middle), and novel antisense
256    or intergenic loci (right); colors and scales as in **B**. **(D)** Classification against TAIR10 of transcripts
257    constructed by four long-read strategies: Iso-seq3 clustering, cluster collapse by ToFU, and FLNC
258    assembly by StringTie2 and Bookend. **(E)** Effect of long-read assembly on the number of transcripts
259    by class (colored as in **D**) by StringTie2 (left) or hybrid assembly with one or more tissue-matched
260    sequencing libraries by Bookend (right). Bars show difference vs. ToFU-collapsed Iso-seq3
261    clusters. **(F)** Integrative Genomics Viewer (IGV) image of the Arabidopsis *DICER-LIKE1* (*DCL1*)
262    locus. From top to bottom: PacBio reads colored by 5′ end label, nanoPARE capped and
263    noncapped read 5′ end frequency, Smart-seq2 read coverage, TAIR10 reference models, and long-
264    read assemblies colored by classification vs. TAIR10 as in **D**.

12

265

266      Another major source of transcript assembly error is truncated 5′ ends due to premature

267    template switching during reverse transcription or amplification of degraded RNA. Although 50%

268    of FLNC alignments matched a full-length TAIR10 transcript, most were copies of a few highly-

269    expressed genes. After collapsing alignments into sets of unique exon chains, full-length

270    reference transcripts accounted for only 18% of all unique chains, and 25% of unique chains were

271    fragments of known TAIR10 transcripts missing one or more exons (Supplemental Table 3).

272    Clustering by Iso-seq3 removes some fragments, and they can be further reduced after alignment

273    by collapsing 5′ truncations with Transcript isOforms: Full-length and Unassembled (ToFU) [45]

274    (Fig3D). However, it was unknown whether an assembly algorithm would further improve the

275    quality of long-read annotations. Surprisingly, passing the FLNC data through StringTie2 yielded

276    1,704 fewer full-length reference matches compared to ToFU, and the number of transcripts

277    classified as fusions of two different genes increased nearly four-fold (Fig3D-E). Because the

278    Arabidopsis genome is compact with an average of only 1.5 kilobases (kb) between adjacent

279    genes, assembly algorithms agnostic to 5′ and 3′ end information risk mis-annotating fused genes

280    due to spurious read-through transcripts (FigS3A). By contrast, end-guided assembly of PacBio

281    FLNCs with Bookend yielded 761 more full-length reference matches than ToFU, fewer than half

282    as many fusions as StringTie2, and over a thousand more putative novel isoforms than both

283    (Fig3D, Table S3).

284      Bookend's assembly model is general enough to mix reads from different sequencing

285    strategies. Therefore, we generated "hybrid assemblies" from combinations of PacBio FLNCs with

286    Smart-seq2 and/or nanoPARE (a 5' end sequencing strategy) from floral bud RNA[16]. All hybrid

287    assemblies had higher precision than assembling long reads alone, and up to 809 more full-length

288    matches could be identified (Fig3E, FigS3C, Supplemental Table 3). For example, *DICER-LIKE1*

289    (*DCL1*) encodes the Arabidopsis *Dicer* homolog required for microRNA biogenesis, and its mRNA

290    is maintained at low cellular abundance through an autoregulatory negative feedback loop
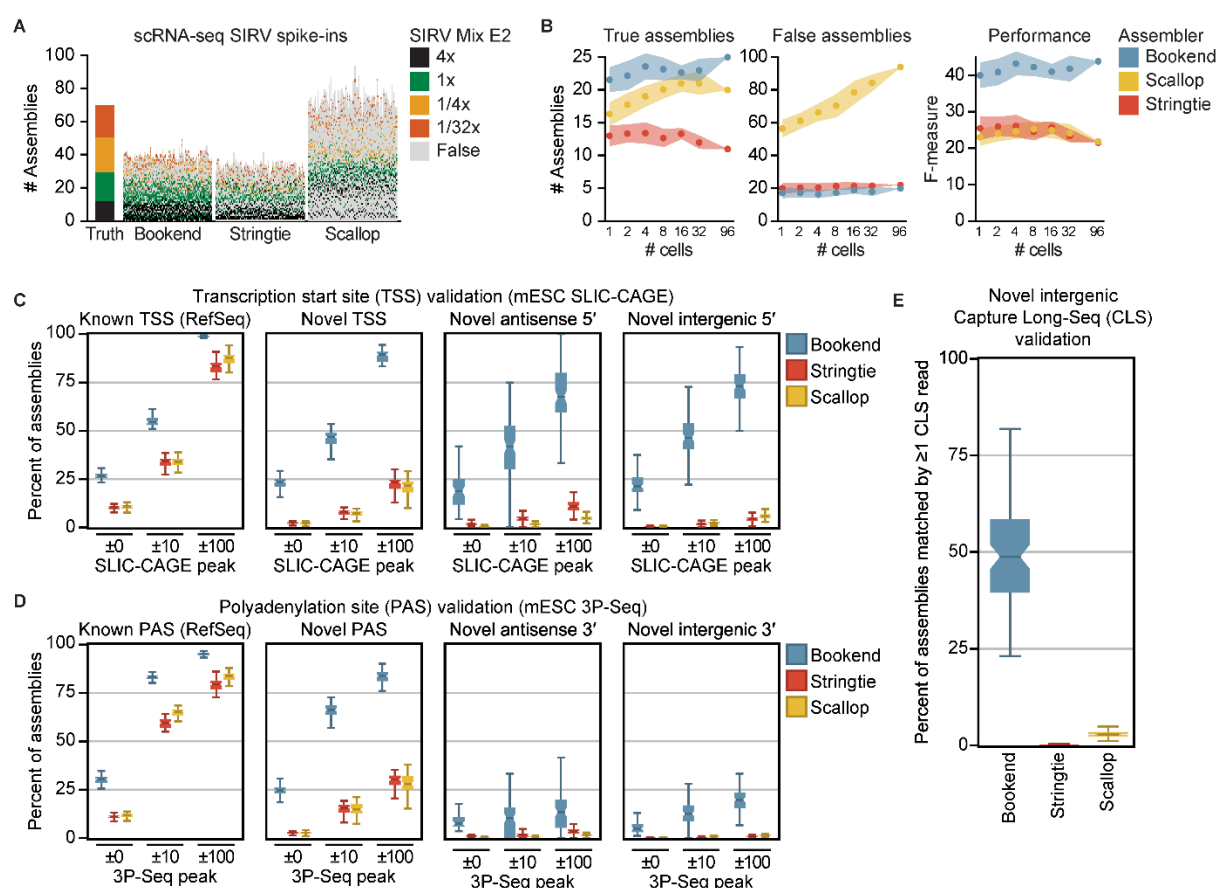
291   involving two microRNAs, miR162 and miR838, the latter of which is encoded in intron 14 of its

292   own gene[46,47]. Long reads alone were not sufficient to define the canonical 6.2 kilobase *DCL1*

293   transcript because 7 of 8 PacBio reads mapping to *DCL1* were non-capped truncations, and intron

294   14 was retained in the only full-length read (Fig3F). By synthesizing information from multiple

295   modes of sequencing, hybrid assembly with Bookend built a more complete transcript catalog

296   that includes both the fully-spliced isoform and the isoform that retains MIR838. As a final

297   refinement, a hybrid assembly that requires the presence of Cap Tags at transcript 5′ ends yielded

298   a transcriptome with a 74.6% global concordance with the TAIR10 annotation. We report this

299   hybrid assembly of long, short and 5′ end reads as the Bookend Floral Bud annotation

300   (Supplemental Dataset 1-2).

301

302   *Transcript discovery from single-cell sequencing*

303   Bookend achieved comparable precision assembling Arabidopsis transcriptomes from either long

304   reads or short reads generated by Smart-seq2, which is a protocol routinely used for single-cell

305   RNA sequencing (scRNA-seq) (FigS3C). However, scRNA-seq poses multiple hurdles to

306   accurate assembly. Amplifying the few picograms of RNA in a single cell exacerbates biases and

307   artifacts during reverse transcription[22], and dropouts from inefficient RNA capture place limits on

308   accurate isoform quantification from scRNA-seq[48]. Additionally, scRNA-seq has been most widely

309   adopted in the study of mammalian systems. The mouse genome (and likewise the human

310   genome) is roughly 30 times larger than the Arabidopsis genome with an average of twice as

311   many introns per gene and nearly three times the number of annotated isoforms. Additionally,

312   mouse introns can exceed 100kb and are on average 36 times longer than in Arabidopsis. Many

313   isoforms per gene and large spans of non-genic sequence make it considerably more challenging

314   both to assemble transcripts and to validate which assemblies are correct. To evaluate Bookend's

315   utility on mammalian scRNA-seq data, we tested it on a dataset designed for single-cell

316   benchmarking[49] which contains a set of synthetic Spike-In RNA Variants (SIRVs) added prior to

14

317   cell lysis. SIRVs were designed to present a challenge to isoform quantification tools by mimicking

318   complex mammalian genes[50]. The 69 synthetic transcripts map to 7 regions on a hypothetical

319   genome in a way that recapitulates canonical and non-canonical splicing variation, antisense

320   transcription and alternative 5′ and 3′ ends with up to 18 isoforms per gene (Fig S4A). SIRV Mix

321   E2 contains molecules in four discrete concentrations so that each locus has major and minor

322   isoforms that vary in relative abundance by up to 128-fold. SMARTer library preparations from 96

323   single mouse embryonic stem cells (mESCs) were deeply sequenced, with an average of 7 million

324   aligned paired-end 100bp reads per cell (Supplemental Table 4) including an average of just over

325   500,000 SIRV-mapping reads per cell. Bookend correctly reconstructed (full splice match and

326   ≤100nt error on both ends) an average of 22.6 transcripts per cell, which was higher than either

327   Scallop (16.3) or StringTie2 (13) (Fig5AB). Moreover, Bookend assembled fewer false SIRVs than

328   StringTie2 and especially Scallop (Fig5B). To test a relationship between performance and

329   sequencing depth, cells were progressively combined into pairs, then sets of 4, 16, 32, and a full

330   merge of reads from all 96 cells. The relative performance of the three assemblers was stable

331   over two orders of magnitude of input with the F-measure (harmonic mean of precision and recall)

332   slightly rising for Bookend as the sequencing depth increased and slightly decreasing for the

333   others (Fig5B). Importantly, Bookend consistently assigned a higher estimated abundance to true

334   transcripts, and false assemblies were more concentrated in the low abundance regime than for

335   other assemblers (Fig5A). Overall precision on SIRVs averaged 55.9% for Bookend (vs. 39.6%

336   StringTie2, 22.5% Scallop), and precision on the most abundant half of assemblies was 74.2%

337   (vs. 48.2% StringTie2, 28.4% Scallop).

338       End-labeled reads mapping to the mouse genome were also assembled for each cell, and

339   transcript models were compared to RefSeq mm39. All matching exon chains were considered

340   matches, and precision was measured as the percent of all assemblies that match RefSeq. Recall

341   was defined by tallying all transcripts correctly assembled at least once and counting the

342   proportion of this transcript set found per cell. Although recall was considerably worse for Bookend

15

343    (average 7.9%) than other methods (StringTie2 16.6%, Scallop 16.5%), precision was multiple

344    times higher (76.3% Bookend, 29.0% StringTie2, 26.5% Scallop; FigS4B). Assemblies were

345    repeated for two replicates of Smart-seq2 data from the same experiment with comparable results

346    demonstrating that end-guided assembly is consistent across full-length sequencing protocols
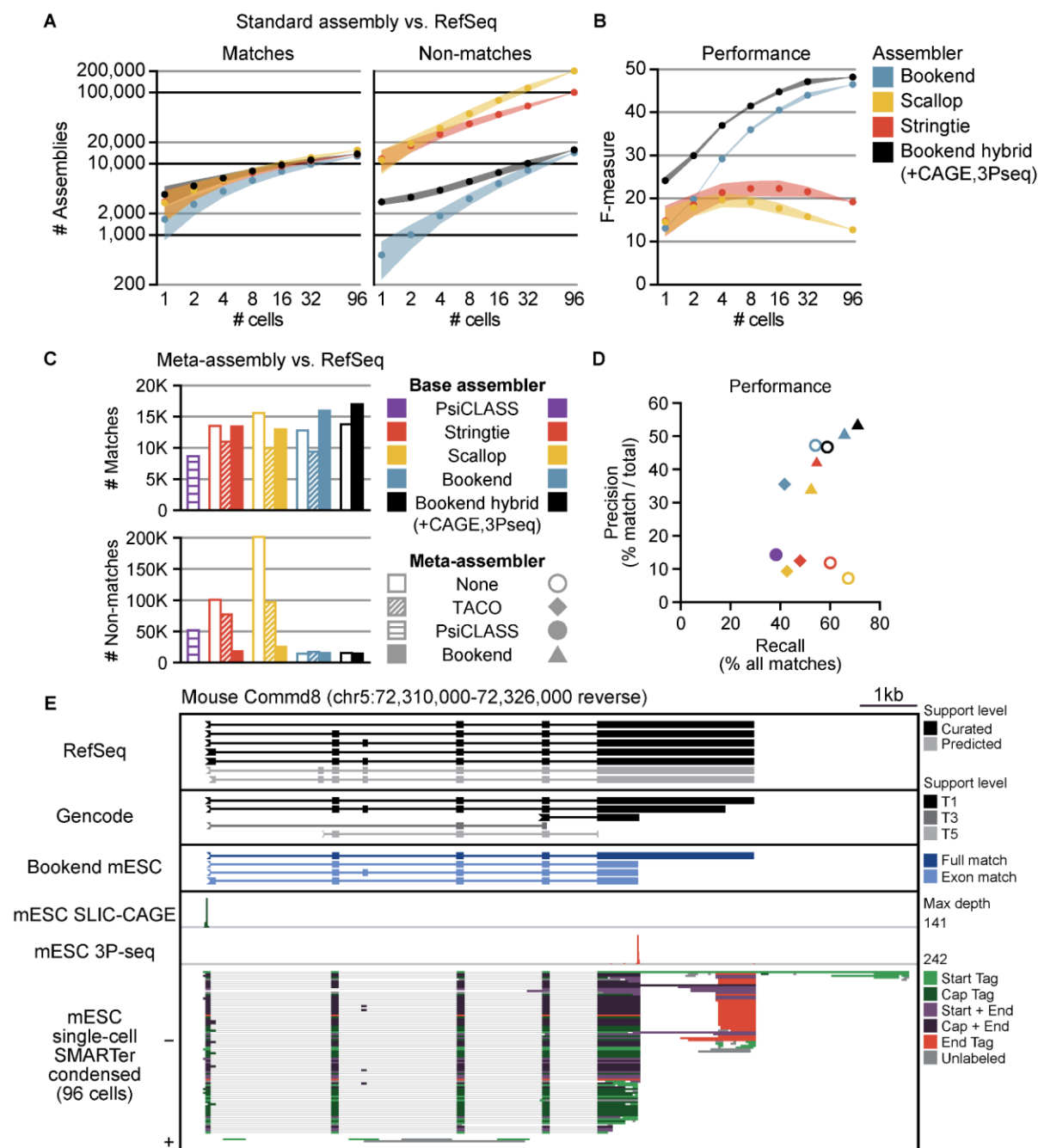
347    (FigS4B).

348



349

350    **Figure 4.** Bookend performance on single mouse cells
351    **(A)** Reconstruction of Spike-In RNA Variants (SIRVs) from 96 paired-end 100bp SMARTer libraries
352    of single mESCs. Each vertical bar depicts the assemblies from one cell, ordered from highest
353    (bottom) to lowest (top) estimated abundance. Colored boxes match a true isoform of the given
354    input concentration; gray boxes are false assemblies. **(B)** SIRV assembly performance as a
355    function of increasing sequencing depth. F1 score (right) is the harmonic mean of sensitivity and
356    precision. **(C)** Boxplots showing percent validation of 5′ ends with SLIC-CAGE support within the
357    given windows for 96 single mESC assemblies. **(D)** Boxplots as in **C** showing 3′ end validation by
358    3P-Seq peaks. **(E)** Percent of intergenic assemblies (no overlap with RefSeq) in single cells which
359    have ≥1 matching Capture Long-Seq read from the mouse CLS atlas.

360

16

361    As with TAIR10, RefSeq is almost certainly incomplete, and non-reference-matching

362    assemblies could still be valid. To experimentally validate non-RefSeq mESC assemblies, three

363    validation datasets were used: uuG-containing SLIC-CAGE[17] reads from mESCs for 5' end

364    validation, mESC 3P-Seq[51] reads for 3' end validation, and a database of long noncoding RNAs

365    identified by intergenic Capture Long-read Sequencing (CLS[52]) for full-length validation of novel

366    intergenic loci. An assembly was considered validated by a method if at least one read directly

367    supported an assembled transcript's respective structure(s). Assemblies with 5' ends ≤100nt

368    away from a RefSeq TSS contained "known" TSSs, and all others possessed "novel" TSSs.

369    Likewise, assemblies with 3' ends ≤100nt from their matching reference polyadenylation sites

370    were considered "known" PASs and all others were "novel". An average of 99.7% of Bookend,

371    83.9% of Scallop and 79.0% of Stringtie2 single-cell assemblies with a known TSS had at least

372    one SLIC-CAGE read within 100nt (Fig4C). Moreover, the majority of novel, antisense and

373    intergenic TSSs from Bookend transcripts were supported by at least 1 capped SLIC-CAGE read,

374    whereas no novel group from StringTie2 or Scallop surpassed a 25% validation rate. The 3P-Seq

375    dataset had fewer total reads and was less sensitive overall, but it still supported 19.9% of

376    intergenic Bookend assembly 3' ends, compared to 1.4% for Scallop and 0.8% for StringTie2

377    (Fig4D). By comparing against the CLS atlas we could validate the full structure of intergenic

378    mESC assemblies. Bookend assembled a very small number of novel intergenic transcripts per

379    cell (average 33 vs. 1209 by StringTie2 and 1073 by Scallop), but 49% of these were supported

380    by one or more reads from the CLS atlas, compared to just 3% for Scallop intergenic assemblies

381    and 0.3% for StringTie2 (Fig4E). Finally, because Cap and End Tags were extremely sparse in

382    each cell (Supplemental Table 4), we hypothesized that the lower sensitivity could be explained

383    by dropout of end labels. Supplying the mESC SLIC-CAGE (5') and 3P-seq (3') datasets to a

384    Bookend hybrid assembly raised recall from 7.9% to 18.2% and retained a precision of 67.2%

385    (FigS4B). Therefore, end-guided assembly of single-cell RNA-seq data can be used to identify

386    genuine transcriptional novelty that is otherwise masked by noise.

387    *Condensed assembly and meta-assembly*

388    A defining feature of single-cell experiments is that many individual cells are profiled in parallel.

389    While sensitivity in an individual cell is low, information across multiple cells can be combined to

390    achieve a more complete view of the experiment. Tools have been developed for transcript "meta-

391    assembly" of reads from multiple sources. By modeling for variation across samples, meta-

392    assemblers achieve higher precision than standard assembly on the same set of reads[53,54]. To

393    measure the impact of meta-assembly, a series of assemblies on subsamples of all 706 million

394    aligned single-cell mESC reads was first performed with StringTie2 and Scallop, as well as

395    Bookend with and without the addition of mESC SLIC-CAGE and 3P-seq libraries (Fig5A). The

396    mean number of reference-matching transcripts varied greatly across assemblers on single cells

397    (1,656 Bookend, 3,711 Bookend hybrid, 2,904 StringTie2, 2,831 Scallop), but the magnitude of

398    difference decreased with progressive doublings, up to the full set of 96 cells (12,794 Bookend,

399    13,762 Bookend hybrid, 13,524 StringTie2, 15,611 Scallop). By contrast, non-matches grew

400    linearly with input. Bookend consistently assembled roughly an order of magnitude fewer non-

401    matching transcripts than other assemblers across all input levels. From the full 96-cell dataset

402    Scallop identified the most matches, but this was dwarfed by nearly 13 times the number of

403    assemblies that failed to match RefSeq (201,631 Scallop, 100,646 StringTie2, 14,301 Bookend,

404    15,711 Bookend hybrid). By assuming non-matches to be mostly false, we calculated recall and

405    precision as before and combined them to track the relationship between overall performance (F-

406    measure) and input. F-measure of Bookend and Bookend hybrid assembly continued to improve

407    with increasing input, but Scallop and StringTie2 began to decline above 4 and 16 cells,

408    respectively, due to the growth of non-matches outpacing matches (Fig 5B). Consistent with

409    previous reports, we see that standard assemblers suffer from an input-dependent decay in

410    precision[53,54].

**Figure 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 ≥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 (top) or the number of non-matches (bottom). **(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.

19

423    As an alternative approach, two published meta-assemblers were used to process the 96-cell

424    dataset. TACO builds a consensus annotation by re-defining transcript boundaries through

425    "change-point detection" on a set of files from any standard assembler[53], whereas PsiCLASS

426    generates the individual assemblies and performs meta-assembly through a consensus voting

427    system[54]. The flexibility of Bookend's framework allows its assembly algorithm to be run on

428    assemblies, including its own output. To test the efficacy of meta-assembly with Bookend, each

429    of the 96 mESC cell datasets were "condensed" by a first pass through Bookend Assemble in

430    which no incomplete transcripts were discarded (FigS5A; see Supporting Notes: "Path Filtering").

431    Assembly was run again on the 96 condensed files, only retaining complete transcript models

432    during the second pass. Bookend was also used to meta-assemble the 96 single-cell assemblies

433    by StringTie2 and Scallop. Compared to standard assembly by StringTie2 or Scallop, all meta-

434    assemblies produced substantially fewer non-matching transcripts (Fig5C). However, single-cell

435    meta-assemblies surprisingly also recalled fewer RefSeq matches than standard assembly, with

436    the exception of Bookend-to-Bookend and hybrid Bookend-to-Bookend meta-assemblies.

437    PsiCLASS and TACO both showed somewhat higher precision than standard assembly, but at

438    the expense of a severe drop in recall (Fig5D). PsiCLASS had the lowest recall of any method,

439    but higher precision than StringTie2-to-TACO or Scallop-to-TACO meta-assembly. Bookend-to-

440    Bookend meta-assemby considerably outperformed PsiCLASS in both recall (relative increase of

441    72%) and precision (relative increase of 253%). PsiCLASS produced an unusually large number

442    of partial transcript fragments, likely due to the fact that scRNA-seq often has substantial 3′ bias

443    that is not adequately accounted for (FigS5A-B). Notably, when TACO was applied to single-cell

444    Bookend assemblies, it showed both a 23% relative reduction in recall and a 25% relative

445    reduction in precision compared to standard Bookend assembly. In contrast, Bookend-to-

446    Bookend meta-assembly increased recall by 22% and precision by 7% (+58% recall and +42%

447    precision vs. Bookend-to-TACO). Across all three base assemblers, TACO reported fewer full

448    reference matches than the standard assembly, while Bookend reported the same number or

449 more full matches with a greater reduction in all non-matching classes than TACO (FigS5C). Of

450 all combinations tested, both sensitivity and precision were highest at the intron chain and full

451 transcript level in a Bookend-to-Bookend hybrid meta-assembly in which SLIC-CAGE and 3P-seq

452 data were supplied alongside the single-cell condensed assemblies[55] (Supplemental Table S5).

453 We report this assembly as the "Bookend mESC" annotation (Supplemental Dataset 3-4).

454 Requiring that both transcript ends are replicable across at least two different samples raised the

455 transcript-level concordance with RefSeq to 54.1%, a relative increase of 271% over the most

456 precise non-Bookend method (PsiCLASS), and a substantially higher agreement than even

457 Gencode, an alternative mouse reference annotation that only shares 31.7% of its transcripts at

458 assembled loci with RefSeq (FigS5C). While Gencode isoforms contain a broader set of

459 alternative TSS and PAS than RefSeq, we noticed that they can be contained in low-confidence

460 or fragmented transcript models, as in the gene Commd8 (Fig5E). By combining multiple unique

461 advantages of end-guided assembly, Bookend could assemble more reference matches than any

462 other strategy while maintaining a majority concordance with known annotations.

463

464 **DISCUSSION**

465 Computational gene annotation pipelines have long struggled to produce a reliable picture of plant

466 and animal transcriptomes at the isoform level[11,29,56]. Studying the details of gene regulation and

467 isoform usage remains restricted to a small number of model organisms in which manually

468 curated accurate transcript models are available. Even with specialized methods for sequencing

469 RNA ends, connecting those ends to a gene model can be computationally challenging, especially

470 for noncoding RNAs[35]. By generating accurate end-to-end transcript assemblies from a range of

471 widely accessible sequencing methods, Bookend enables the automated annotation of promoter

472 architecture, alternative polyadenylation and splicing dynamics in tissues in response to

473 developmental, environmental and disease state cues.

474     Despite rapid advancements in scale and sensitivity of single-cell RNA sequencing, the

475     accurate detection of transcript isoforms is still an outstanding challenge[48]. Full-length cDNA can

476     be amplified from single cells with the Smart-seq family of "full-length sequencing" methods,

477     including the recently developed Smart-seq3 that more efficiently captures 5'-labeled ends and

478     gene body reads simultaneously[22,23]. Multiple approaches to apply long-read sequencing to single

479     cells have been developed, but limits on throughput, error rate, and cost restricts their use[57–59].

480     Notably, large-scale Smart-seq2 experiments across multiple organisms have already been

481     sequenced, including tens of thousands of cells from 20 mouse tissues and 24 human tissues by

482     the Tabula Muris and Tabula Sapiens Consortia, respectively[60,61]. Through meta-assembly of full-

483     length scRNA-seq data, Bookend enables the wholesale reannotation of genomes at single-cell

484     resolution using existing and future datasets.

485 **METHODS**

486 *PacBio Sequencing*

487 Two PacBio Iso-seq libraries were generated each using 10 μg of total RNA from Arabidopsis

488 inflorescences containing unopened floral buds. Total RNA was extracted with TRIzol following

489 the method described in Schon et al. 2018[16] to yield two biological replicates with an RNA integrity

490 number (RIN) of 9.0 and 9.2, respectively. SMRTbell libraries were constructed by the Vienna

491 BioCenter Core Facilities (VBCF) and sequenced on a Sequel SMRT Cell 1M.

492

493 *Published RNA sequencing data*

494 Smart-seq2 datasets from 5ng Arabidopsis thaliana floral bud RNA and tissue-matched

495 nanoPARE libraries from 10ug total RNA were downloaded from the NCBI Gene Expression

496 Omnibus (GEO), series accession GSE112869. Single-cell RNA-seq of mouse embryonic stem

497 cells and SIRVs from Natarajan et al. 2019 was downloaded from EMBL-EBI ArrayExpress,

498 accession E-MTAB-7239. SLIC-CAGE samples from 100ng mESC total RNA were downloaded

499 from ArrayExpress, accession E-MTAB-6519. One 3P-Seq library from 75ug mESC RNA was

500 downloaded from GEO, sample accession GSM1268958.

501

502 *Short read data processing*

503 Prior to alignment, reads were preprocessed with cutadapt[62] to remove sequencing adapters. End

504 labels were identified and trimmed using the utility *bookend label*, with settings tailored to each

505 library. For Arabidopsis single-end Smart-seq2 reads, the arguments *--strand unstranded -S*

506 *AAGCAGTGGTATCAACGCAGAGTACGGG                                                                  -E*

507 *AAGCAGTGGTATCAACGCAGAGTACTTTTTTTTTTTTTTTTTTTTTT+ --min_start 7 --min_end 9 -*

508 *-minlen 18 --minqual 25 --qualmask 16 --mismatch_rate 0.06* were used. Paired-end mouse

509 SMARTer     reads     used     the     same     arguments     except     for     *-S*

510 *AAGCAGTGGTATCAACGCAGAGTACATGGG*. 5' end reads from nanoPARE libraries were

511    labeled with the arguments *--strand forward --minstart 20*. After end labeling, short reads were

512    aligned using STAR[63]. Arabidopsis reads were aligned to the TAIR10 genome, and mouse reads

513    were aligned to mm39 (GRCm39). Short reads in both species were aligned using an identical

514    two-pass alignment strategy except for allowed intron lengths. First, reads were aligned with the

515    command *STAR --runMode alignReads --alignEndsType EndToEnd --outFilterMatchNmin 20 --*

516    *outFilterMismatchNmax    6    --outFilterMismatchNoverLmax    .05    --outFilterIntronMotifs*

517    *RemoveNoncanonicalUnannotated --alignSJoverhangMin 20 --alignSJDBoverhangMin 1 --*

518    *outFilterMultimapNmax 2 --outSJfilterOverhangMin -1 15 20 20 --outSJfilterCountUniqueMin -1 2*

519    *3 3 --outSJfilterCountTotalMin -1 2 3 3.* Arabidopsis alignments used the additional arguments *--*

520    *alignIntronMax 5000 --alignMatesGapMax 5100*, and mouse alignments instead used *--*

521    *alignIntronMax 100000 --alignMatesGapMax 100100*. Splice junctions from all samples were

522    aggregated across all samples for each species with *bookend sj-merge --new --min_reps 2* to

523    retain only novel splice junctions that were detected in multiple samples. Second pass mapping

524    was performed with the settings above, except the merged splice junction file was provided with

525    *--sjdbFileChrStartEnd*, and the following arguments were modified: *--alignEndsType Local --*

526    *outFilterMatchNminOverLread 0.9 --outFilterType BySJout --outFilterMultimapNmax 10 --*

527    *outSAMtype BAM Unsorted --outSAMorder Paired --outSAMprimaryFlag AllBestScore --*

528    *outSAMattributes NH HI AS nM NM MD jM jI XS*. Unsorted BAM files were converted to End-

529    Labeled Read (ELR) files with the command *bookend elr --genome [genome.fa]* with library-

530    specific    settings.    Arabidopsis    Smart-seq2:    *--start_seq    ACGGG    --end_seq*

531    *RRRRRRRRRRRRRRRRRRRRRRRRRRRRRR --mismatch_rate .2;* Arabidopsis nanoPARE: --

532    stranded *-s --start_seq ACGGG --mismatch_rate .2*; mouse SMARTer: *--start_seq ACATGGG --*

533    *end_seq AAAAARRRRRRRRRRRRRRRRRRRRRRRRRRR --mismatch_rate .25.*

534

535 *Long read data processing*

536 Raw Arabidopsis PacBio reads were converted to Circular Consensus Sequences using Iso-seq3

537 software with the command *ccs --min-passes 2 --min-rq .9*, and CCS reads were converted to

538 full-length non-chimeric (FLNC) reads using *lima* and *isoseq3 refine --require-polya --min-rq -1 --*

539 *min-polya-length 10*. FLNC reads were aligned to the Arabidopsis genome with the command

540 *minimap2 -G 5000 -H -ax splice --MD -C 5 -u f -p 0.9 --junc-bed [TAIR10 transcript BED12]*.

541 Aligned unsorted SAM files were converted to ELR with the command *bookend elr --stranded -s*

542 *-e --start_seq ATGGG --genome [TAIR10.fa]*.

543

544 *Assembly*

545 To make assembly setting maximally uniform across Bookend, StringTie2, Scallop, and Cufflinks,

546 the following arguments were used. For Arabidopsis assemblies: *bookend --max_gap 50 --*

547 *min_cov 2 --min_len 60 --min_proportion 0.02 --min_overhang 3 --cap_bonus 5 --cap_filter 0.02*;

548 *stringtie -g 50 -c 2 -m 60 -f 0.02 -a 3 -M 1 -s 5; scallop --min_bundle_gap 50 --*

549 *min_transcript_coverage 2 --min_transcript_length_base 60 --min_flank_length 3 --*

550 *min_single_exon_coverage 5 --min_transcript_length increase 50; cufflinks -F 0.02 --overhang-*

551 *tolerance 3 --min-frags-per-transfrag 10 -j 0.15 -A 0.06*. For mouse assemblies the same settings

552 were used with the following exceptions: *--min_proportion* was set to 0.01, --min_len to 200, and

553 *--require_cap* was enforced on mouse assemblies except when assembling spike-in transcripts,

554 which do not possess caps. For meta-assembly, Bookend was run with the same settings as

555 above for mouse. TACO was run with the arguments *--filter-min-expr 2 --filter-min-length 200 --*

556 *isoform-frac 0.01*, and PsiCLASS was run with default settings

557

558 *Assembly algorithms*

559 A brief overview of the end-guided assembly process implemented in Bookend is below. For a full

560 breakdown of the algorithms used, see the "Bookend Algorithms" Supplemental Note.

25

561 (*Generate Chunks*) First, reads are streamed in from an ELR file in sorted order and separated

562 into overlapping chunks. (*Tag Clustering*) In each chunk, Start Tags and End Tags are clustered

563 on each strand by grouping tags by genomic position and assigning each position a signal score

564 of counts $\times$ proportion of total coverage. A signal threshold is set and positions below the

565 threshold are discarded. Remaining positions are grouped within a user-specified distance to yield

566 Start and End clusters on each strand. (*Calculate Membership Matrix*) Start/End clusters are

567 added to a catalog of boundaries, which include splice donor/acceptor sites that are also filtered

568 by a threshold of total overlapping coverage. Adjacent boundary pairs define a "frag", and each

569 read is assigned a Membership array that describes whether the read overlaps or excludes each

570 frag. Redundant membership arrays are combined, and the unique set of elements is stored as

571 the Membership Matrix. (*Calculate Overlap Matrix*) A matrix describing the relationship between

572 each element pair *a* and *b* is generated by asking (from left to right in genomic coordinates): can

573 *a* extend into *b*? Can *b* extend into *a*? Each comparison returns a pair of Overlaps, $O_{ab}$ and $O_{ba}$,

574 respectively: 1 = extends, -1 = excludes, 2 = is contained by, 0 = does not overlap. The values -

575 1 and 0 are symmetric, but 1 and 2 are directed relationships that can be used as edges in a

576 directed graph. (*Collapse Linear Chains*) It is possible to identify and collapse non-branching sets

577 of elements ("linear chains") prior to assembly. Two graphs are constructed with elements as

578 nodes: a directed graph with extensions as edges, and an undirected graph with exclusions as

579 edges. A depth-first search is conducted by visiting each element in increasing order of

580 information content (number of non-zero memberships). During a visit, the element's edges are

581 traversed recursively to record all traversed nodes' exclusions. An element with no edges is

582 assigned to a new chain. Otherwise, when an element's edges are all traversed, the element is

583 compared against its outgroup, the set of all elements reached. If all outgroup elements belong

584 to one chain and the element and outgroup have the same set of exclusions, then the element is

585 added to the same chain. If the element's outgroup is assigned to multiple chains, the element

586 begins a new chain. After completion of the search, each chain is combined to form a single

587　reduced element. (*Generate Overlap Graph*) From the set of reduced elements a second directed

588　graph is constructed with a global source (Start+/End-) and sink (Start-/End+), where each node

589　records the element weight (sequenced bases / genomic length), outgroup (extends to), ingroup

590　(extends from), containments and exclusions. (*Resolve Containment*) All elements contained by

591　one or more longer elements have their weight redistributed proportionally to their container as

592　long as not all containers exclude any single node the element doesn't already exclude. (*Greedy*

593　*Paths*) All elements begin unassigned. Starting with the heaviest unassigned element, choose an

594　extension (ingroup/outgroup pair) that maximizes a score that equally combines the following:

595　maximal weight of the extension, maximal similarity of weight distribution across samples between

596　element and extension, minimal coverage variance across covered frags, and does not cause the

597　source or sink to become unreachable. The highest-scoring extension is iteratively added to a

598　path until both source and sink are reached. Paths are generated in this manner until the total

599　weight of unassigned elements falls below a given signal threshold.

600

601　**Contributions**

602　M.A.S. and M.D.N. conceived the project; M.A.S. developed the methodology; M.A.S and S.L.

603　performed the experiments; M.A.S. and F.H. analyzed data; M.A.S. prepared figures; M.A.S wrote

604　the article; M.A.S. and M.D.N. edited the article; M.D.N. acquired funding and supervised the

605　project.

606

607　**Acknowledgements**

608　We thank the Next Generation Sequencing Facility at Vienna BioCenter Core Facilities GmbH

609　(VBCF) for their outstanding services and technical support.

610

611　**Competing interests**

612　The authors declare that they have no conflicts of interests.

**Data access**

Bookend software is available on the Python Package Index and can be installed with the command *pip install bookend-rna*. Source code is available as a repository on GitHub at https://github.com/Gregor-Mendel-Institute/bookend. 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.

## REFERENCES

1. Liu, F., Marquardt, S., Lister, C., Swiezewski, S. & Dean, C. Targeted 3' processing of antisense transcripts triggers Arabidopsis FLC chromatin silencing. *Science* **327**, 94–97 (2010).

2. Rhinn, H. *et al.* Alternative α-synuclein transcript usage as a convergent mechanism in Parkinson's disease pathology. *Nat. Commun.* **3**, 1084 (2012).

3. Solana, J. *et al.* Conserved functional antagonism of CELF and MBNL proteins controls stem cell-specific alternative splicing in planarians. *Elife* **5**, (2016).

4. Mudge, J. M. & Harrow, J. The state of play in higher eukaryote gene annotation. *Nat. Rev. Genet.* **17**, 758–772 (2016).

5. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res.* **47**, D766–D773 (2019).

6. McGarvey, K. M. *et al.* Mouse genome annotation by the RefSeq project. *Mamm. Genome* **26**, 379–390 (2015).

7. Berardini, T. Z. *et al.* The Arabidopsis information resource: Making and mining the 'gold standard' annotated reference plant genome. *Genesis* **53**, 474–485 (2015).

8. FANTOM Consortium and the RIKEN PMI and CLST (DGT) *et al.* A promoter-level mammalian expression atlas. *Nature* **507**, 462–470 (2014).

9. Wu, P.-Y., Phan, J. H. & Wang, M. D. Assessing the impact of human genome annotation choice on RNA-seq expression estimates. *BMC Bioinformatics* **14 Suppl 11**, S8 (2013).

10. SEQC/MAQC-III Consortium. A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium. *Nat. Biotechnol.* **32**, 903–914 (2014).

11. Guigó, R. *et al.* EGASP: the human ENCODE Genome Annotation Assessment Project. *Genome Biol.* **7 Suppl 1**, S2.1–31 (2006).

12. Stark, R., Grzelak, M. & Hadfield, J. RNA sequencing: the teenage years. *Nat. Rev. Genet.* **20**, 631–656 (2019).

13. Levin, J. Z. *et al.* Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nat. Methods* **7**, 709–715 (2010).

14. Murata, M. *et al.* Detecting expressed genes using CAGE. *Methods Mol. Biol.* **1164**, 67–85 (2014).

15. Adiconis, X. *et al.* Comprehensive comparative analysis of 5'-end RNA-sequencing methods. *Nat. Methods* **15**, 505–511 (2018).

16. Schon, M. A., Kellner, M. J. & Plotnikova, A. NanoPARE: parallel analysis of RNA 5′ ends from low-input RNA. *Genome Res.* (2018).

17. Cvetesic, N. *et al.* SLIC-CAGE: high-resolution transcription start site mapping using nanogram-levels of total RNA. *Genome Res.* **28**, 1943–1956 (2018).

18. Jan, C. H., Friedman, R. C., Ruby, J. G. & Bartel, D. P. Formation, regulation and evolution of Caenorhabditis elegans 3'UTRs. *Nature* **469**, 97–101 (2011).

19. Moll, P., Ante, M., Seitz, A. & Reda, T. QuantSeq 3′ mRNA sequencing for RNA quantification. *Nat. Methods* **11**, i–iii (2014).

20. Pelechano, V., Wei, W. & Steinmetz, L. M. Extensive transcriptional heterogeneity revealed by isoform profiling. *Nature* **497**, 127–131 (2013).

21. Wang, J. *et al.* TIF-Seq2 disentangles overlapping isoforms in complex human transcriptomes. *Nucleic Acids Res.* **48**, e104 (2020).

22. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat. Methods* **10**, 1096–1098 (2013).

23. Hagemann-Jensen, M. *et al.* Single-cell RNA counting at allele and isoform resolution using Smart-seq3. *Nat. Biotechnol.* **38**, 708–714 (2020).

24. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian organogenesis.

675      *Nature* **566**, 496–502 (2019).

676  25. Cao, J. *et al.* Comprehensive single-cell transcriptional profiling of a multicellular organism.
677      *Science* **357**, 661–667 (2017).

678  26. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat.*
679      *Commun.* **8**, 14049 (2017).

680  27. Garalde, D. R. *et al.* Highly parallel direct RNA sequencing on an array of nanopores. *Nat.*
681      *Methods* **15**, 201–206 (2018).

682  28. Wan, Y. *et al.* Systematic identification of intergenic long-noncoding RNAs in mouse retinas
683      using full-length isoform sequencing. *BMC Genomics* **20**, 559 (2019).

684  29. Steijger, T. *et al.* Assessment of transcript reconstruction methods for RNA-seq. *Nat.*
685      *Methods* **10**, 1177–1184 (2013).

686  30. Tardaguila, M. *et al.* SQANTI: extensive characterization of long-read transcript sequences
687      for quality control in full-length transcriptome identification and quantification. *Genome Res.*
688      (2018) doi:10.1101/gr.222976.117.

689  31. Kuo, R. I. *et al.* Normalized long read RNA sequencing in chicken reveals transcriptome
690      complexity similar to human. *BMC Genomics* **18**, 323 (2017).

691  32. Tang, A. D. *et al.* Full-length transcript characterization of SF3B1 mutation in chronic
692      lymphocytic leukemia reveals downregulation of retained introns. *Nat. Commun.* **11**, 1438
693      (2020).

694  33. Kovaka, S. *et al.* Transcriptome assembly from long-read RNA-seq alignments with
695      StringTie2. *Genome Biol.* **20**, 278 (2019).

696  34. Cumbie, J. S., Ivanchenko, M. G. & Megraw, M. NanoCAGE-XL and CapFilter: an
697      approach to genome wide identification of high confidence transcription start sites. *BMC*
698      *Genomics* **16**, 597 (2015).

699  35. de Rie, D. *et al.* An integrated expression atlas of miRNAs and their promoters in human
700      and mouse. *Nat. Biotechnol.* **35**, 872–878 (2017).

701  36. Thieffry, A. *et al.* Characterization of Arabidopsis thaliana Promoter Bidirectionality and
702      Antisense RNAs by Inactivation of Nuclear RNA Decay Pathways. *Plant Cell* **32**, 1845–
703      1867 (2020).

704  37. Sherstnev, A. *et al.* Direct sequencing of Arabidopsis thaliana RNA reveals patterns of
705      cleavage and polyadenylation. *Nat. Struct. Mol. Biol.* **19**, 845–852 (2012).

706  38. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-
707      seq reads. *Nat. Biotechnol.* **33**, 290–295 (2015).

708  39. Shao, M. & Kingsford, C. Accurate assembly of transcripts through phase-preserving graph
709      decomposition. *Nat. Biotechnol.* **35**, 1167–1169 (2017).

710  40. Trapnell, C. *et al.* Differential gene and transcript expression analysis of RNA-seq
711      experiments with TopHat and Cufflinks. *Nat. Protoc.* **7**, 562–578 (2012).

712  41. Amarasinghe, S. L. *et al.* Opportunities and challenges in long-read sequencing data
713      analysis. *Genome Biol.* **21**, 30 (2020).

714  42. Goodwin, S., McPherson, J. D. & McCombie, W. R. Coming of age: ten years of next-
715      generation sequencing technologies. *Nat. Rev. Genet.* **17**, 333–351 (2016).

716  43. Balázs, Z. *et al.* Template-switching artifacts resemble alternative polyadenylation. *BMC*
717      *Genomics* **20**, 824 (2019).

718  44. Tang, D. T. P. *et al.* Suppression of artifacts and barcode bias in high-throughput
719      transcriptome analyses utilizing template switching. *Nucleic Acids Res.* **41**, e44 (2013).

720  45. Gordon, S. P. *et al.* Widespread Polycistronic Transcripts in Fungi Revealed by Single-
721      Molecule mRNA Sequencing. *PLoS One* **10**, e0132628 (2015).

722  46. Xie, Z., Kasschau, K. D. & Carrington, J. C. Negative Feedback Regulation of Dicer-Like1
723      in Arabidopsis by microRNA-Guided mRNA Degradation. *Current Biology* vol. 13 784–789
724      (2003).

725  47. Rajagopalan, R., Vaucheret, H., Trejo, J. & Bartel, D. P. A diverse and evolutionarily fluid

set of microRNAs in Arabidopsis thaliana. *Genes Dev.* **20**, 3407–3425 (2006).

48. Westoby, J., Artemov, P., Hemberg, M. & Ferguson-Smith, A. Obstacles to detecting isoforms using full-length scRNA-seq data. *Genome Biol.* **21**, 74 (2020).

49. Natarajan, K. N. *et al.* Comparative analysis of sequencing technologies for single-cell transcriptomics. *Genome Biol.* **20**, 70 (2019).

50. Paul, L. *et al.* SIRVs: Spike-In RNA Variants as External Isoform Controls in RNA-Sequencing. *bioRxiv* 080747 (2016) doi:10.1101/080747.

51. Nam, J.-W. *et al.* Global analyses of the effect of different cellular contexts on microRNA targeting. *Mol. Cell* **53**, 1031–1043 (2014).

52. Lagarde, J. *et al.* High-throughput annotation of full-length long noncoding RNAs with capture long-read sequencing. *Nat. Genet.* **49**, 1731–1740 (2017).

53. Niknafs, Y. S., Pandian, B., Iyer, H. K., Chinnaiyan, A. M. & Iyer, M. K. TACO produces robust multisample transcriptome assemblies from RNA-seq. *Nat. Methods* **14**, 68–70 (2017).

54. Song, L., Sabunciyan, S., Yang, G. & Florea, L. A multi-sample approach increases the accuracy of transcript assembly. *Nat. Commun.* **10**, 5000 (2019).

55. Pertea, G. & Pertea, M. GFF Utilities: GffRead and GffCompare. *F1000Res.* **9**, 304 (2020).

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

57. Gupta, I. *et al.* Single-cell isoform RNA sequencing characterizes isoforms in thousands of cerebellar cells. *Nat. Biotechnol.* (2018) doi:10.1038/nbt.4259.

58. Philpott, M. *et al.* Nanopore sequencing of single-cell transcriptomes with scCOLOR-seq. *Nat. Biotechnol.* (2021) doi:10.1038/s41587-021-00965-w.

59. Zheng, Y. F., Chen, Z. C., Shi, Z. X., Hu, K. H. & Zhong, J. Y. HIT-scISOseq: High-throughput and high-accuracy single-cell full-length isoform sequencing for corneal epithelium. *bioRxiv* (2020).

60. Tabula Muris Consortium *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature* **562**, 367–372 (2018).

61. Quake, S. R. & Sapiens Consortium, T. The Tabula Sapiens: a single cell transcriptomic atlas of multiple organs from individual human donors. *bioRxiv* (2021).

62. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**, 10–12 (2011).

63. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).