1    **Enhancing transcriptome expression quantification through accurate assignment of long RNA**

2    **sequencing reads with TranSigner**

3

4    **Hyun Joo Ji[1,2,*] and Mihaela Pertea[1,2,3,*]**

5    [1]Center for Computational Biology, Johns Hopkins University; Baltimore, MD

6    [2]Department of Computer Science, Johns Hopkins University; Baltimore, MD

7    [3]Department of Biomedical Engineering, Johns Hopkins University; Baltimore, MD

8    *corresponding authors: hji20@jh.edu, mpertea@jhu.edu

9

10    **Keywords:** long-read RNA sequencing, transcriptomics, expression quantification

11

12    **Abstract**

13

14    Recently developed long-read RNA sequencing technologies promise to provide a more accurate and

15    comprehensive view of transcriptomes compared to short-read sequencers, primarily due to their

16    capability to achieve full-length sequencing of transcripts. However, realizing this potential requires

17    computational tools tailored to process long reads, which exhibit a higher error rate than short reads.

18    Existing methods for assembling and quantifying long-read data often disagree on expressed transcripts

19    and their abundance levels, leading researchers to lack confidence in the transcriptomes produced using

20    this data. One approach to address the uncertainties in transcriptome assembly and quantification is by

21    assigning the long reads to transcripts, enabling a more detailed characterization of transcript support at

22    the read level. Here, we introduce TranSigner, a versatile tool that assigns long reads to any input

23    transcriptome. TranSigner consists of three consecutive modules performing: read alignment to the given

24    transcripts, computation of read-to-transcript compatibility based on alignment scores and positions, and

25    execution of an expectation-maximization algorithm to probabilistically assign reads to transcripts and

26    estimate transcript abundances. Using simulated data and experimental datasets from three well-studied

27    organisms — *Homo sapiens*, *Arabidopsis thaliana,* and *Mus musculus* — we demonstrate that TranSigner

28    achieves accurate read assignments, obtaining higher accuracy in transcript abundance estimation

29    compared to existing tools.

30

31    **Background**

32

33    Long-read RNA sequencing (RNA-seq) represents a remarkable advancement towards achieving full-

34    length sequencing of transcripts, offering novel insights into transcriptomes previously characterized only

35    with short reads. Short-read sequencing data has limitations in several applications such as transcript

36    assembly, primarily due to its fragmented nature and inherent biases (e.g., GC content, amplification) that

37    add noise to downstream analyses (Benjamini & Speed, 2012; Hansen et al., 2010; Li et al., 2009). Long-

38    read sequencing technologies address these limitations by substantially increasing the read lengths,

39    allowing each read to generally cover a full-length transcript, and employing strategies such as direct

40    RNA sequencing to reduce biases. Consequently, long-read data can provide more comprehensive and

41    accurate profiles of complex transcriptomes.

42

43    However, despite their potential, the full capabilities of long-read RNA-seq remain untapped due to the

44    limited inventory of tools optimized for analyzing long-read data. Although tools such as FLAIR (Tang et

45    al., 2020), Bambu (Chen et al., 2023), ESPRESSO (Gao et al., 2023), and StringTie2 (Kovaka et al.,

46    2019) are designed to characterize transcriptomes by both identifying novel isoforms and quantifying

47    transcripts using long-read RNA-seq data, their results often lack agreement (Chen et al., 2023; Gao et al.,

48    2023; Pardo-Palacios et al., 2023; Tang et al., 2020).

49

50    One way to address uncertainties in transcriptome assemblies is by assigning specific long reads to

51    transcripts. This allows for a more in-depth evaluation of the read-level support for transcripts, as opposed

52    to relying on read counts only. Given read-to-transcript assignments, transcripts can be directly associated

53    with a distribution of supporting read lengths, quality scores, alignment positions, and more. These

54    expanded sets of features can be used to derive a more confident set of transcripts and improve the

55    accuracy of transcript abundance estimates.

56

57    Few tools, including FLAIR and Bambu, track read-to-transcript assignments, but this functionality is

58    integrated into more complex pipelines that also identify novel isoforms in addition to quantifying known

59    transcripts. A standalone tool capable of performing read assignment and quantification on any input

60    transcriptome can be paired with other methods focusing on transcriptome assembly and could therefore

61    enable users to investigate any transcriptome of their choice. However, this need remains largely unmet,

62    with only a few recent methods, namely NanoCount (Gleeson et al., 2021), attempting to address it by

63    quantifying transcripts, yet still lacking the ability to assign specific reads to transcripts.
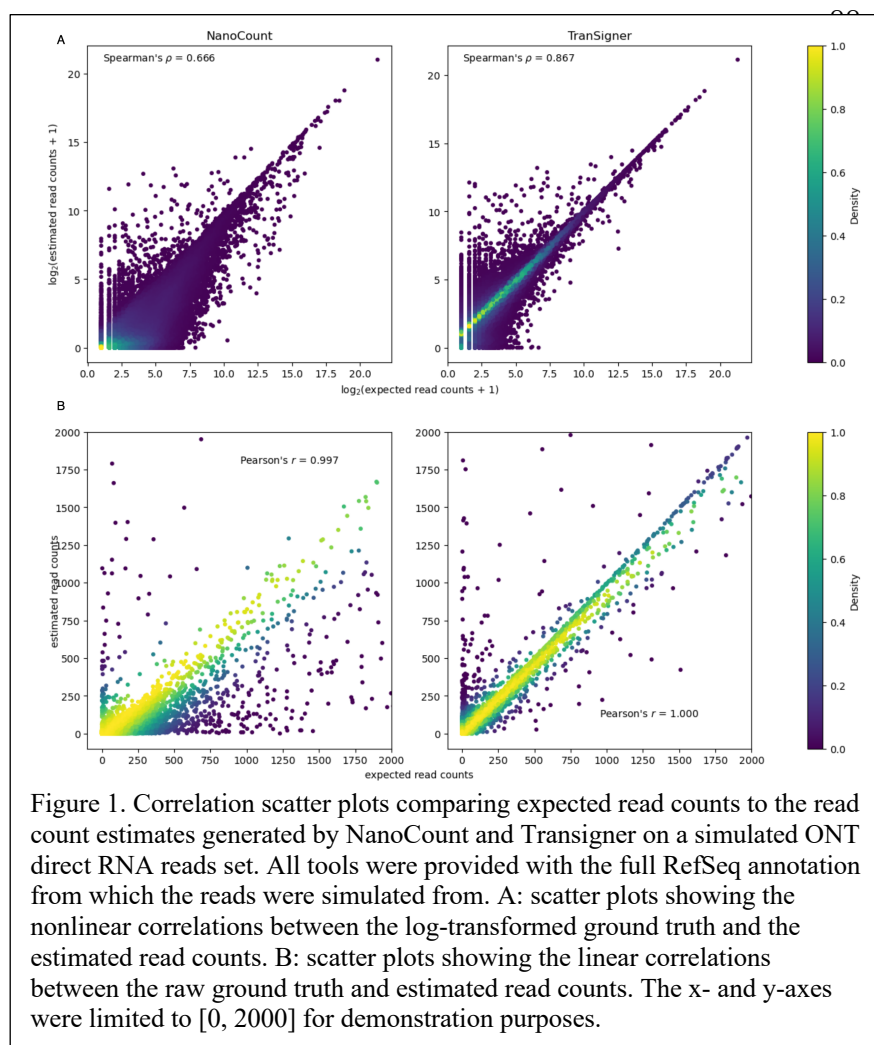
64

65    Here we introduce TranSigner, a novel transcript quantification-only method that accurately assigns long

66    RNA-seq reads to any given transcriptome. TranSigner first maps reads onto the transcriptome using

67    minimap2 (Li, 2018, 2021) and extracts specific features from the alignments, such as alignment scores or

68    the 3' and 5' end read positions on a transcript. These features are then utilized to compute compatibility

69    scores between read and transcript pairs, which indicate the likelihood of a read to originate from a

70    specific transcript. TranSigner then employs an expectation-maximization (EM) algorithm to derive

71    maximum likelihood (ML) estimates for both the read-to-transcript assignments and transcript

72    abundances simultaneously. We show that by guiding the EM algorithm in the expectation step with

73    precomputed compatibility scores, TranSigner generates high-confidence read-to-transcript mappings and

74    improves transcript abundance estimates.

75

76    **Results**

77

3

78    **Simulated data performance**. We first compared TranSigner against an existing quantification-only

79    tool, NanoCount (Gleeson et al., 2021). We benchmarked all three tools using five sets of simulated ONT

80    reads: three sets of direct RNA reads and two sets of cDNA reads. The reads were simulated from

81    protein-coding and long non-coding transcripts in the GRCh38 RefSeq annotation (release 110), and then

82    each tool was provided with both the simulated reads as well as the full RefSeq annotation as the target

83    transcriptome (see Methods for a full description of the simulated datasets). For simplicity, we will refer

84    to the transcripts from which the reads were simulated as the origin transcripts. To estimate how

85    accurately a tool assigns a read to its respective origin, we conducted both linear and nonlinear correlation

86    analyses between the expected read counts and each tool's estimates, using Pearson's correlation

87    coefficients (PCCs) between raw read counts and Spearman's correlation coefficients (SCCs) between

log-transformed read counts, respectively. A linear correlation analysis evaluates the ability of a tool to assign each read to a transcript, while a nonlinear correlation analysis assesses how well estimates capture monotonic trends in gene expression patterns.



Figure 1. Correlation scatter plots comparing expected read counts to the read count estimates generated by NanoCount and Transiger on a simulated ONT direct RNA reads set. All tools were provided with the full RefSeq annotation from which the reads were simulated from. A: scatter plots showing the nonlinear correlations between the log-transformed ground truth and the estimated read counts. B: scatter plots showing the linear correlations between the raw ground truth and estimated read counts. The x- and y-axes were limited to [0, 2000] for demonstration purposes.

In both analyses, we observed that TranSigner's estimates had stronger correlations with the ground truth compared to NanoCount's, as illustrated in

4

104    Figure 1, which shows results from one dataset typical of all three simulated ONT direct RNA datasets

105    (see Supplementary Table S3 for the SCC and PCC values on each read set). In both log-transformed

106    (Figure 1A) and raw (Figure 1B) read count correlation scatter plots, TranSigner shows higher

107    concentrations of dots near the diagonal. However, this feature is not as pronounced in the plots of

108    NanoCount's results; the accumulations of dots well below the diagonal in the case of NanoCount reveal

109    the tool's tendency to underestimate the read counts. On the simulated ONT direct RNA datasets,

110    TranSigner's average SCC and PCC values were 0.867 and 0.999, whereas NanoCount's were 0.667 and

111    0.997. TranSigner also achieves higher correlations with the ground truth when applied to the simulated

112    ONT cDNA datasets (see Supplementary Figure S1, Supplementary Tables S4).

113

114    Even for extensively studied species, gene annotation catalogs are often incomplete, missing both

115    potential gene loci and many transcript

116    isoforms (Amaral et al., 2023;

117    Varabyou et al., 2023). This is one

118    reason why most long-read processing

119    tools identify which transcripts are

120    present before quantification.

121    Identifying novel isoforms not present

122    in the annotation, as well as

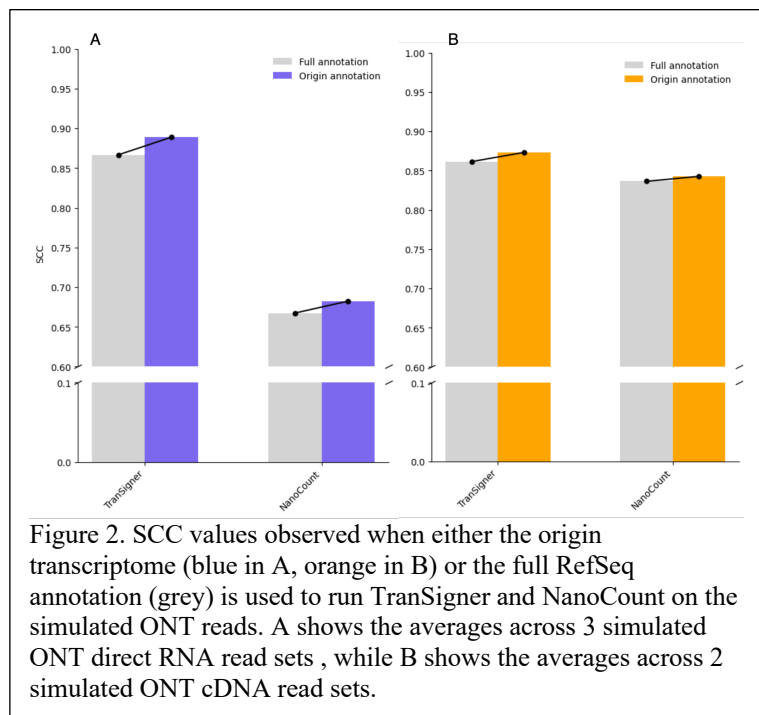123    determining which of the known

124    mRNA variants are expressed can lead

125    to better quantification of expressed



Figure 2. SCC values observed when either the origin transcriptome (blue in A, orange in B) or the full RefSeq annotation (grey) is used to run TranSigner and NanoCount on the simulated ONT reads. A shows the averages across 3 simulated ONT direct RNA read sets , while B shows the averages across 2 simulated ONT cDNA read sets.
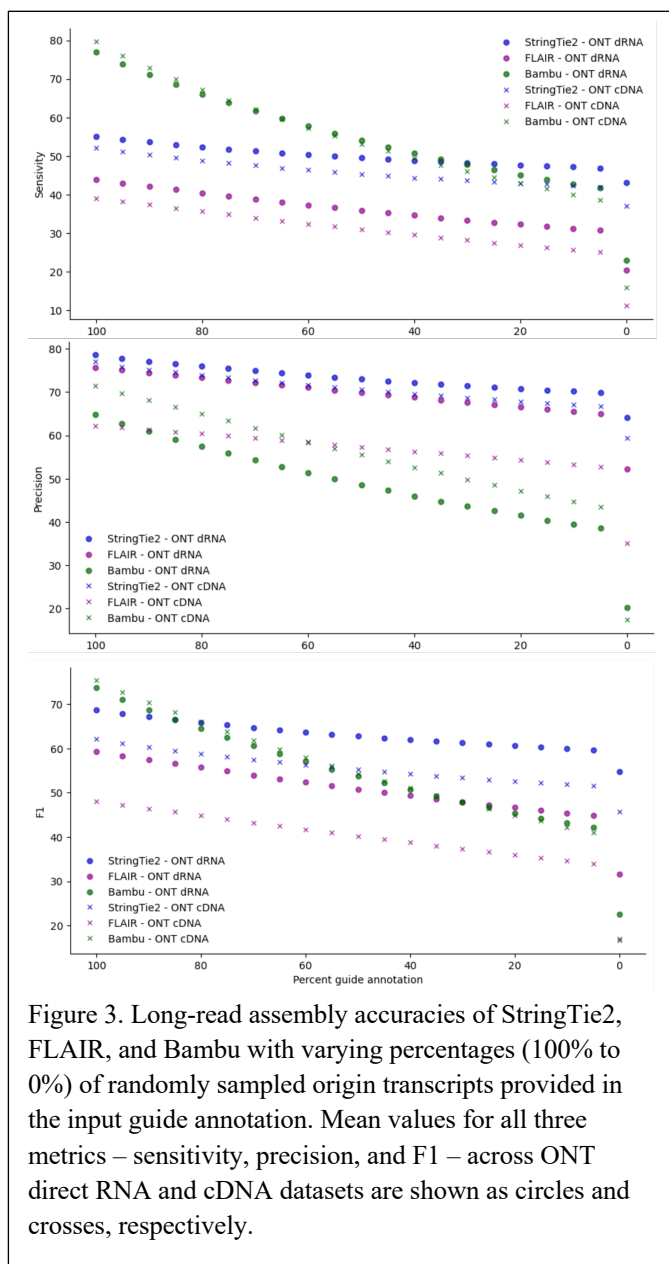
126    transcripts. This is illustrated by our results in Figure 2, where we show that the average nonlinear

127    correlation coefficients between estimated and true read counts improve for both TranSigner and

128    NanoCount when just the origin transcripts are provided in the input instead of the full reference

129    annotation (see Supplementary Tables S3 and S4 for SCC and PCC values across all simulated ONT direct RNA and cDNA data sets).



Figure 3. Long-read assembly accuracies of StringTie2, FLAIR, and Bambu with varying percentages (100% to 0%) of randomly sampled origin transcripts provided in the input guide annotation. Mean values for all three metrics – sensitivity, precision, and F1 – across ONT direct RNA and cDNA datasets are shown as circles and crosses, respectively.

Achieving an accurate transcriptome remains a challenging problem, with different tools obtaining varying accuracies in this task, while also relying to varying degrees on the input reference annotation. Using the same simulated ONT data sets (3 direct RNA, 2 cDNA) we used to benchmark TranSigner and NanoCount, we evaluated existing tools' ability to handle incompleteness in the input guide annotations. To do this, we randomly sampled the full RefSeq annotation to include varying percentages–between 0% and 100% with increments of 5%–of the origin transcripts and provided the resulting annotations as guides to StringTie2, FLAIR, and Bambu. We did not include ESPRESSO in this comparison, as processing a single simulated data set took

149    more than 24h to process. We also randomly sampled each percentage of retained origin transcripts three

150    times (see Methods for further details).

151

152    Genome-guided transcriptome assemblers like StringTie2 (Kovaka et al., 2019) can reliably profile a

153    transcriptome even in the absence of an input guide annotation, while methods like Bambu (Chen et al.,

154    2023) or FLAIR (Tang et al., 2020) demonstrate a substantial decrease in both sensitivity and precision of

155    transcript identification when the percentage of origin transcripts in the input guide annotation is

156    progressively reduced. Figure 3 shows that while Bambu outperforms StringTie2 and FLAIR in terms of

157    average sensitivity when a substantial portion of the origin transcriptome is provided in the input,

158    StringTie2 consistently outperforms the rest of the tools in precision across all percentages of origin

159    transcripts kept in the input annotation. Bambu achieved highest F1 scores when the guides retained most

160    of origin transcripts, but StringTie2 gradually surpassed others as guides became increasingly incomplete

161    (Figure 3C). Such resilience to varying degrees of incompleteness in the input transcriptome is critical,

162    especially for studies involving poorly annotated organisms or in cases where the RNA-seq sample

163    contains many novel isoforms (see Supplementary Tables S1 and S2 for the metric values on each

164    dataset). However, StringTie2 does not assign individual reads to the transcripts it assembles, making it

165    difficult for the user to check the reliability of the isoforms it assembles using long reads. By introducing

166    TranSigner, we aimed to also address this gap, in addition to improving transcript quantification

167    accuracies.

168

169    Next, we compared TranSigner's quantification accuracies against those of several other tools –

170    StringTie2, NanoCount, Bambu, and FLAIR – when provided with guide annotations containing varying

171    percentages of the origin transcripts. Since TranSigner is not capable of identifying novel transcripts, we

172    also ran TranSigner on the transcriptome assembled by StringTie2 (denoted as StringTie2 + TranSigner)

173    to investigate its performance against other tools, such as FLAIR or Bambu, which are capable of novel

174    isoform identification. For this experiment, we re-used the same sets of simulated ONT reads and 5% ~
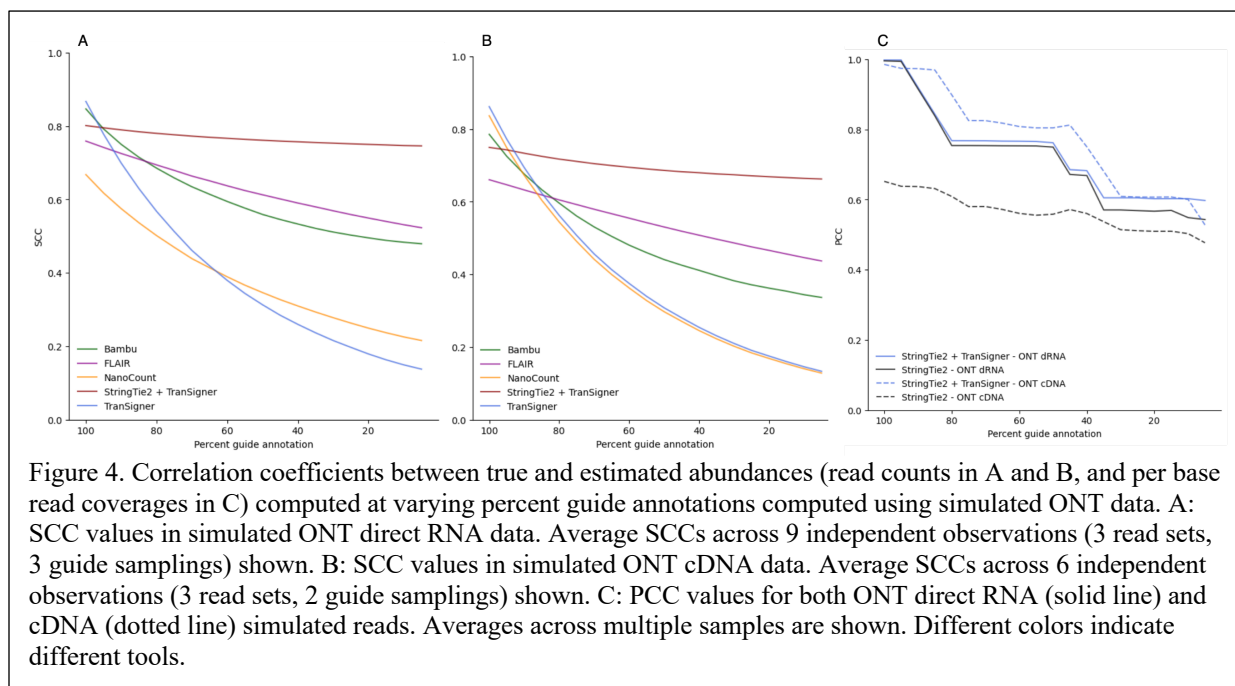
175    100% guide annotations sampled before.



Figure 4. Correlation coefficients between true and estimated abundances (read counts in A and B, and per base read coverages in C) computed at varying percent guide annotations computed using simulated ONT data. A: SCC values in simulated ONT direct RNA data. Average SCCs across 9 independent observations (3 read sets, 3 guide samplings) shown. B: SCC values in simulated ONT cDNA data. Average SCCs across 6 independent observations (3 read sets, 2 guide samplings) shown. C: PCC values for both ONT direct RNA (solid line) and cDNA (dotted line) simulated reads. Averages across multiple samples are shown. Different colors indicate different tools.

176

177    Average correlation coefficients between the true and estimated read counts are shown in Figure 4 (also

178    see Supplementary Tables S5 and S6 for results on all input datasets). Except for StringTie2 +

179    TranSigner, every tool experienced a drastic drop in SCC values as the percentage of origin transcripts

180    decreased. TranSigner had the highest correlation values when the input guide annotation contained

181    nearly all origin transcripts. However, when 90% or fewer of the origin transcripts were retained in the

182    guide annotation, StringTie2 + TranSigner yielded the best SCC values in both ONT direct RNA and

183    cDNA benchmarks  (Figure 4A, 4B), demonstrating that this combination is the best in preserving the

184    rank of the expression values across most levels of incompleteness in the available annotation. This same

185    pattern holds for PCC values (Supplementary Figure S2A, S2B). StringTie2 does not output read counts

186    for its transcript abundance estimates, so it was excluded from this initial correlation analysis. As

187    StringTie2 outputs read per base coverages, we post-processed TranSigner's read-to-transcript

188    assignments to generate read per base coverages (see Methods). TranSigner + StringTie2 obtains better
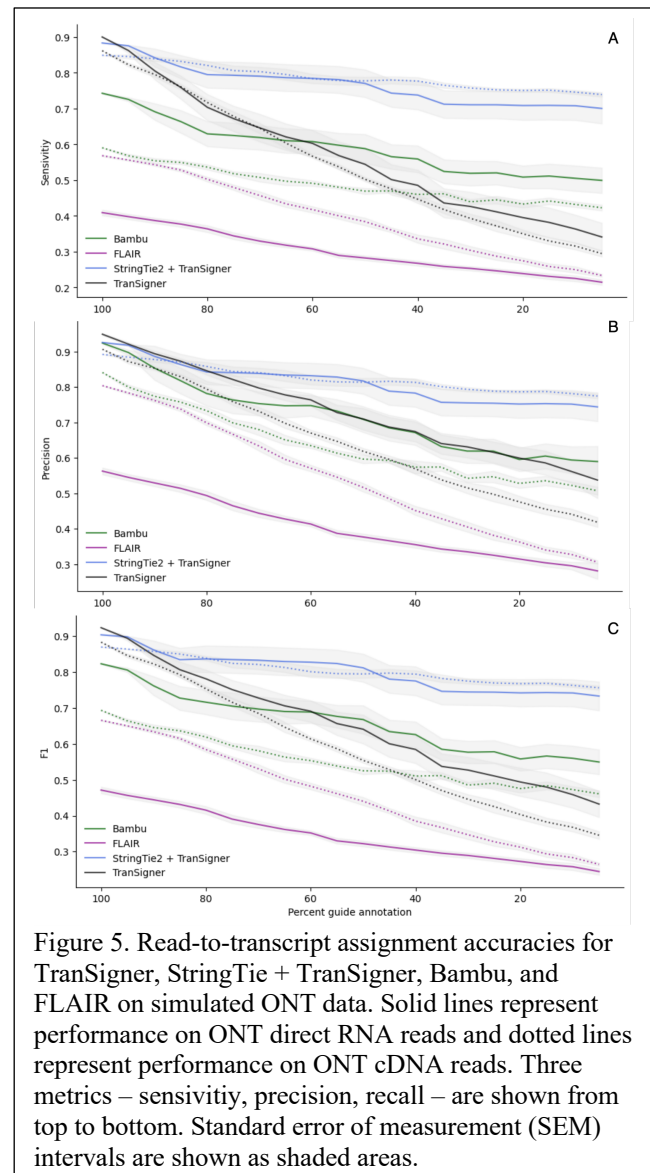
189  read per base coverage PCCs (Figure 4C) and SCCs (Supplementary Figure S2C) correlation values than

190  StringTie2. The improvement is more notable in PCCs than in SCCs.

191

192  One key feature of TranSigner is its ability to

193  assign specific reads to transcripts, particularly

194  useful in experiments where users need to

195  identify reads originating from specific

196  transcripts of interest. In this context, we

197  compared TranSigner and StringTie2 +

198  TranSigner with FLAIR and Bambu, which also

199  output read-to-transcript assignments. Their

200  performance was evaluated using recall,

201  precision, and F1 scores, computed by counting

202  the number of correctly versus incorrectly

203  assigned reads (see Methods). When all origin

204  transcripts are provided (i.e., 100% complete

205  guide annotation), TranSigner demonstrated the

206  highest sensitivity, recall, and hence F1 score

207  (see Figure 5, Supplementary Tables S7 and S8).

208  However, as soon as the guides become even

209  slightly incomplete, StringTie2 + TranSigner



Figure 5. Read-to-transcript assignment accuracies for TranSigner, StringTie + TranSigner, Bambu, and FLAIR on simulated ONT data. Solid lines represent performance on ONT direct RNA reads and dotted lines represent performance on ONT cDNA reads. Three metrics – sensitivity, precision, recall – are shown from top to bottom. Standard error of measurement (SEM) intervals are shown as shaded areas.

210  had the highest performance, making it the preferred choice when the target transcriptome is 95% or less

211  complete.

212

213  Although TranSigner achieved the highest F1 scores with nearly complete guides, its performance

214  declined rapidly as the number of origin transcripts in the guides decreased, as expected (Figure 5C). A

9

215    similar pattern of decline is observed in every tool across all metrics. Bambu experienced a greater drop

216    in precision than StringTie2 + TranSigner, despite both starting at a similar value. Note that both Bambu

217    and FLAIR showed fluctuations in performance depending on the ONT read types. In contrast, StringTie2

218    + TranSigner showed the least amount of variation in performance across different read types.

219

220    **Real data performance**. To evaluate the performance of TranSigner and StringTie2 + TranSigner using

221    experimental data, we utilized the ONT RNA-seq data sets provided by the Singapore Nanopore

222    Expression Project (SG-NEx) (Chen et al., 2021), which include synthetic spike-in transcripts, known as

223    sequins, with known annotation and concentrations. We selected 12 ONT direct RNA and cDNA samples

224    from three different human cell lines: HCT116, K562, and MCF7. As ground truth for this experiment,

225    we used the counts per million (CPM) values provided by SG-NEx and compared them with the estimates

226    obtained by TranSigner, StringTie2 + TranSigner, and Bambu, the next best performer on the simulated

227    data. We ran StringTie2 + TranSigner and Bambu twice, each time providing two different input guides:

228    one including the full sequin annotation in addition to the GRCh38 reference annotation and the other
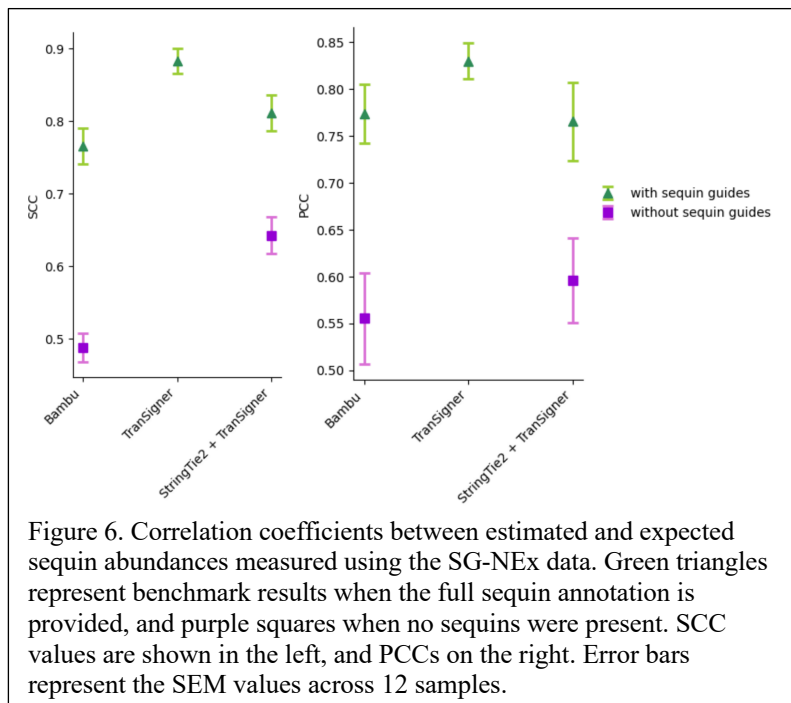


Figure 6. Correlation coefficients between estimated and expected sequin abundances measured using the SG-NEx data. Green triangles represent benchmark results when the full sequin annotation is provided, and purple squares when no sequins were present. SCC values are shown in the left, and PCCs on the right. Error bars represent the SEM values across 12 samples.
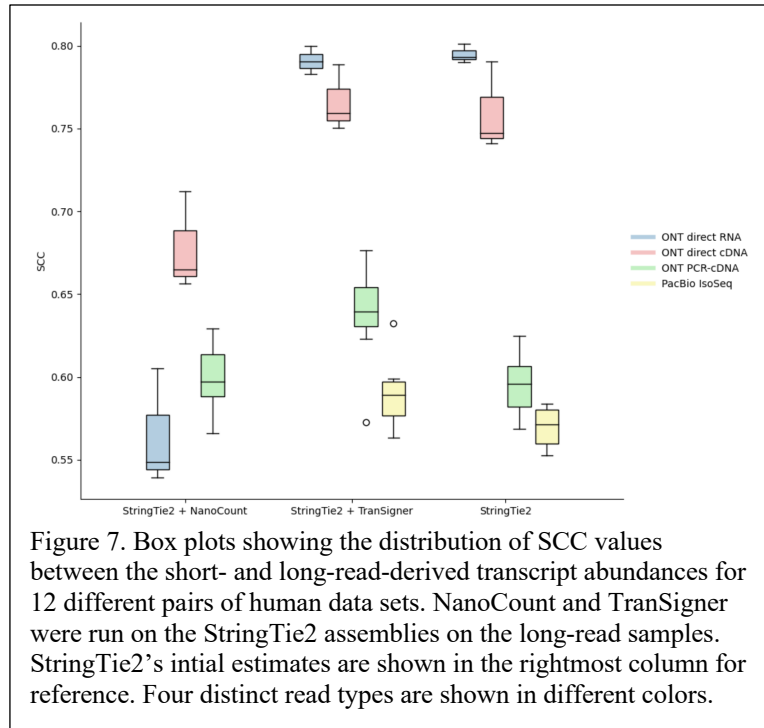
containing only the GRCh38 reference transcripts without the sequins. The guide annotation without the sequins reflects real-world scenarios where transcript annotations are absent from the reference. TranSigner was only run with the full sequin annotation, as it cannot assemble any novel transcripts itself.

239

10

240    TranSigner achieved an average SCC of 0.88 between ground truth and estimated values, surpassing both

241    Bambu (0.76) and StringTie2 + TranSigner (0.81) when provided with the full sequin annotation, as

242    displayed in Figure 6 (also see Supplementary Tables S9). However, when no sequin annotation was

243    provided, StringTie2 + TranSigner outperformed Bambu, obtaining an average SCC value of 0.64,

244    compared to Bambu's SCC value of 0.49. This trend persisted in linear correlation analyses, with

245    TranSigner achieving the highest PCC value with full annotation (0.83), while StringTie2 + TranSigner's

246    was the best performer in the absence of sequin annotation (Figure 6 and Supplementary Tables S9).

247    Overall, these results suggest that StringTie2 + TranSigner may be preferable in scenarios where

248    numerous unannotated or novel isoforms are anticipated, while TranSigner is optimal when the reference

249    is presumed to be nearly complete. Note that with complete sequin annotation, TranSigner outperformed

250    both Bambu and StringTie2 + TranSigner, on all three different long-read types available in the data:

251    direct RNA, direct cDNA, PCR-cDNA (average and per-sample SCC values shown in Supplementary

252    Figure S3 and Supplementary Tables S9).

253

254    We also evaluated the correlation between short-read-based and long-read-based abundance estimates

255    using publicly available paired short and long-read datasets, sequenced from the same biological sample.

256    In all following results, the short-read libraries were all generated through poly-A selection and

257    sequenced with Illumina sequencers, while the long reads were mostly generated using ONT direct RNA

258    or cDNA sequencing protocols. Unlike the sequin samples or simulated long reads, the ground truth is

259    unknown for these datasets as we lack information on which transcripts are expressed and their relative

260    abundances. However, it is generally assumed that short reads provide more accurate abundance estimates

261    compared to long reads, as they are less error-prone and typically yield more reads.

262

263 Specifically, we assessed the long

264 read-based abundance estimates by

265 two quantification-only tools we

266 benchmarked with simulated data:

267 NanoCount and TranSigner. All tools

268 were provided with a StringTie2-

269 assembled transcriptome, which

270 represents a typical use for these

271 tools where users provide

272 transcriptomes assembled from

273 samples of their interest. We used



Figure 7. Box plots showing the distribution of SCC values between the short- and long-read-derived transcript abundances for 12 different pairs of human data sets. NanoCount and TranSigner were run on the StringTie2 assemblies on the long-read samples. StringTie2's intial estimates are shown in the rightmost column for reference. Four distinct read types are shown in different colors.

274 each tool's abundance estimates to conduct nonlinear correlation analyses between the short read-derived

275 TPM estimates and long read-derived CPM. As previously done for benchmarking long-read

276 quantification tools (Pardo-Palacios et al., 2023), we assumed that a higher correlation between long read-

277 and short read-derived abundance estimates is indicative of a higher quantification accuracy. Since none

278 of the three quantification-only tools we used include TPMs in their output, we processed the read counts

279 they provide to obtain counts per million (CPM) estimates, which are equivalent to TPMs in a long-read

280 RNA-seq experiment where each read is considered to represent a transcript (see Methods for the read

281 counts to CPM conversion equation). We used Salmon (Patro et al., 2017) to obtain TPM estimates on

282 StringTie2 assemblies, using the Illumina short-read datasets (see Supplementary Text 3). As transcripts

283 with low abundances are prone to misassembly and are often excluded from downstream analyses, we

284 only included in our results transcripts with > 1 TPM as estimated by Salmon.

285

286 For our first experiment, we chose 21 short and long read paired datasets: 9 pairs from two normal human

287 cell lines, A549 and HCT116, included in the SG-NEx datasets (Chen et al., 2021), and 12 pairs from two

288 human cancer cell lines, H1975 and HCC827, provided by the long-read benchmarking of human lung

289 cancer cell lines (Dong et al., 2023). The human lung cancer cell lines data sets also included PacBio

290 reads, which are not present in the SG-NEx data sets. As shown in Figure 7, TranSigner consistently

291 achieved higher correlations than NanoCount as well as StringTie2, across all read types (see

292 Supplementary Tables S10 for the SCC and PCC values on each pair). TranSigner improved StringTie2's

estimates to varying degrees, with the highest improvements observed in the ONT PCR-cDNA data sets. Note that NanoCount was not evaluated on PacBio data as it was designed specifically to work with ONT data only.
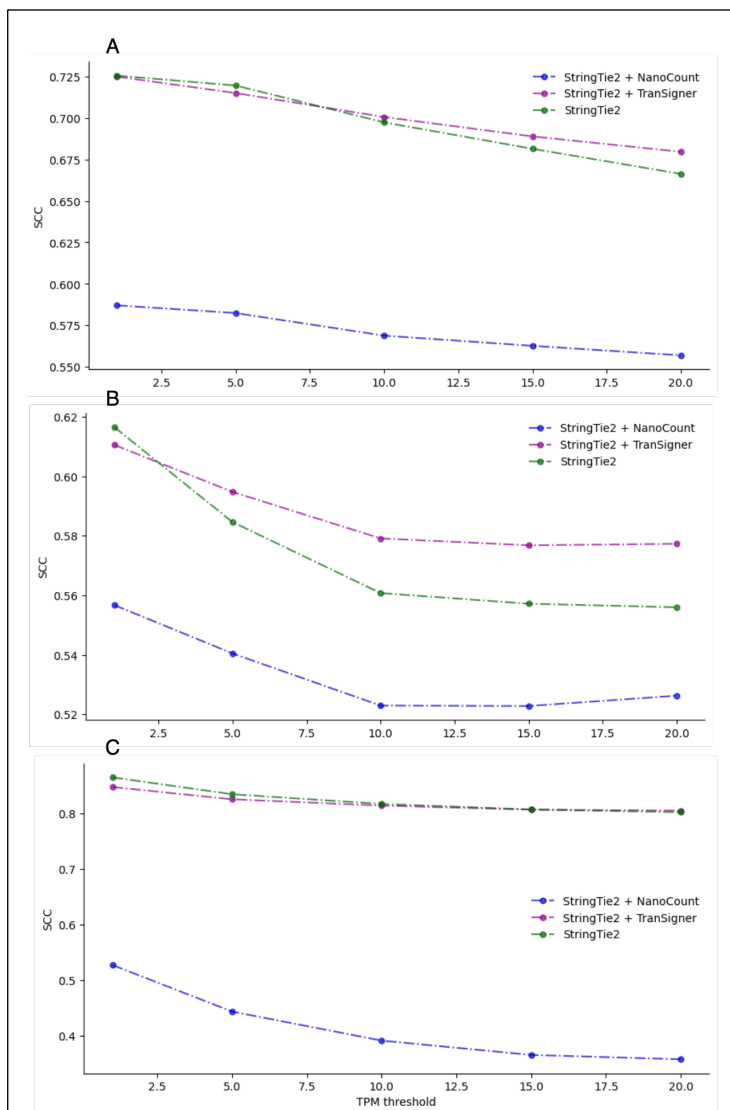


Figure 8. Correlation coefficient values between short- and long-read-derived transcript abundances estimated by NanoCount and TranSigner when run on StringTie2 assemblies, as well as StringTie2 itself, on paired *M. musculus* (A and B) and *A. thaliana* data sets (C). Each plot is showing a different organism and a different read type. A: average SCC values across increasing TPM thresholds on *M. musculus* ONT direct RNA data sets. B: average SCC values across increasing TPM thresholds on *M. musculus* ONT PCR-cDNA data sets. C: average SCC values across increasing TPM thresholds on *A. thaliana* ONT direct RNA data sets.

Finally, we further expanded our benchmark to include paired short- and long-read data sets from two well-studies species: *A. thaliana* and *M. musculus*. To investigate how quantification accuracies vary at different levels of expression, we evaluated the performance of StringTie2 and StringTie2 + <a quantification-only tool> at progressively increasing TPM thresholds: 1, 5, 10, 15, and 20. For this experiment, we selected eight *M. musculus* pairs (four ONT direct RNA, four ONT cDNA) and three *A. athaliana* pairs (all ONT direct RNA). We benchmarked TranSigner's and

13

315    NanoCount's performances when run on unguided StringTie2 assemblies, consistent with the previous

316    analysis. As illustrated in Figure 8, when TranSigner was applied to StringTie2's output, it achieved

317    higher nonlinear correlations between short- and long-read TPM estimates than NanoCount, with the best

318    improvements in SCC values obtained on the *M. musculus* ONT PCR-cDNA reads. These improvements

319    were more pronounced for higher TPM thresholds.

320

321    **Discussion and Conclusions**

322

323    Assigning long reads to transcripts is a challenging task that involves the effective resolution of multi-

324    mapping reads. Recent studies have unveiled the growing complexity of eukaryotic transcriptomes,

325    revealing numerous isoforms across gene loci. The introduction of long-read RNA-seq technologies

326    promises to uncover even more novel isoforms, as reads produced by these methodologies can capture

327    full-length transcripts, overcoming the limitations of short reads. Although long reads cover transcripts at

328    greater lengths, technical artifacts such as base calling errors and end truncations prevent these reads from

329    being accurately mapped to their origins. With TranSigner, we have developed several strategies to

330    address this challenge, facilitating the correct assignment of reads that ambiguously map to multiple

331    isoforms.

332

333    Additionally, we designed TranSigner to complement another method capable of transcriptome assembly.

334    As gene annotation is still an unresolved issue, determining the accuracy and completeness of a profiled

335    transcriptome remains difficult. Users often struggle to select the appropriate reference for their analyses,

336    leading to unpredictable impacts on their results. In our study, we observed a significant drop in assembly

337    quality when less complete guides were provided. This suggests that tools heavily reliant on high-quality

338    reference annotations may struggle in real-world scenarios where many novel isoforms are expected. By

339    introducing a standalone tool for read-to-transcript assignments, we made these assignments easier to

340    obtain regardless of the input transcriptome. Integrating this step into long-read RNA-seq data processing

341    pipelines will improve the accuracy of transcriptomes identified using long reads by allowing users to

342    inspect the quality of the reads supporting the transcripts and filter out less-supported transcripts. This, in

343    turn, will lead to more accurate abundance estimates, as our results demonstrate the significant influence

344    of assembly accuracy on correctly identifying transcript abundances.

345

346    **Methods**

347

348    **Long-read RNA-seq model for read assignment.** We describe the long-read RNA-seq process using a

generative model (Figure 9). The conceptualization of RNA-seq as a generative process in which reads are sampled from a pool of transcripts has already been used in models for short-read quantification. We adopted the general framework proposed by others (Li et al., 2009; Pachter, 2011) but introduced necessary modifications to tailor the model to long read data. Given a read, we assume that three unobserved events in the RNA-seq experiment determine a read's sequence: (1) the transcript from which that read was sequenced, (2) the position within the transcript of the 3' end of the read, and (3) the transcript position of the reads' 5' end. Our model, thus, associates each observed read with three latent variables: the transcript ($T$) from which the read was generated, its 3' end position ($S$), and 5' end position ($E$) in $T$.
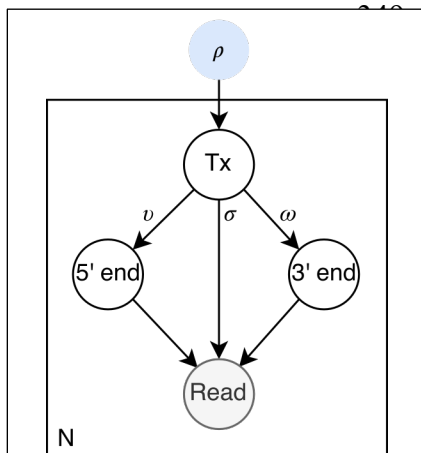


Figure 9. Graphical representation of TranSigner's long-read RNA-seq model. Empty circles denote latent variables, the shaded circle represents the observed variable, and the blue circle indicates the primary parameter of the model – specifically, the relative abundance of the transcript. Parameters $v$, $\omega$ approximate the likelihood of the specific 5' and 3' end positions of the read on the transcript, while parameter $\sigma$ models the likelihood of observing a specific read sequence given a transcript and the read's end positions. $N$ represents the total number of reads generated in a single long-read RNA-seq experiment.

Existing RNA-seq quantification methods focus on accurately

365    estimating $\rho$, the relative transcript abundances (Jousheghani & Patro, 2024; Li et al., 2009; Pachter,

366    2011). In contrast, our primary goal here is to assign reads to transcripts, which is solved by finding the

15

367 most probable distributions over the latent variables, not $\rho$. However, deriving a maximum likelihood

368 (ML) estimate on $\rho$ also gets us ML estimates on the latent variable distributions, as they get repeatedly

369 updated in the process of optimization. Hence, $\rho$ is still the main parameter to optimize, and we define our

370 objective with respect to $\rho$ as follows. Given a set of transcripts $T = \{t\}$ where $|T| = M$, the complete

371 data likelihood function of our RNA-seq model is:

372

373
$$\mathcal{L}(\rho) = \prod_{r \in R} \sum_{t \in T} P(r \in t|\rho)P(s_{rt}|r \in t)P(e_{rt}|r \in t)P(r|r \in t, s_{rt}, e_{rt})$$

374
$$(1)$$

375 where $\rho = \{\rho_t\}_{t \in T}$ with $\sum_{t \in T} \rho_t = 1$, $R$ is the set of mapped reads defined as $R = \{r\}$ with the cardinality

376 of $N$, $s_{rt}$ and $e_{rt}$ are the 3' and 5' end positions of a read $r$ in a transcript $t$, and $r \in t$ indicates that $r$

377 comes from $t$. Note $P(O_r = t|\rho) = \rho_t$, since in an RNA-seq experiment the probability of selecting a

378 transcript $t$ to sequence depends on its relative abundance. We'll approximate the 5' end and 3' end

379 positions of a read in a transcript as the positions where the read alignment starts and ends on that

380 transcript, respectively. The relationship between this likelihood function and read assignment estimates

381 is easier to understand when Eq. 1 is rewritten as:

382

383
$$\mathcal{L}(\rho) = \prod_{r \in R} P(r) = \prod_{r \in R} \sum_{t \in T} P(r|r \in t) = \prod_{r \in R} \sum_{t \in T} \alpha_{rt}$$

384
$$(2)$$

385 where $\alpha_{rt}$ is the relative fraction of read $r$ assigned to transcript $t$. $P(r)$ can also be written as a sum of

386 conditional probabilities $P(r|r \in t)$, which represents the likelihood of $r$ given that it comes from $t$. This

387 conditional probability is also easily interpretable as the fraction of $r$ that ought to be assigned to $t$,

388 implying that a lower $P(r|r \in t)$ corresponds to a smaller $\alpha_{rt}$. Moreover, optimizing $\mathcal{L}$ involves driving

389 $P(r)$ to the maximum possible value in a probability distribution – 1, which is also equal to the sum of

390 relative fractions of a read's assignments to the set of transcripts (i.e., $\sum_{t \in T} \alpha_{rt} = 1$).

16

391

392  Different long-read RNA-seq technologies show various biases towards the ends of the transcripts

393  (Amarasinghe et al., 2020; Chen et al., 2021; Grünberger et al., 2022; Wongsurawat et al., 2022).

394  Nonetheless, long reads are more likely to cover all bases of a transcript, compared to short reads, which

395  are generated from fragments of the transcript. The likelihood of a read's end position should decrease as

396  its distance from the transcript end increases. We model this expectation using two indicator variables– $v$

397  and $\omega$ for the 3' and 5' ends, respectively – to control how far apart the ends of a read can be from the

398  ends of a transcript. For an alignment between a read $r$ and a transcript $t$, we will refer to the distances

399  between the alignment ends and transcript ends as 'end distances' and denote them as $\delta_s^{rt}$ and $\delta_e^{rt}$ for the

400  5' and 3' ends, respectively. Then we define $v$ and $\omega$ as:

401

$$P(s_{rt} = i | r \in t) \approx v_{rt} = 1 \text{ if } \left|\delta_s^{rt'} - \delta_s^{rt}\right| \leq \beta_s, 0 \text{ o.w.}$$

$$P(e_{rt} = j | r \in t) \approx \omega_{rt} = 1 \text{ if } \left|\delta_e^{rt'} - \delta_e^{rt}\right| \leq \beta_e, 0 \text{ o.w.}$$

404                                                                                                (3)

405  where $\delta_s^{rt'}$ and $\delta_e^{rt'}$ represent the end distances for the primary alignment of read $r$ and transcript $t'$.

406

407  Here, $t'$ represents the transcript to which read $r$ aligns on the primary alignment, which might not be the

408  same as transcript $t$. Since alignment positions are indexed from the 5' to 3' direction on transcript $t$, end

409  distances are computed as $\delta_s^{rt} = s_{rt} = i$ and $\delta_e^{rt} = |t| - e_{rt} = |t| - j$ where $|t|$ is the length of transcript

410  $t$. Parameter $\beta$ represents the tolerance threshold on how much greater the end distances can be compared

411  to the primary alignment's end distances for a given read $r$. This relative thresholding on end distances

412  ($\delta$) ensures that each read is compatible with at least one transcript (i.e., $t'$) after this filtering step since

413  the primary alignment will always be considered "good," which would not be true if a constant threshold

414  was uniformly applied for all reads. When either $v$ or $\omega$ is set to 0, $P(r | r \in t)$ in Eq. 2 is also set to 0, and

17

415    no fraction of $r$ is assigned to $t$, guaranteeing that the corresponding $(r, t)$ pair will be considered entirely

416    incompatible, filtering it out from any downstream analysis.

417

418    Moreover, the parameters for the 3' end are treated separately from those for the 5' end because

419    sequencing behaves differently at these ends. For example, there is a stronger coverage bias towards the

420    3' end when nanopore-based direct RNA sequencing protocols are employed (Amarasinghe et al., 2020;

421    Chen et al., 2021; Grünberger et al., 2022; Wongsurawat et al., 2022). We set the $\beta$ parameter values

422    based on both prior knowledge and a grid search (Supplementary Text 1). For the ONT direct RNA data,

423    the current default values are $\upsilon = -\infty$ (i.e., no filter) and $\omega = -800$, while for ONT cDNA and PacBio

424    data, they are $\upsilon = -500$ (i.e., unset) and $\omega = -550$ for ONT cDNA and PacBio data.

425

426    The probability of observing a read $r$ given all the latent variables is modeled using the alignment score

427    between read $r$ and transcript $t$ (denoted by $x_{rt}$) as:

428

429
$$P(r|r \in t, s_{rt} = i, e_{rt} = j) \approx \sigma_{rt} = \frac{x_{rt}}{\max\limits_{k \in T} x_{rk}}$$

430                                                                           (4)

431    Note that if multiple alignments exist between read $r$ and transcript $t$, we only retain the alignment with

432    the maximum score. Using the above definitions, we can redefine the likelihood function as:

433

434
$$\mathcal{L}(\rho) = \prod_{r \in R} \sum_{t \in T_r} \rho_t \upsilon_{rt} \omega_{rt} \sigma_{rt}$$

435                                                                           (5)

436    where $T_r$ is the set of transcripts aligned to read $r$, with $\upsilon_{rt}$, $\omega_{rt}$, and $\sigma_{rt}$ set to zero for any unaligned pair

437    of read $r$ and transcript $t$. By combining Eqs 2 and 5 we obtain that:

438

18

439
$$\rho_t \nu_{rt} \omega_{rt} \sigma_{rt} = \alpha_{rt}$$

440
(6)

441 which shows how $\alpha_{rt}$ can be computed from the alignments between reads and transcripts, assuming that

442 the relative transcript abundances are given.

443

444 **Alignment**. We used minimap2 with parameter -N 181 to align the long reads to the set of input

445 transcripts (Li, 2018, 2021). By default, minimap2 limits the maximum number of secondary alignments

446 to 5. We observed that the number of true positives (correct read to transcript alignments) increases when

447 we retain more secondary alignments, so we set -N to 181, the highest number of transcripts in a single

448 gene locus according to the RefSeq release 110 annotation on the human GRCh38 genome, assuming this

449 is the maximum number of secondary alignments a read can have. This strategy provides rough,

450 preliminary estimates on the compatibility between reads and transcripts, without excluding any read and

451 transcript pair for having suboptimal alignment scores. The user can freely adjust this parameter by

452 specifying it in TranSigner's input, which will then pass it to minimap2.

453

454 **Alignment-guided expectation-maximization algorithm (AG-EM)**. Our primary goal is to accurately

455 assign reads to their respective transcript origins. We previously introduced $\alpha$ as a variable representing

456 read-to-transcript assignments and established that the distribution over $\alpha$ is equivalent to that over the

457 latent variables of our long-read RNA-seq model (Figure 9 and Eqs. 1, 2, 3). An expectation-maximum

458 (EM) algorithm finds a maximum likelihood (ML) estimate for a main parameter (e.g., $\rho$) through

459 iterative updates to the distribution over a set of latent variables (e.g., $\alpha$). Hence, TranSigner employs an

460 EM algorithm to obtain the most probable–in the sense that the complete data likelihood is maximized–

461 distribution over $\alpha$ and presents the corresponding expected values as read-to-transcript assignments. It

462 also outputs the ML estimates on $\rho$.

463

464  *Update rules.* The EM algorithm consists of alternating expectation (E) and maximization (M) steps,

465  repeated until convergence. During the E step, the expected values for $\alpha_{rt}^{(n)}$–at some iteration $n$–are

466  computed as follows:

467
$$\alpha_{rt}^{(n)} = \frac{\rho_t^{(n)} v_{rt} \omega_{rt} \sigma_{rt}}{\sum_{t' \in T_r} \rho_{t'}^{(n)} v_{rt'} \omega_{rt'} \sigma_{rt'}}$$

468                                                                                                              (7)

469  where $\alpha = \{\alpha_{rt}\}_{r,t \in A}$ and $A$ is the set of alignments between all reads and transcripts. In the following M

470  step, then, the fragments of reads assigned to each transcript are summed up and then normalized by the

471  total number of transcripts to get the relative transcript abundances, expressed as:

472

473
$$\rho_t = \frac{\sum_{r \in R_t} \alpha_{rt}}{\sum_{r',t' \in A} \alpha_{r't'}}$$

474                                                                                                              (8)

475  where $R_t$ is the set of reads aligned to transcript $t$. The denominator is constant across iterations and is

476  equivalent to the total number of reads in a long-read RNA-seq experiment where each read represents a

477  transcript, so we precompute this value before EM.

478

479  *Initialization.* Before the EM iterations, the relative transcript abundances ($\rho$) are initialized to the

480  uniform distribution:

481
$$\rho_t = \frac{1}{|T_A|}$$

482

483  where $T_A$ is the set of transcripts with at least one alignment to a read in $R$. Additionally, the values for $v$,

484  $\omega$, and $\sigma$ don't change during iterations, so we precompute their values and store them separately in a

485  matrix $X$ of dimensions $N$ rows and $M$ columns. For simplicity, we'll refer to $X$ as the compatibility score

486  matrix. The computation specified in Eq. 7 is further simplified as:

20

487

$$\alpha_{rt}^{(n)} = \frac{\rho_t^{(n)} X_{rt}}{\sum_{t' \in T_r} \rho_{t'}^{(n)} X_{rt'}}$$

489                                                                                                      (9)

490     The pre-computation step involves a single scan over the alignment results, extracting values such as the

491     alignment scores and alignment start/end positions, and then applying the definitions provided in Eqs. 3

492     and 4.

493

494     *Optimization.* Once $X$ is precomputed and $\rho$ is initialized, EM iterations are repeated until convergence,

495     i.e., until the total sum of changes in the relative transcript abundances is less than a predefined threshold,

496     by default set at 0.005. The user can adjust this threshold to increase the accuracy of the ML estimates at

497     the expense of speed.

498

499     The novelty of our method comes from guiding the EM algorithm with the priors extracted from the

500     alignment results, as detailed in the E-step update rule shown in Eq. 9. To further amplify the impact of

501     these priors, we implemented an algorithm called the `drop`. The `drop` algorithm (Supplementary Figure

502     S4) sets $X_{rt} = 0$ if the fraction of read $r$ that is assigned to transcript $t$ (i.e., $\alpha_{rt}$) gets below a threshold,

503     $\tau \in [0,1]$. This effectively drops the compatibility relationship between read $r$ and transcript $t$ and

504     ensures that no fraction of $r$ gets assigned to $t$ in any iterations following the drop, as $\alpha_{rt}$ will always be

505     0 since its computation involves multiplication by $X_{rt}$ (Eq. 9). After the drop, another E-step is performed

506     with the updated $X$ scores to recompute the new $\alpha_{rt}$ values. The $\tau$ value depends on the read $r$ considered,

507     and by default:

508     $$\tau_r = \frac{1}{|T_r|}$$

509                                                                                                      (10)
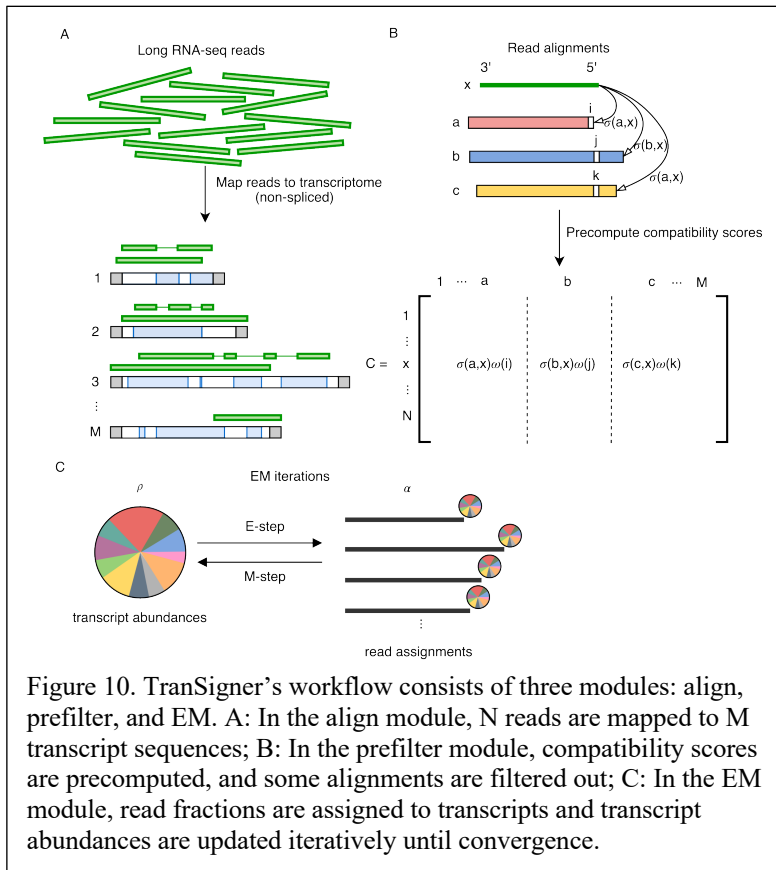
21

510    where $T_r$ is the set of transcripts that are compatible with $r$. The `drop` algorithm is called only right after

511    the first E-step calculation, and its purpose is to discard minimap2 alignments that are not robust. The

512    drop algorithm offers the potential to achieve a higher optimum compared to a naïve EM algorithm

513    (Pachter, 2011), which relies solely on the relative transcript abundances ($\rho$) in its E-step update. We also

514    allow users to increase this threshold (i.e., make it stricter) using the $-f$ parameter that'll increment $\tau_r$ by

515    a fraction of its own value as follows:

516    $$\tau_r' = \tau_r + (\tau_r * f)$$

517                                                                                                        (11)

518    where $f$ is a fractional value within the range [0, 1].

519

520    *Read assignment.* We can use the $\alpha$ values estimated by the EM algorithm to infer read assignments to

521    transcripts. Raw $\alpha$ values represent fractional read assignments, where a single read may be distributed

522    among multiple transcripts. These assignments might be challenging to interpret, as we assume each read

523    to originate from a single transcript. To increase the interpretability and usability of the $\alpha$ values, we

524    implemented the push algorithm (Supplementary Figure S5). This algorithm processes raw $\alpha$ values,

525    converting them into hard assignments where each read is assigned to exactly one transcript. The push

526    algorithm iterates through the reads and pairs each of them to the transcript with the highest read fraction

527    as shown by the corresponding $\alpha$ value. It then recomputes the relative transcript abundances based on

528    these hard assignments. These new $\alpha$ and $\rho$ values may deviate from their EM-derived ML estimates,

529    potentially resulting in reduced accuracy. We tested this using simulated data and observed only

530    negligible reductions in accuracy.

531

532    **Implementation**. TranSigner requires two inputs: a GTF file containing a reference gene annotation of

533    the target transcriptome and a FASTQ file containing long RNA-seq reads. The reference annotation can

534    be obtained from public sources such as RefSeq (O'Leary et al., 2016), GENCODE (Frankish et al.,

22

Figure 10. TranSigner's workflow consists of three modules: align, prefilter, and EM. A: In the align module, N reads are mapped to M transcript sequences; B: In the prefilter module, compatibility scores are precomputed, and some alignments are filtered out; C: In the EM module, read fractions are assigned to transcripts and transcript abundances are updated iteratively until convergence.

2019), or CHESS (Varabyou et al., 2023), or it can be derived from transcriptome assemblies produced by programs like StringTie2. The latter annotations have the advantage of including novel isoforms while restricting the annotated transcripts to only those found to be expressed in the analyzed sample.

As illustrated in Figure 10, TranSigner consists of three modules: align, prefilter, and em. In

548 the align module, input long reads are aligned to the target transcriptome using minimap2. The resulting

549 alignment file becomes the input for the next module. Next, in the prefilter module, TranSigner extracts

550 features such as the 3' and 5' end alignment positions and the ms alignment scores computed by

551 minimap2. These features are used to compute the compatibility score matrix between transcripts and

552 reads, as well as an index of the IDs of the transcripts found to be compatible with reads in the align

553 module, which represent a subset of the target transcriptome.

554

555 Finally, the EM module takes as inputs the compatibility score matrix and the target transcriptome index

556 from the prefilter module. It estimates the transcript coverage abundances using an expectation-

557 maximization (EM) algorithm. The EM algorithm converges when the total change in the relative

558 transcript abundances ($\rho$) is less than a specified threshold, by default set to 0.05. The drop algorithm,

559 described above and in Supplementary Figure S5, is implemented as a component of this module. It

560 allows users to use the `--drop` flag to remove low compatibility relations between reads and transcripts

23

561    immediately after the first E-step update. Read-to-transcript assignments (i.e., $\alpha$ estimates) and relative

562    transcript abundances (i.e., $\rho$ estimates) are outputted as TSV files at the end of the EM module. Users

563    also have the option to further process the assignments and output hard 1-to-1 assignments between reads

564    and transcripts for increased interpretability by specifying the `--push` flag, whose algorithm is described

565    in Supplementary Figure S5.

566

567    **Simulated data**. Three sets of Oxford Nanopore Technologies (ONT) direct RNA reads and two sets of

568    ONT cDNA reads were simulated using NanoSim (Gleeson et al., 2021). Expression levels were derived

569    from protein-coding and long non-coding transcripts located on the main chromosomes (i.e.,

570    chromosomes 1 – 22, X, and Y) of the GRCh38 genome, extracted from the RefSeq annotation (release

571    110). We supplied the NA12878 direct RNA and cDNA reads from Workman et al. to NanoSim's read

572    characterization module to first construct two separate read profiles, one for generating direct RNA and

573    the other for generating cDNA reads (Workman et al., 2019). We then estimated the transcript

574    abundances of the direct RNA and cDNA samples by aligning each sample to the GRCh38 genome using

575    minimap2 and providing the alignment results to salmon (Patro et al., 2017) in its alignment-based mode.

576    We used the RefSeq annotation as the target transcriptome. Salmon estimates were then used as input for

577    the NanoSim simulation module. For each direct RNA read set, we generated ~14 million ONT direct

578    RNA reads, and ~25 million for each cDNA read set (Supplementary Text 5).

579

580    **Spiked-in data**. We used an ONT direct-RNA dataset, which was released as part of the Singapore

581    Nanopore Expression Project (SG-NEx) (Chen et al., 2021). This dataset was sequenced from three

582    different human cell lines, HCT116, K562, and MCF7, and includes synthetic sequencing spike-in RNAs,

583    also known as sequin RNAs. We used the SG-NEx-provided genome, which includes the in silico

584    chromosome on which sequins are defined, to align these datasets. We also obtained the sequin transcripts

585    annotation, their raw abundances, and the sample-wise spike-in concentration (i.e., from the SG-NEx

586    AWS repository). To obtain sequin counts per million (CPM) levels, we followed the same method as in

24

587    Chen et al..The ground truth sequin CPM for a sequin transcript $x$ in a given sample $s$ was computed as

588    follows:

589
$$\mathrm{CPM}_x = \frac{a_x}{\sum_{t \in T} a_t} * c_s * 1000000$$

590                                                                                              (12)

591    where $a$ is the set of raw abundances provided by SG-Nex, $t$ iterates through the entire set of transcripts

592    to get the sum of all abundances, and $c_s$ is the spike-in concentration in sample $s$.

593

594    **Paired short- and long-read RNA-seq data.** For humans, we employed paired short- and long-read

595    RNA-seq data from the SG-NEx collection and long-read transcriptome profiling of human lung cancer

596    cell lines data sets. Short- and long-read datasets are considered paired if they were obtained by

597    sequencing the same biological sample. A subset of these samples included spike-in RNAs, and their

598    reads were aligned to augmented versions of the GRCh38 genome that also includes the sequin-

599    containing in silico chromosomes, provided by the original authors. All other samples (i.e., not spiked)

600    were aligned to the regular GRCh38 p13 genome.

601

602    The goal with paired RNA-seq data sets is to compute the correlation between the short- and long-read-

603    derived transcript abundance estimates. Long reads are first aligned to the GRCh38 genome using

604    minimap2 and the resulting alignments are provided to StringTie2 for a transcriptome assembly. Short

605    reads are then quantified on the long-read-derived StringTie2 transcripts using Salmon. Afterward, we ran

606    quantification-only methods – NanoCount and TranSigner – on the StringTie2 assembly to obtain long-

607    read-derived abundance estimates. We evaluated these tools' estimates based on their nonlinear

608    correlation with Salmon's short-read-derived estimates (see Supplementary Text 3 for the commands used

609    for short-read quantification). We repeated the same steps for two other organisms: *A. thaliana* and *M.*

610    *musculus*. None of the samples from these two species contained sequins, so all reads were aligned to

611    their respective reference genomes.

612

613   **Read assignments evaluation**. For simulated and sequin data, we can define the following values based

614   on the known origin transcript of each read:

- True positive (TP): a read is correctly assigned to its true origin.

- False positive (FP): a read is incorrectly assigned to a transcript that is not its true origin.

- False negative (FN): a read is not assigned to its true origin.

618   If a read is assigned to multiple transcripts without specifying the fraction allocated to each transcript,

619   then the read is evenly distributed among those transcripts, with these fractions contributing to TP and FP

620   values as appropriate. If the exact fraction of a read assigned to a transcript is provided, those fractions are

621   used instead.

622   For each sample, the recall value of a method for the read-to-transcript assignment is calculated as the

623   number of TPs divided by the total number of reads sequenced from that sample. The precision value is

624   computed as the number of TPs divided by the sum of TPs and FPs. F1 score is defined as 2 * precision *

625   recall / (precision + recall).

626

627   **Transcript abundance estimates evaluation**. By default, TranSigner outputs read counts and relative

628   transcript abundances as its quantification estimates. The read count of a transcript $t$ (denoted as $\mathrm{rc}_t$) is

629   the sum of all positive read fractions assigned to transcript $t$, while the relative transcript abundance of $t$

630   (denoted as $\rho_t$) is equal to $\mathrm{rc}_t$ normalized by the sum of all transcript read counts, ensuring that $\sum_{t \in T} \rho_t =$

631   1. Note that in a long-read RNA-seq experiment, each read counts as a transcript, making the sum of the

632   read counts equivalent to the total number of transcripts identified from the long-read data.

633

634   TranSigner's read count estimates can be converted to counts per million (CPM) estimates by calculating

635   $\mathrm{CPM}_t = \mathrm{rc}_t / l * 10^6$ where $t$ is a transcript and $l$ is the total number of reads (aligned and unaligned).

636   TranSigner also outputs read-to-transcript assignments where each read is assigned to one or more

26

637    transcripts. More precisely, TranSigner outputs a list of transcripts to which a read $r$ is assigned along

638    with the fraction of $r$ assigned to each transcript in that list, or the $\alpha$ estimates. These assignments can be

639    used to compute coverage estimates for transcripts as $\lambda_t = \left(\sum_{r \in R_t} \alpha_{rt} * l(r)\right) \Big/ l(t)$ where $\alpha_{rt}$ is the

640    fraction of $r$ assigned to transcript $t$, $R_t$ is the set of reads whose fractions were assigned to $t$, and $l$ is a

641    function that returns the length of a read or a transcript.

642

643    We performed both linear and nonlinear correlation analyses to evaluate the correlation between

644    estimated and ground truth values, each assessing different qualities of the read assignment and

645    quantification methods. While nonlinear correlation analysis, utilizing log-transformed read counts and

646    Spearman's correlation coefficient (SCC), evaluates monotonic trends in the data, linear correlation

647    analysis, utilizing Pearson's correlation coefficient (PCC), assesses a tool's accuracy in assigning all reads

648    to transcripts, valuing each read equally regardless of its source. It's worth noting that log transformation

649    is typically applied to reduce variance in gene expression values. However, log transformation may

650    compress differences in data points with large magnitudes, potentially diminishing the impact of errors in

651    assigning reads to high abundance transcripts.

652

653    **Evaluation of tools capable of transcriptome assembly**. We assessed the quality of assemblies

654    generated by StringTie2, Bambu, and FLAIR using the intron chain-level sensitivity and precision values

655    computed by GffCompare (Pertea & Pertea, 2020). We initially wanted to include ESPRESSO in this

656    comparison, but we were unable to run it as it took more than 24 hours to process a single sample

657    containing ~14 million reads.

658

659    We benchmarked each tool using random samples of the RefSeq annotation to observe how well the

660    completeness of the guides impacts the accuracy of the assembled transcriptome and the simulated ONT

661    data.  More precisely, we randomly sampled a percentage of the origin transcriptome, referring to the set

27

662    of transcripts from which a set of reads are simulated, to remove from RefSeq. The guides were sampled

663    to contain 21 different percentages between 0% and 100% of the origin transcriptome. For each

664    percentage, we independently sampled the guides three times, yielding 63 different guides per read set.

665    StringTie2, Bambu, and FLAIR were provided with the same guide annotations. Additionally, StringTie2

666    and Bambu were provided with the same minimap2 alignment results produced using the recommended

667    options for processing ONT RNA-seq data (`-x splice -uf -k14` for direct RNA reads and `-x`

668    `splice` for cDNA reads); FLAIR had its own align module. Unlike StringTie2 and FLAIR which output

669    an annotation containing only the identified expressed transcripts, Bambu outputs both expressed and

670    unexpressed transcripts in the guide annotation (see Supplementary Text 2). Therefore, for our

671    evaluations, we removed any transcript that was assigned a zero read count from Bambu's output.

672

673    **Declarations**

674

675    **Ethics approval and consent to participate.** Not applicable**.**

676

677    **Consent for publication.** All authors have consented for publication.

678

679    **Availability of data and materials.** The *A. thaliana* and *M. musculus* datasets are available from the

680    European Nucleotide Archive (ENA) under accession numbers PRJEB32782 and PRJEB27590. Specific

681    ENA sample accession IDs for each pair of short- and long-read data sets are made available in

682    Supplementary Table S11. The SG-NEx samples containing spike-in RNAs are available from GitHub,

683    ENA, and AWS open data registry. The long-read benchmarking on the human lung cancer cell lines data

684    sets are made available from Gene Expression Omnibus (GEO) under accession number GSE172421. The

685    values used to generate plots in this manuscript are made available as Supplementary Tables S1 ~ S11.

686    Supplementary Tables S0 contains the captions for each table. TranSigner is implemented in Python and

687    is publicly available at https://github.com/haydenji0731/transigner and is also archived in Zenodo at

688   https://doi.org/10.5281/zenodo.13334738. All code used to generate all figures (either in the main

689   manuscript or in the supplementary materials) and the scripts and data files (e.g., ground truths for

690   simulated and sequin data) used for benchmarking are available in Zenodo at

691   https://doi.org/10.5281/zenodo.13334733. The transcript abundances used for read simulation are also

692   available at the same address.

693

699

700   **Authors' contributions.** HJJ and MP designed the study. HJJ wrote the software and code used for

701   benchmarking. HJJ and MP evaluate analysis results and wrote / revised the manuscript. HJJ prepared all

702   (main and supplementary) figures. All authors read and approved the final manuscript.

703

707

708   **References**

709
710   Amaral, P., Carbonell-Sala, S., De La Vega, F. M., Faial, T., Frankish, A., Gingeras, T., Guigo, R.,
711       Harrow, J. L., Hatzigeorgiou, A. G., Johnson, R., Murphy, T. D., Pertea, M., Pruitt, K. D., Pujar,
712       S., Takahashi, H., Ulitsky, I., Varabyou, A., Wells, C. A., Yandell, M.,…Salzberg, S. L. (2023).
713       The status of the human gene catalogue. *Nature*, *622*(7981), 41-47.
714       https://doi.org/10.1038/s41586-023-06490-x
715   Amarasinghe, S. L., Su, S., Dong, X., Zappia, L., Ritchie, M. E., & Gouil, Q. (2020). Opportunities and
716       challenges in long-read sequencing data analysis. *Genome Biol*, *21*(1), 30.
717       https://doi.org/10.1186/s13059-020-1935-5

718 Benjamini, Y., & Speed, T. P. (2012). Summarizing and correcting the GC content bias in high-
719        throughput sequencing. *Nucleic acids research*, *40*(10), e72-e72.
720 Chen, Y., Davidson, N. M., Wan, Y. K., Patel, H., Yao, F., Low, H. M., Hendra, C., Watten, L., Sim, A.,
721        Sawyer, C., Iakovleva, V., Lee, P. L., Xin, L., Ng, H. E. V., Loo, J. M., Ong, X., Ng, H. Q. A.,
722        Wang, J., Koh, W. Q. C.,…consortium, S.-N. (2021). A systematic benchmark of Nanopore long
723        read RNA sequencing for transcript level analysis in human cell lines. *bioRxiv*,
724        2021.2004.2021.440736. https://doi.org/10.1101/2021.04.21.440736
725 Chen, Y., Sim, A., Wan, Y. K., Yeo, K., Lee, J. J. X., Ling, M. H., Love, M. I., & Göke, J. (2023).
726        Context-aware transcript quantification from long-read RNA-seq data with Bambu. *Nature*
727        *Methods*, *20*(8), 1187-1195. https://doi.org/10.1038/s41592-023-01908-w
728 Dong, X., Du, M. R. M., Gouil, Q., Tian, L., Jabbari, J. S., Bowden, R., Baldoni, P. L., Chen, Y., Smyth,
729        G. K., Amarasinghe, S. L., Law, C. W., & Ritchie, M. E. (2023). Benchmarking long-read RNA-
730        sequencing analysis tools using in silico mixtures. *Nature Methods*, *20*(11), 1810-1821.
731        https://doi.org/10.1038/s41592-023-02026-3
732 Frankish, A., Diekhans, M., Ferreira, A. M., Johnson, R., Jungreis, I., Loveland, J., Mudge, J. M., Sisu,
733        C., Wright, J., Armstrong, J., Barnes, I., Berry, A., Bignell, A., Carbonell Sala, S., Chrast, J.,
734        Cunningham, F., Di Domenico, T., Donaldson, S., Fiddes, I. T.,…Flicek, P. (2019). GENCODE
735        reference annotation for the human and mouse genomes. *Nucleic Acids Res*, *47*(D1), D766-d773.
736        https://doi.org/10.1093/nar/gky955
737 Gao, Y., Wang, F., Wang, R., Kutschera, E., Xu, Y., Xie, S., Wang, Y., Kadash-Edmondson, K. E., Lin,
738        L., & Xing, Y. (2023). ESPRESSO: robust discovery and quantification of transcript isoforms
739        from error-prone long-read RNA-seq data. *Science Advances*, *9*(3), eabq5072.
740 Gleeson, J., Leger, A., Prawer, Y. D. J., Lane, T. A., Harrison, P. J., Haerty, W., & Clark, M. B. (2021).
741        Accurate expression quantification from nanopore direct RNA sequencing with NanoCount.
742        *Nucleic acids research*, *50*(4), e19-e19. https://doi.org/10.1093/nar/gkab1129
743 Grünberger, F., Ferreira-Cerca, S., & Grohmann, D. (2022). Nanopore sequencing of RNA and cDNA
744        molecules in Escherichia coli. *Rna*, *28*(3), 400-417. https://doi.org/10.1261/rna.078937.121
745 Hansen, K. D., Brenner, S. E., & Dudoit, S. (2010). Biases in Illumina transcriptome sequencing caused
746        by random hexamer priming. *Nucleic acids research*, *38*(12), e131-e131.
747 Kovaka, S., Zimin, A. V., Pertea, G. M., Razaghi, R., Salzberg, S. L., & Pertea, M. (2019). Transcriptome
748        assembly from long-read RNA-seq alignments with StringTie2. *Genome biology*, *20*, 1-13.
749 Li, B., Ruotti, V., Stewart, R. M., Thomson, J. A., & Dewey, C. N. (2009). RNA-Seq gene expression
750        estimation with read mapping uncertainty. *Bioinformatics*, *26*(4), 493-500.
751        https://doi.org/10.1093/bioinformatics/btp692
752 Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, *34*(18), 3094-
753        3100. https://doi.org/10.1093/bioinformatics/bty191
754 Li, H. (2021). New strategies to improve minimap2 alignment accuracy. *Bioinformatics*, *37*(23), 4572-
755        4574. https://doi.org/10.1093/bioinformatics/btab705
756 O'Leary, N. A., Wright, M. W., Brister, J. R., Ciufo, S., Haddad, D., McVeigh, R., Rajput, B., Robbertse,
757        B., Smith-White, B., Ako-Adjei, D., Astashyn, A., Badretdin, A., Bao, Y., Blinkova, O., Brover,
758        V., Chetvernin, V., Choi, J., Cox, E., Ermolaeva, O.,…Pruitt, K. D. (2016). Reference sequence
759        (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation.
760        *Nucleic Acids Res*, *44*(D1), D733-745. https://doi.org/10.1093/nar/gkv1189
761 Pachter, L. (2011). Models for transcript quantification from RNA-Seq. *arXiv preprint arXiv:1104.3889*.
762 Pardo-Palacios, F. J., Wang, D., Reese, F., Diekhans, M., Carbonell-Sala, S., Williams, B., Loveland, J.
763        E., María, M. D., Adams, M. S., Balderrama-Gutierrez, G., Behera, A. K., Gonzalez, J. M., Hunt,
764        T., Lagarde, J., Liang, C. E., Li, H., Meade, M. J., Amador, D. A. M., Prjibelski, A. D.,…Brooks,
765        A. N. (2023). Systematic assessment of long-read RNA-seq methods for transcript identification
766        and quantification. *bioRxiv*, 2023.2007.2025.550582. https://doi.org/10.1101/2023.07.25.550582

767  Patro, R., Duggal, G., Love, M. I., Irizarry, R. A., & Kingsford, C. (2017). Salmon provides fast and bias-
768      aware quantification of transcript expression. *Nature Methods*, *14*(4), 417-419.
769      https://doi.org/10.1038/nmeth.4197
770  Pertea, G., & Pertea, M. (2020). GFF Utilities: GffRead and GffCompare [version 2; peer review: 3
771      approved]. *F1000Research*, *9*(304). https://doi.org/10.12688/f1000research.23297.2
772  Tang, A. D., Soulette, C. M., van Baren, M. J., Hart, K., Hrabeta-Robinson, E., Wu, C. J., & Brooks, A.
773      N. (2020). Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic
774      leukemia reveals downregulation of retained introns. *Nature communications*, *11*(1), 1438.
775  Varabyou, A., Sommer, M. J., Erdogdu, B., Shinder, I., Minkin, I., Chao, K.-H., Park, S., Heinz, J.,
776      Pockrandt, C., Shumate, A., Rincon, N., Puiu, D., Steinegger, M., Salzberg, S. L., & Pertea, M.
777      (2023). CHESS 3: an improved, comprehensive catalog of human genes and transcripts based on
778      large-scale expression data, phylogenetic analysis, and protein structure. *Genome biology*, *24*(1),
779      249. https://doi.org/10.1186/s13059-023-03088-4
780  Wongsurawat, T., Jenjaroenpun, P., Wanchai, V., & Nookaew, I. (2022). Native RNA or cDNA
781      Sequencing for Transcriptomic Analysis: A Case Study on Saccharomyces cerevisiae [Original
782      Research]. *Frontiers in Bioengineering and Biotechnology*, *10*.
783      https://doi.org/10.3389/fbioe.2022.842299
784  Workman, R. E., Tang, A. D., Tang, P. S., Jain, M., Tyson, J. R., Razaghi, R., Zuzarte, P. C., Gilpatrick,
785      T., Payne, A., Quick, J., Sadowski, N., Holmes, N., de Jesus, J. G., Jones, K. L., Soulette, C. M.,
786      Snutch, T. P., Loman, N., Paten, B., Loose, M.,…Timp, W. (2019). Nanopore native RNA
787      sequencing of a human poly(A) transcriptome. *Nat Methods*, *16*(12), 1297-1305.
788      https://doi.org/10.1038/s41592-019-0617-2
789
790