


METHOD

Open Access

LIQA: long-read isoform quantification and analysis



Yu Hu¹, Li Fang¹, Xuelian Chen², Jiang F. Zhong², Mingyao Li³ and Kai Wang^{1,4*} 

* Correspondence: wangk@email.chop.edu

¹Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children's Hospital of Philadelphia, Philadelphia, PA 19104, USA

⁴Department of Pathology and Laboratory Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA

Full list of author information is available at the end of the article

Abstract

Long-read RNA sequencing (RNA-seq) technologies can sequence full-length transcripts, facilitating the exploration of isoform-specific gene expression over short-read RNA-seq. We present LIQA to quantify isoform expression and detect differential alternative splicing (DAS) events using long-read direct mRNA sequencing or cDNA sequencing data. LIQA incorporates base pair quality score and isoform-specific read length information in a survival model to assign different weights across reads, and uses an expectation-maximization algorithm for parameter estimation. We apply LIQA to long-read RNA-seq data from the Universal Human Reference, acute myeloid leukemia, and esophageal squamous epithelial cells and demonstrate its high accuracy in profiling alternative splicing events.

Introduction

RNA splicing is a major mechanism for generating transcriptomic variations, and misregulation of splicing is associated with a large array of human diseases caused by hereditary and somatic mutations [1–5]. Over the past decade, RNA sequencing (RNA-seq) has revolutionized transcriptomics studies and facilitated the characterization and understanding of transcriptomic variations in an unbiased fashion. With RNA-seq, we can quantitatively measure isoform-specific gene expression, discover novel and unique transcript isoform signature, and detect differential alternative splicing (DAS) events [6–8]. Until now, short-read RNA-seq has been the method of choice for transcriptomics studies due to its high coverage and single-nucleotide resolution [8]. However, due to limited read length, it is difficult to accurately characterize transcripts using short reads, as 81% of isoforms have length greater than 500 bp in the GENCODE annotation (median = 1543 bp and mean = 2121 bp). This fragmented sequencing of the RNA/cDNA molecules results in biases and has become a barrier for short reads to be correctly mapped to the reference genome, which is crucial for gene or isoform expression estimation and novel or unique isoform detection. As a consequence, transcriptome profiling using short-read RNA-seq is highly biased by read coverage heterogeneity across isoforms. To tackle these challenges, a number of computational tools, including RSEM [9], eXpress [10], TIGAR 2 [11], Salmon [12], Sailfish



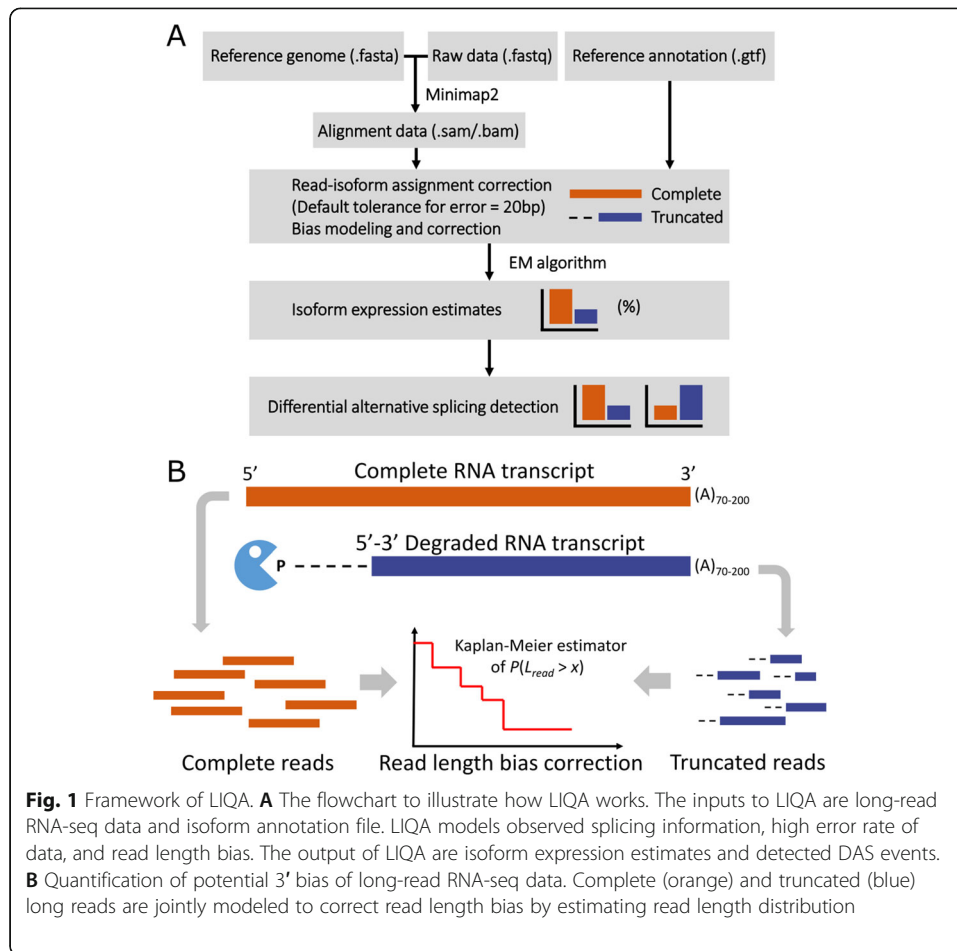
© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

[13], Kallisto [14], Cufflinks [15], CEM [16], PennSeq [17], IsoEM [18], and RD [19], have been developed to quantify isoform expression from short-read RNA-seq data, but different bias correction algorithms can result in conflicting estimates [17]. Overall, quantifying isoform expression using fragmented short reads remains challenging, especially at complex gene loci [20, 21].

In recent years, long-read RNA sequencing has gained popularity due to its potential to overcome the above-mentioned issues when compared to short-read RNA-seq [22, 23]. Previous studies have utilized both single-molecule long-read PacBio Iso-Seq and synthetic long-read MOLECULO methods [24–27]. For Oxford Nanopore sequencing, there are two types of RNA-seq technologies: direct mRNA sequencing and cDNA sequencing. Recently, the Oxford Nanopore Technologies (ONT) MinION has been used to analyze both full-length cDNA samples and mRNA samples derived from tissue cells [28]. Nanopore sequencing is able to generate reads as long as 2 Mbp, which allows a large portion or the entire mRNA or cDNA molecule to be sequenced. Compared to short reads, this advantage of long reads greatly facilitates rare isoform discovery, isoform expression quantification, and DAS event detection.

However, there are still a few unique challenges to analyze long-read RNA-seq data because existing methods developed for Illumina short-read RNA-seq do not have optimal performance when directly used on long-read RNA-seq. This is because parametric bias correction of short-read approaches requires high read coverage and isoform-read assignment is not robust to small range misalignment from long-read data [16, 18, 19, 29]. Methods designed specifically for isoform expression estimation in long-read RNA-seq have only emerged recently. For example, Byrne et al. [30] demonstrated the feasibility of quantifying complex isoform expression using Nanopore RNA-seq data. Tang et al. [31] characterized mutated gene *SF3B1* at isoform level in chronic lymphocytic leukemia cells by leveraging full-length transcript sequencing data generated by Nanopore. While long-read RNA-seq has great potential, the isoform quantification accuracy is still constrained by high error rates and sequencing biases [32], which has yet to be thoroughly accounted for. Specifically, high sequencing error rates (~ 15%) of Nanopore data can result in misalignment of sequencing reads, but current methods assume equal weight for each single-molecule read without accounting for error rate differences when estimating isoform expression. This may complicate isoform usage quantification. In addition, potential read coverage biases are not explicitly taken into account by existing long-read transcriptomic tools [32]. In Nanopore direct mRNA sequencing protocol, pore block and fragmentation can result in truncated reads, leading to biased coverage toward the 3' end of a transcript [32]. These biases are also shown in data sequenced from cDNA. In the presence of such biases, the accuracy of isoform expression quantification inference can be severely affected, leading to over estimation of expression for isoforms with short length.

In this article, we present LIQA, a statistical method that allows each read to have its own weight when quantifying isoform expression. Rather than counting single-molecule reads equally, we give a different weight to each read to account for read-specific error rate and alignment bias at the gene (Fig. 1). To evaluate the performance of LIQA, we simulated long data with known ground truth and also sequenced two real samples using Oxford Nanopore sequencing. Our results demonstrate that LIQA is an



accurate approach for isoform expression quantification accounting for read coverage bias and high error rate of long-read data.

Results

Overview of LIQA

Figure 1 shows the workflow of LIQA and highlights the read length bias correction step. LIQA requires aligned long-read RNA-seq files in BAM or SAM format and isoform annotation file as input. For estimation steps, LIQA first feeds read alignment information to a complete likelihood function and corrects biases for each long read by accounting for quality score and read coverage bias. Second, given that isoform origins are unobserved for some reads, an expectation maximization (EM) algorithm is utilized to achieve the optimal solution of isoform relative abundance estimation. The output values of LIQA are isoform expression estimates. Moreover, LIQA can further detect DAS events given estimated isoform expression values.

To evaluate the performance of LIQA, we compared it with existing long-read based quantification algorithms, including FLAIR [31], Mandalorion [30], TALON [33], and the Oxford Nanopore Pipeline (ONP; <https://github.com/nanoporetech/pipeline-transcriptome-de>). These methods use long-read RNA-seq data to detect novel

isoforms and quantify transcript expression by counting the number of reads, which give equal weight for each read. To make the comparisons fair, we ran LIQA, FLAIR, TALON, Mandalorion, and ONP in quantification mode only with isoform annotation information provided by GENCODE. We benchmarked the performance of each method on both simulated and real data. In addition, we simulated more data to evaluate the performance of LIQA in detecting DAS events between conditions.

Nanopore RNA-seq data simulation

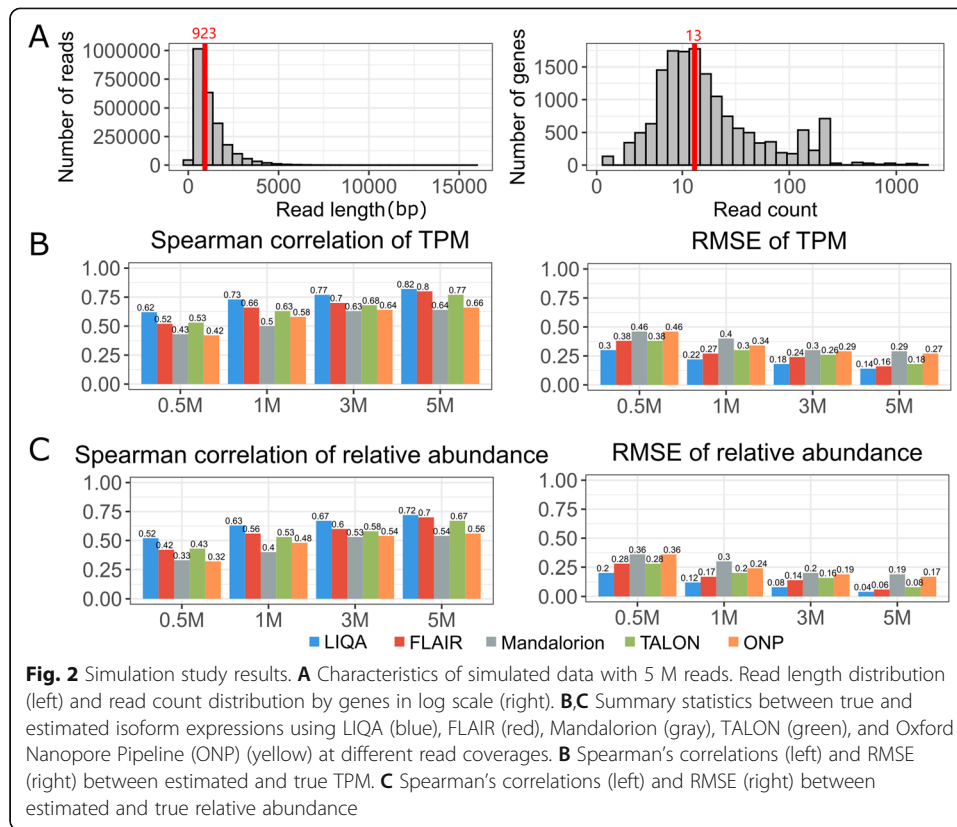
We conducted a simulation study to evaluate the performance of LIQA and compared it with other state-of-the-art algorithms for isoform expression estimation and DAS detection based on GENCODE v24 annotation. To simulate a realistic dataset with known ground truth, we used NanoSim [34] to generate the ONT RNA-seq data. NanoSim is a fast and scalable read simulator that captures the technology-specific features of ONT data and allows for adjustment upon improvement of Nanopore sequencing technology. This simulator first characterizes Nanopore reads and models features of the library preparation protocols in silico for read simulation. The human genome sequence (GRCh38), transcriptome sequence, and GTF annotation file were downloaded from GENCODE. To make the simulated data close to real studies, we assigned abundance values for each isoform obtained from a real human eye RNA-seq dataset. Using NanoSim, we generated 5 million (5 M) Nanopore reads. To evaluate the impact of sequencing depth on isoform expression quantification, we down-sampled 3 million (3 M), 1 million (1 M), and 0.5 million (0.5 M) reads for the simulated data, respectively. These reads were aligned to the reference genome using minimap2 [35]. Then, we selected genes with 2 or more isoforms (67.2%) to evaluate the performance of LIQA in isoform expression quantification. For each isoform, we compared it with Mandalorion, FLAIR, TALON, and ONP. All methods were run with the same set of simulated aligned data in BAM format as input.

The characteristics of the simulated data are shown in Fig. 2A and Additional file 1: Fig. S1(A). The median lengths of ONT reads in the 0.5 M, 1 M, 3 M, and 5 M datasets are 896, 922, 1010, and 923 base pairs, respectively. Among the evaluated genes with multiple isoforms (67.2%) based on GENCODE annotation, 13% have two isoforms, 14% have three isoforms, and 73% have four or more isoforms. The simulated isoforms have a wide range of relative abundance (interquartile range = (0.002, 0.75), median = 0.041). In addition, by training the statistical model of NanoSim with a real long-read RNA-seq dataset, the coverage plots of the simulated data capture the features of real ONT RNA-seq data, demonstrating 3' bias (Additional file 1: Fig. S1(B)). These simulated data thus provide an ideal basis to evaluate the performance of LIQA as the ground truth is known.

Gene or isoform expression quantification accuracy

For each simulated dataset, we computed a set of measurements to evaluate the estimation accuracy of each method. First, we measured the similarity between the estimated isoform relative abundance and the ground truth by calculating Spearman's correlation. Second, we measured the estimation accuracy by calculating the root mean squared

error (RMSE), defined as $\sqrt{\frac{\sum_g \sum_i (\hat{\theta}_{g,i} - \theta_{g,i})^2}{n}}$, where the summation is taken over all genes



and all isoforms within each gene and n is the total number of isoforms across all genes. Both statistics were computed at three levels: global gene expression, global isoform expression, and within-gene isoform relative abundances.

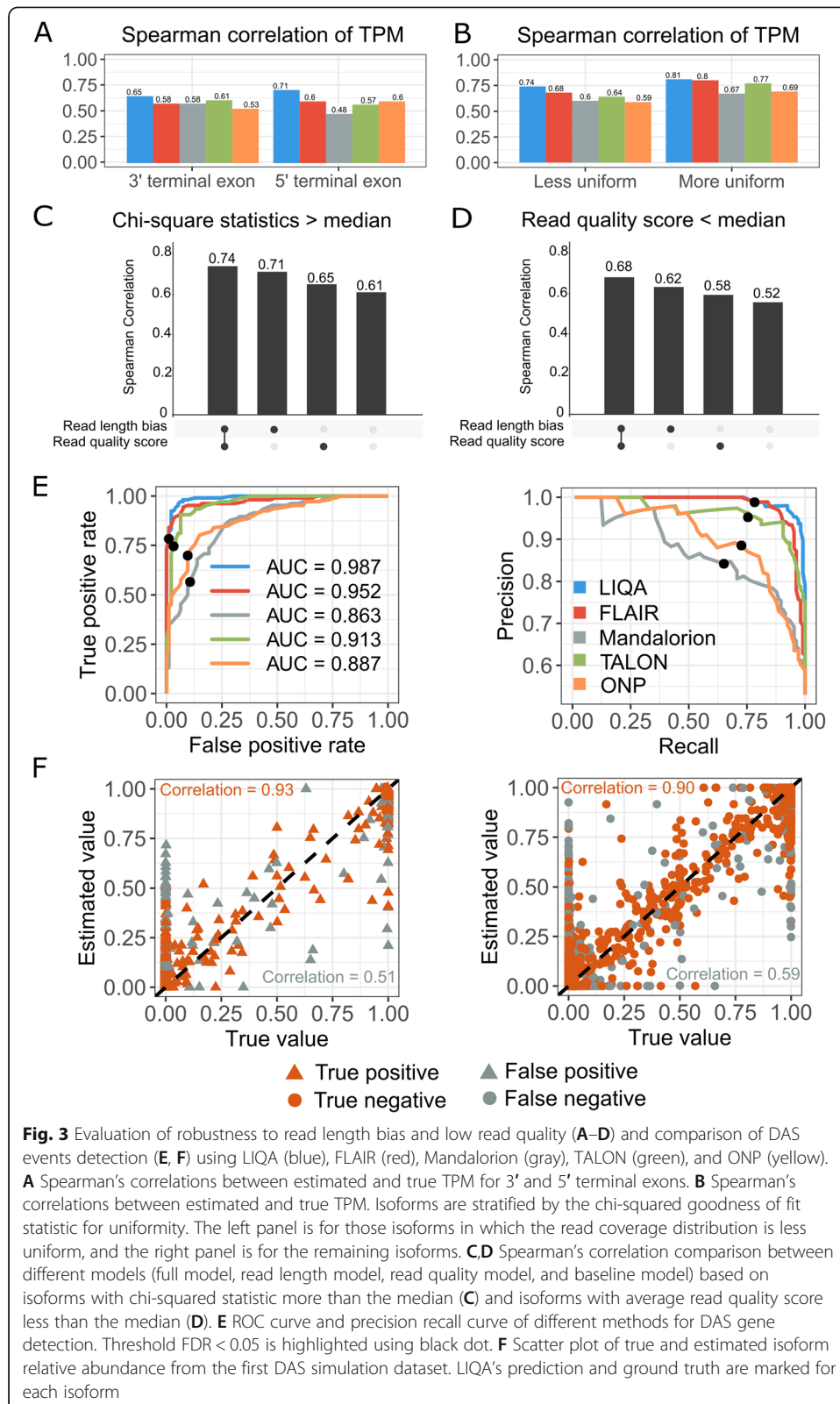
Figure 2B,C (Additional file 1: Fig. S2) shows the summary statistics between estimated and true values of isoform expression (global isoform expression and isoform relative abundances) at different read coverages. Spearman's correlation and RMSE were calculated for all five methods. LIQA has higher Spearman's correlation than other methods for simulated datasets with low sequencing depth (0.5 M) (Additional file 1: Table S2). For simulated data with high sequencing depth (3 M, 5 M), Spearman's correlation differences between LIQA, FLAIR, and TALON are not significant (Additional file 1: Table S2). Figure 2C gives summary statistics of relative abundance estimates for the five methods. For relative abundance estimation, LIQA outperforms FLAIR and TALON with 6.6% and 7% lower RMSE on average, respectively. Comparison results at gene level reveal a similar pattern (Additional file 1: Fig. S4, Table S5). The improved performance of LIQA is likely due to its use of the EM algorithm, which assigns unequal weight to each read to better account for mapping uncertainty and read mapping bias (Fig. 3C,D and Additional file 1: Table S6, Table S7). In contrast, FLAIR, TALON, and Mandalorion provide discrete estimations by directly counting the number of reads aligned to each corresponding gene or isoform. Due to the limited read coverage of ONT RNA-seq data, it is not surprising that they yield lower estimation accuracy.

To evaluate the robustness of LIQA to a more complex isoform annotation, we analyzed 5 M simulation dataset based on GENCODE v37 annotation. For major use

isoforms, as shown in Additional file 1: Fig. S4(D), LIQA still yields 7% higher Spearman correlation of TPM and relative abundance estimates than second best approach FLAIR. For over-annotation, we simulated RNA-seq reads based on 66% of the GENCODE v37 annotation. We then analyzed the simulated data with various methods using 100% of the GENCODE annotation, corresponding to 50% more of the true annotation. Additional file 1: Fig. S4(E) shows the Spearman correlation results of over-annotation. We find that the quantification accuracy of LIQA is nearly unchanged (1% less). For under-annotation, we simulated RNA-seq reads based on 100% of the GENCODE v37 annotation. We then analyzed the simulated data using 50% of the GENCODE annotation, corresponding to 50% less of the true annotation. Additional file 1: Fig. S4(E) shows Spearman's correlation results of under-annotation. The estimation accuracy is 10% lower when 50% less of the true annotation was used in the analysis.

Next, we evaluated the robustness of LIQA to 3' read coverage bias (Fig. 3 and Additional file 1: Fig. S3). First, we compared the accuracy statistics for 5' terminal exon and 3' terminal exon of each isoform. Isoform expression with non-uniform read coverage is more challenging to estimate because the 5' end is less likely to be covered by sequencing reads compared to 3' end. Figure 3A shows the comparison of Spearman's correlation for five methods with 0.5 M read coverage. LIQA is more accurate than the other four methods at 5' terminal exon, especially when sequencing depth is low (Additional file 1: Table S4). Spearman's correlation coefficient of LIQA is 11% higher than the second best performing method FLAIR for 5' terminal exons, while only 6% higher for 3' exons. This improved performance of LIQA in terminal exon quantification is also demonstrated by RMSE values. LIQA has 8–15% improvement of RMSE values compared to other methods. Second, we considered the chi-square statistics that measures the goodness of fit of coverage uniformity. Then, we divided the isoforms into two categories based on median of the corresponding measure (chi-square statistic > median, chi-square statistic < median) and summarized Spearman's correlation coefficient and RMSE for each group of isoforms. For isoforms with more uniform read coverage, Spearman's correlations of LIQA and FLAIR are close. However, despite reduced Spearman's correlation value, LIQA is more accurate than the other four methods for isoforms with less uniform read coverage (chi-square statistic > median) (Additional file 1: Table S4). The improvement of LIQA compared to FLAIR is 5% higher for these isoforms. This is likely because LIQA models potential truncated reads which result in 3' coverage bias when quantifying isoform expression.

Moreover, we assessed the impact of modeling read length bias and read quality score on the accuracy of isoform expression estimation. Figure 3C,D shows the comparison of isoform estimation accuracy using different models. For isoforms with less uniform read coverage (chi-square statistics > median), model with read length bias correction has 9% (full model vs read quality model) and 10% (read length model vs baseline model) higher Spearman's correlation. For genes with less average read quality, model with read quality score has 6% higher Spearman's correlation. Overall, isoform estimation accuracy drops noticeably when using baseline model (Additional file 1: Table S7). This comparison demonstrates the advantage of LIQA in handling read length bias and 3' bias correction over other approaches.



Differential alternative splicing (DAS) detection

Next, we evaluated the performance of LIQA in DAS detection. More ONT RNA-seq data across multiple samples (10 cases and 10 controls) were simulated for 10 times. NanoSim generated 3 million reads based on the GENCODE annotation per sample. To make true DAS events more realistic, we sampled relative abundances of isoforms from a Dirichlet distribution with mean and variance parameters estimated from a human eye RNA-seq dataset. Similarly, these simulated data were mapped to the hg38 human reference genome using minimap2. Isoform expression and usage difference between conditions were quantified using LIQA, FLAIR, TALON, Mandalorion, and ONP, respectively. We first compared the performance of DAS detection between these methods using three summary statistics. After FDR control, we measured the recalls (power) of our method by calculating the proportion of correctly predicted DAS events among true DAS events. Second, we obtained precisions by calculating the proportion of correctly predicted DAS events among predicted DAS events. Additionally, F1 score ($F1 \text{ score} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$) was summarized to average the precision and recall values. As shown in Additional file 1: Fig. S5(B), LIQA, FLAIR, and TALON are more powerful than others for all three evaluation metrics. This is not surprising because Mandalorion and ONP have lower accuracy in isoform expression estimation. For recall value, FLAIR (mean = 0.809, SD = 0.041) gives better and more consistent performance across 10 simulations than LIQA (mean = 0.776, SD = 0.058). However, in terms of precision value, LIQA (mean = 0.915, SD = 0.043) yields less false positives than FLAIR (mean = 0.884, SD = 0.051). LIQA, FLAIR, and TALON had similar performance in detecting DAS events based on F1 score. Furthermore, we generated ROC curve and precision recall curve to compare the performance between methods at different FDR thresholds (Fig. 3E and Additional file 1: Fig. S6, S7). As shown in Fig. 3E, LIQA achieved AUC = 0.94 after FDR control (FDR < 0.05). Given FDR threshold equals to 0.05, LIQA gave the best performance with precision = 0.98 and recall = 0.78. Compared to LIQA, the second best performing method FLAIR yields 0.1%, 0.3%, and 3.5% less in precision, recall, and AUC respectively. In addition, we examined isoform relative abundance estimation accuracy from correct and incorrect detected DAS genes by LIQA (Fig. 3F). After FDR control, we identified that 537 out of 2465 genes are significantly differential spliced, which 431 are true positives and 1836 are true negatives. For these correctly predicted genes, true isoform relative abundance is highly correlated with LIQA's estimates (Spearman's correlation = 0.91). For false positive and negative genes, Spearman's correlation is 38% lower compared to true positive and negative. This is not surprising because accurate estimation of isoform expression level leverages the power of regression model in detecting DAS events.

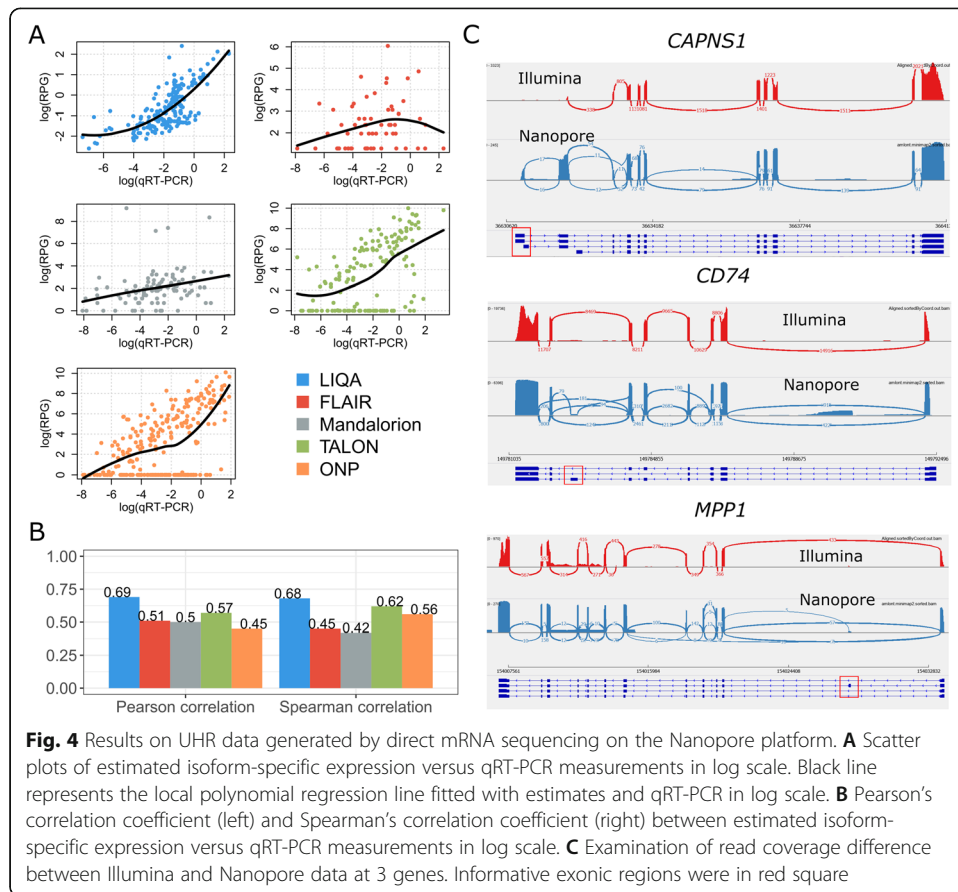
Application to the Universal Human Reference (UHR) RNA-seq data

As NanoSim generates ONT RNA-seq data based on trained parametric statistical model, we recognized that simulated data is hardly a full representation of reality. To evaluate the performance of LIQA in a real setting, we sequenced the Universal Human Reference sample with Nanopore Direct mRNA sequencing (Additional file 1: Fig. S12). Then, the resulting ONT-RNA-seq data were analyzed using all five long-read-based methods (LIQA, FLAIR, Mandalorion, TALON, ONP). As quantitative real-time

polymerase chain reaction (qRT-PCR) is considered as the most reliable technology for measuring true isoform abundance, we downloaded the qRT-PCR measurements from the MAQC project under Gene Expression Omnibus with accession number GSE5350. As part of the MAQC project, the expression levels of 894 isoforms were measured by TaqMan Gene Expression Assay based qRT-PCR. Additionally, we downloaded the UHR short-read RNA-seq data generated using the Illumina platform. This dataset was analyzed using Cufflinks [15], CEM [16], Salmon [36], and Kallisto [14] to compare the performance in isoform quantification between long reads and short reads. Specifically, we mapped ONT and Illumina sequenced reads to the reference genome using Minimap2 [35] and STAR [37], respectively, and applied each quantification method to the RNA-seq data. qRT-PCR measurements were treated as gold standard to compare the performance across methods. We note that 563 of the 894 transcripts with qRT-PCR measurements are from genes with a single isoform. Estimation results from these genes were served as positive controls (Additional file 1: Fig. S8(A)) because estimating isoform-specific expression for these single-transcript genes is trivial. To compare the performance across different methods, we considered those transcripts that are derived from genes with two or more isoforms.

To assess the accuracy between estimates and qRT-PCR measurements, we summarized similarity metrics (Spearman's correlation and Pearson's correlation) of the isoform abundance values in log scale. As shown in Fig. 4A,B, the estimation accuracy of all methods is lower than simulated data, especially for those transcripts with qRT-PCR measures close to 0. Nevertheless, we observed consistent results in terms of relative performance of different methods with simulation data. LIQA is more accurate than other methods with stronger linear relationship between logarithm estimates and qRT-PCR measurements. However, many of the lowly to moderately expressed isoforms were underestimated using the other methods with their TPM values being compacted toward 0. For ONT data, Spearman's correlation of LIQA is 0.68, whereas the corresponding values from second best method TALON is 0.62. For Illumina data, Cufflinks seems to correlate with the qRT-PCR measurements better than others (Additional file 1: Fig. S8(B)). The main reason for the better performance of LIQA is likely due to quantifying isoform expression by accounting for isoform length bias and base quality scores. Moreover, we randomly selected 3 genes and generated sashimi plots in Fig. 4C to show the read coverage difference between direct mRNA sequencing and Illumina data. Overall, read distribution of long-read data is less heterogeneous than short-read. Specifically, for gene *CAPNS1*, there is clearly severe 5' degradation in Illumina data, but with full length and more even coverage across the transcripts for long-read data. Terminal exons at 5' end in red square are crucial informative regions for splicing analysis, which enable us to differentiate read origin from different isoforms. As shown in Fig. 4C, these exonic regions were captured by Nanopore reads but missed by Illumina reads, which significantly facilitates isoform expression quantification using long-read RNA-seq data. Similarly, sashimi coverage plots of other two genes showed the same pattern, which demonstrates the advantage of long-read data over short-read in alternative splicing study.

Moreover, we conducted additional analysis of another long-read data on UHR with much higher coverage (5.6 million reads), generated on the PacBio sequencing platform [38]. As shown in Additional file 1: Fig. S8(E), Pearson and Spearman's correlations for



each method are generally improved with increasing sequencing depth compared to our Nanopore-based UHR dataset (Fig. 4B). For example, Spearman's correlation of FLAIR is increased by 0.28 (from 0.45 to 0.73), whereas the corresponding values of LIQA are increased by 0.11. Nevertheless, LIQA still has the best performance among all methods. Based on this real dataset with increasing sequencing depth, we found that LIQA is more robust to low read coverage compared to FLAIR, which performs well when sequencing depth is high. These observations from these two real UHR datasets are consistent with the simulation-based datasets with different sequencing depths (0.5 M, 1 M, 3 M, 5 M).

Application to Nanopore cDNA sequencing data on a patient with acute myeloid leukemia

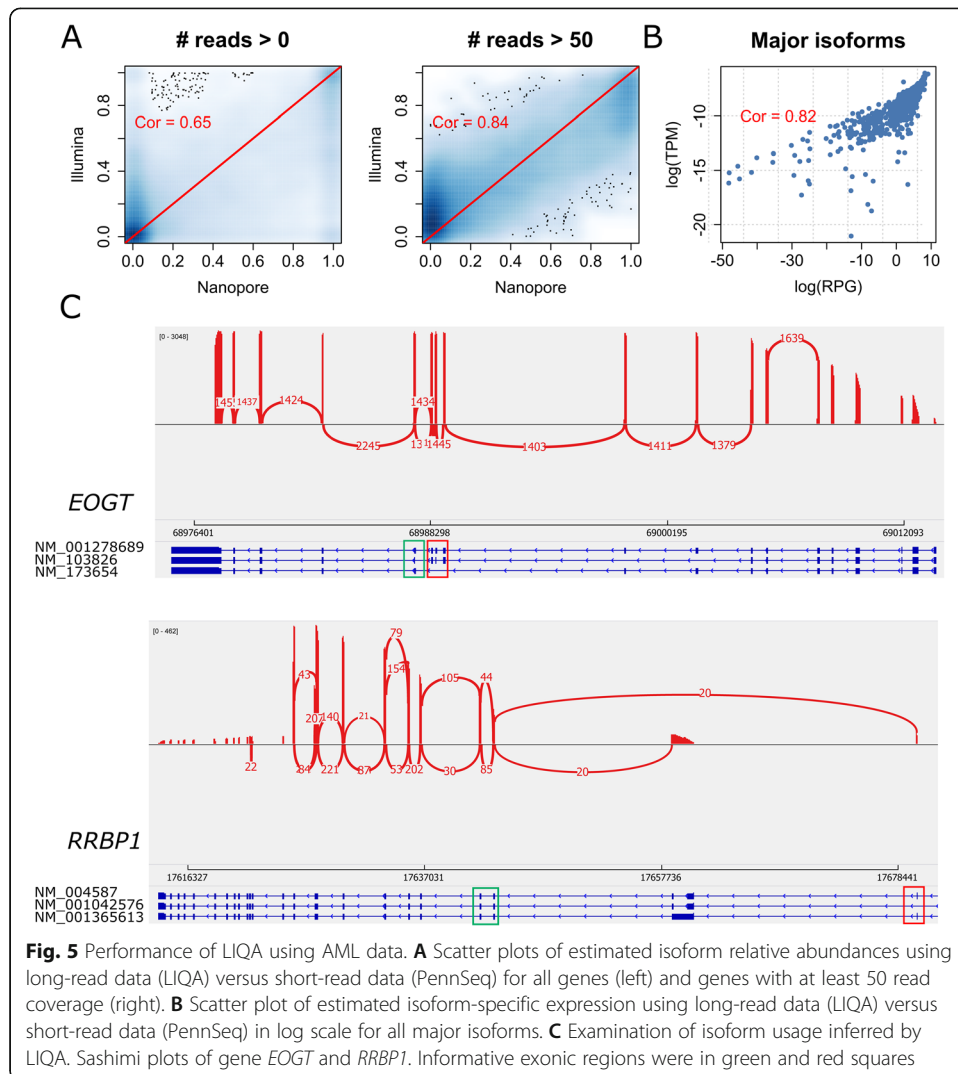
AML is a type of blood cancer where abnormal myeloblasts are made by bone marrow [39]. In this study, we sequenced peripheral blood from an acute myeloid leukemia (AML) patient using GridION Nanopore sequencer with Guppy basecalling (<https://denbi-nanopore-training-course.readthedocs.io/en/latest/basecalling/basecalling.html>). In total, we generated 8,061,683 long reads with 6.6 GB bases (Additional file 1: Fig. S13). We aligned the data against a reference genome (hg38) using minimap2 [35], and 63% long reads (73% bases) are mapped. Then, we analyzed this ONT RNA-seq data with LIQA for genes with at least two isoforms.

We considered two ways to benchmark the performance of LIQA. First, we used PennSeq to analyze short-read sequencing data for the same AML sample and treated the estimates as gold standard. This dataset included 440 M short read with 150 bp in length. Figure 5A shows the scatter plots of isoform relative abundance estimates between long- and short-read data. Spearman's correlation coefficients were calculated. We found that correlation was improved significantly for genes with at least 50 reads compared to all genes without filtration. Then, we examined the major isoforms (with the highest expression level in a gene) inferred by LIQA. As shown in Fig. 5B, long-read and short-read shared consistent estimates for the major isoforms. This is not surprising because major isoforms were more likely to be sequenced, leading to higher read coverage at unique exonic regions. Second, we visually examined the read coverage plots at unique exonic regions with at least 100 reads to benchmark the performance of LIQA. We generated sashimi plots for two randomly selected genes, *EOGT* and *RRBP1* (Fig. 5C). For gene *EOGT*, the read coverage ratio between exons in red and green squares suggests that isoforms NM_103826 and NM_001278689 expressed much higher than NM_173654. This is consistent with LIQA's estimates, with relative abundance of NM_173654 less than 0.01. A similar pattern is observed for gene *RRBP1*, where isoform NM_004587 (relative abundance estimates = 0.68) is the major isoform. Results from this AML data demonstrate the robust performance of LIQA to 3' coverage biases.

Application to PacBio data on esophageal squamous epithelial cell (ESCC)

Next, we evaluated the performance of LIQA in differential alternative splicing (DAS) detection using an RNA-seq dataset generated from esophageal squamous epithelial cell (ESCC) [40]. This dataset includes PacBio SMRT reads generated from normal immortalized and cancerous esophageal squamous epithelial cell lines. The RNA-seq data were downloaded from Gene Expression Omnibus (PRJNA515570). We applied LIQA to detect differential isoform usage between normal-like and cancer cells. Known splicing differences in existing studies were treated as ground truth to evaluate LIQA's performance in characterizing isoform usage across samples. In addition, short-read data from these two samples were sequenced using the Illumina platform, which allows us to compare the consistency and accuracy of DAS detection between long-read and short-read data.

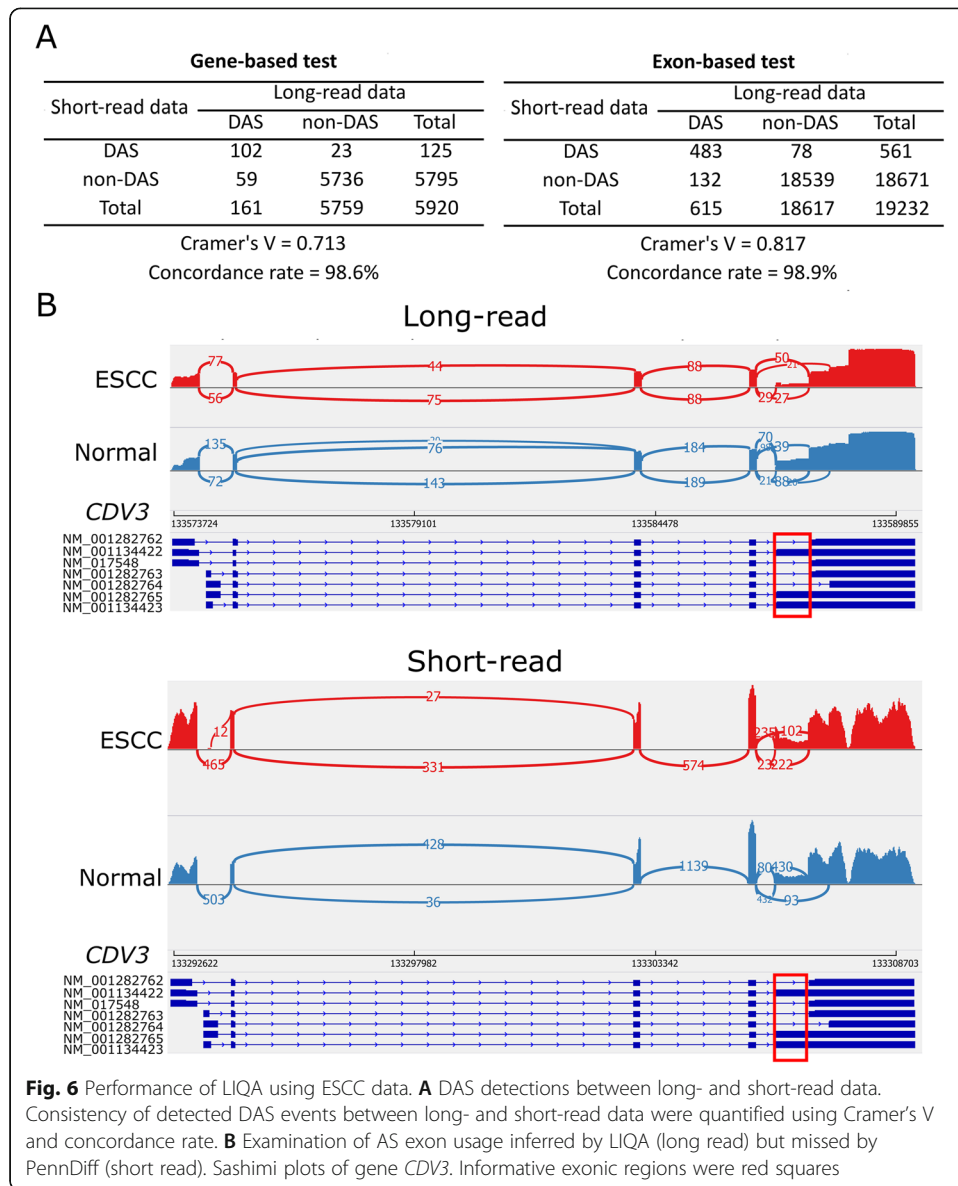
Employing LIQA and PennDiff [41], PacBio, and Illumina data were analyzed to detect DAS events, which are classified into different types, such as skipped exon (SE), alternative 5' splice site (A5SS), alternative 3' splice site (A3SS), mutually exclusive exon (MXE), and retained intron (RI). Our results showed that SE is the most frequent type of event among detected DAS between normal-like and cancerous cells, followed by RI, A5SS, and A3SS. MXE is the most infrequent splicing type. As shown in Fig. 6A, detected DAS events by long- and short-read share strong association at both exon and gene level (Cramer's $V > 0.5$). Also, the concordance rate between long- and short-read data is greater than 98%. Compared to short-read data, long-read data shows preference in detecting more differential splicing events at both exon and gene level. This is not surprising because read coverage heterogeneity, which might bias DAS detection, is alleviated in long-read data by capturing full-length transcript in each read.



The expression of alternatively spliced isoforms from gene *CDV3* shows difference in cancerous ESCC cells compared to non-cancerous [42, 43]. Figure 6B provides the sashimi plots of a DAS exon at gene *CDV3* detected by LIQA, but was missed by Penn-Diff using short-read data. From long-read data, it is clear that the relative expression of exon in the red square is lower in cancerous cells than normal-like ESCC. However, this event is missed by short-read data. The read coverage difference between normal-like and cancerous ESCC in sashimi plots indicates the less usage of isoforms (NM_001134422, NM_001134423, NM_001282765) which include this exon in ESCC cells, suggesting better performance of long-read data.

Discussion

Accurate estimation of isoform-specific gene expression is a critical step for transcriptome profiling. The emergence of long-read RNA-seq has made it possible to discover complex novel isoforms and quantify isoform usage based on full-length sequenced fragments without amplification bias. However, there are still issues for long-read data, which if not taken into account, can affect the estimations. The major challenges in the



analysis of long-read RNA-seq data are the presence of high error rate and potential coverage bias. In this article, we propose LIQA, a statistical method that allows read-specific weight in estimating isoform-specific gene expression. The central idea of our method is to extract error rate information and model non-uniformity read coverage distribution of long-read data. LIQA is the first long-read transcriptomic tool that takes these limitations of long-read RNA-seq data into account. Results of our simulation study and analyses of real data demonstrated that LIQA is more effective in bias correction than the limited existing approaches (Additional file 1: Table S2, S3).

However, we note that there is still room to improve LIQA. LIQA is computationally intensive because the approximation of nonparametric Kaplan-Meier estimator of function $f(L_r)$ relies on empirical read length distribution and the parameters are estimated using EM algorithm. Based on the analysis of the UHR and AML data, we found that running LIQA is slower than FLAIR and Mandalorion (Additional file 1: Table S1).

Currently, we are evaluating the impact of possible parametric functions such as exponential or Weibull distributions for read distribution modeling. This will sacrifice the robustness of isoform expression estimates but the running time can be significantly reduced. We believe it may be worth making this trade-off between computational efficiency and estimation accuracy for LIQA.

We have benchmarked the performance of LIQA with the use of minimap2 for long-read alignment, while there have been several approaches supporting RNA-seq long-read alignment, such as STAR [37], GMAP [44], BLAT [45], BMap (<https://sourceforge.net/projects/bbmap/>), and GraphMap 2[46]. LIQA can take SAM or BAM files generated from any listed aligner as input. Nevertheless, we recognize that it is important to evaluate whether LIQA's superior performance is robust to different aligners. Therefore, we plan to explore more long-read aligner options and settings to benchmark LIQA in the future.

As LIQA is EM algorithm-based, the robustness to parameter initialization is a potential issue. Read-specific weight of LIQA extracts more information from observed data than direct read count strategy as implemented in Mandalorion and FLAIR. Especially, more read coverage is needed for stable approximation of function $f(L_r)$. For genes with limited reads coverage (less than 5), the likelihood function of LIQA will be flattened, then optimal points are harder to be reached by EM algorithm and estimates may be sensitive to initial values of parameter. Therefore, the sensitivity of LIQA to parameter initialization should be further evaluated and improved.

With full-length transcript sequencing, long-read RNA-seq data (ONT and PacBio) are expected to facilitate transcriptomic studies by offering number of advantages over short reads. For PacBio, HiFi reads are generated with circular consensus sequencing (CCS) using single-molecule consensus, which increases their accuracy over traditional multi-molecule consensus. Compared to Nanopore sequencing, this protocol yields much lower per-base error rate compared to Nanopore sequencing, but potentially shorter reads. Smaller read length may introduce much larger biases in 5' or 3' coverage ratio, which requires further adjustment for LIQA to derive more accurate isoform expression estimates. LIQA has custom settings that allow users to flexibly adjust such parameters to handle these platforms. Compared to PacBio (either with traditional library or HiFi library preparation protocols), ONT may be a more promising platform in quantifying isoform expression while generating data with much higher error rate. This is because ONT is currently more affordable with lower per-based cost of data generation, and sequencing data with high read coverage can improve estimation accuracy of isoform usage. For ONT-RNA-seq, there are two types: direct mRNA sequencing and cDNA sequencing. Compared to direct mRNA sequencing, cDNA sequencing allows samples to be barcoded, amplified and requires less amounts of starting materials. Our studies showed that the decrease of read coverage had less impact on LIQA compared to other existing approaches.

In summary, long-read RNA-seq data offer advantages and can help us better understand transcriptomic variations than short-read data. However, better utilizing informative single-molecule sequencing read is not straightforward. LIQA is a robust and effective computational tool to estimate isoform-specific gene expression from long-read RNA-seq data. With the increasing adoption of long-read RNA-seq in biomedical

research, we believe LIQA will be well-suited for various transcriptomics studies and offer additional insights beyond gene expression analysis in the future.

Methods and materials

Complete likelihood function of LIQA

Given a gene of interest, let \mathbf{R} denote the set of reads that are mapped to the gene of interest, and \mathbf{I} denote the set of known isoforms. For a specific isoform $i \in \mathbf{I}$, let θ_i denote its relative abundance, with $0 \leq \theta_i \leq 1$ and $\sum_{i \in \mathbf{I}} \theta_i = 1$ and l_i denote its length. For each single-molecule long-read r , let L_r denote its length. The probability that a read originates from isoform i is $P(\text{iso.} = i) = \theta_i$. For read-isoform assignment, LIQA accounts for incorrect alignment at splice site. We define parameter $\mathbf{Z}_{\mathbf{R},\mathbf{I}}$ as a $|\mathbf{R}| \times |\mathbf{I}|$, a read-isoform compatibility matrix with $Z_{\mathbf{R},\mathbf{I}}(r, i) = 1$ if long-read r is generated from a molecule that is originated from isoform i (number of mismatch base pairs < 20 bp instead of exact match), and $Z_{\mathbf{R},\mathbf{I}}(r, i) = 0$ otherwise. For isoform quantification, our goal is to estimate $\Theta = \{\theta_i, i \in \mathbf{I}\}$ based on RNA-seq long reads mapped to the gene.

With the notation above, the complete data likelihood of the RNA-seq data can be written as

$$\begin{aligned} L(\tilde{\Theta}|\mathbf{R}, \mathbf{Z}) &= \prod_{r \in \mathbf{R}} \prod_{i \in \mathbf{I}} (P(\text{read} = r, \text{read len.} = L_r, \text{iso.} = i))^{Z_{\mathbf{R},\mathbf{I}}(r,i)} \\ &= \prod_{r \in \mathbf{R}} \prod_{i \in \mathbf{I}} (P(\text{read} = r, \text{read len.} = L_r | \text{iso.} = i) \cdot P(\text{iso.} = i))^{Z_{\mathbf{R},\mathbf{I}}(r,i)} \\ &= \prod_{r \in \mathbf{R}} \prod_{i \in \mathbf{I}} (P(\text{read} = r, \text{read len.} = L_r | \text{iso.} = i) \cdot \theta_i)^{Z_{\mathbf{R},\mathbf{I}}(r,i)} \end{aligned}$$

This formula is based on the fact that given the isoform origin, the probability of observing read alignment can be inferred. The conditional probability of read r derived from isoform i with length L_r is

$$P(\text{read} = r, \text{read len.} = L_r | \text{iso.} = i) = q(r, i) \cdot f(L_r | \text{iso.} = i)$$

where $q(r, i)$ is isoform-specific read quality score and $f(L_r | \text{iso.} = i)$ is isoform-specific read length probability. Essentially, we quantify isoform relative abundance with weighted read assignment. To account for the error-prone manner of Nanopore sequencing data, we consider isoform-specific read quality score $q(r, i) = \prod_{j=1}^m q_j(x_j, y_{(j)})$ where x is the sequence of the long-read r , y is the sequence of the corresponding isoform i in the reference genome, and $q_j(a, b)$ is the probability that we observe base a at position j of the read given that the true base is b , which can be calculated as $1 - 10^{-Q_j/10}$, with Q_j being the per-based Phred quality score at position j .

Estimation of isoform-specific read length probability $f(L_r | \text{iso.} = i)$

Because read length L_r is not fixed and short prone in Nanopore sequencing, we treat it as a random variable with right skewed distribution density function $f(\cdot)$. Given an isoform, existing long-read methods assume fixed read length for all sequenced reads, and this is equivalent to setting $f(L_r)$ at 1. However, this assumption does not hold as recent studies suggest that potential 3' coverage bias exists in long-read RNA-seq data [24, 32, 47]. To offer flexibility in modeling read length distribution, we employ a non-parametric approach. For all long reads mapped to the genome, we categorize them

into two groups: complete reads and truncated reads. Accounting for misalignment due to high error rate, the read is treated as complete when the distance between its ending alignment position and any known isoform 5' end is less than a tolerance threshold (default = 20 bp) (Fig. 1A). This indicates that this read is completely sequenced from a known isoform. Otherwise, the read is considered as truncated. The presence of truncated reads is due to incomplete sequencing or novel isoforms. As known annotated isoforms are treated as gold standard during estimation, we assume true length of truncated read is censored. Given the observed lengths of all complete and truncated reads, we fit them into a survival model, a natural modeling approach for censored data (Additional file 1: Fig. S9, S10, S11). Function $\hat{F}(l) = P(\text{read len.} < l)$ can be estimated based on Kaplan-Meier estimator [48], hence we have $f(l) = \hat{F}(l+1) - \hat{F}(l)$.

Given a gene of interest with $I = \{\text{isoform } i : 1 \leq i \leq I\}$, isoform-specific read length probability $f(L_r | \text{iso.} = i)$ can be written as

$$f(L_r | \text{iso.} = i) = \frac{f(L_r) \cdot P(\text{iso.} = i | L_r)}{P(\text{iso.} = i)} = \frac{f(L_r) \cdot \theta_i / \sum_{l_j > L_r} \theta_j}{\theta_i} = \frac{f(L_r)}{\sum_{l_j > L_r} \theta_j}$$

This isoform-specific read length probability $f(L_r | \text{iso.} = i)$ captures the sequencing biases due to fragmentation during library preparation or pore-blocking for nanopore data.

Quantification of isoform expression level

Given that isoform indicators $Z_{R,I}(r, i)$ for some reads are not observed from read data, Θ are estimated using EM algorithm. Then, we have isoform relative abundance $\hat{\theta}_i$. In addition to relative abundance, it is also important to quantify the absolute expression level of an isoform. At gene level, we consider read per gene per 10 K reads (RPG 10 K) as the standard for long-read RNA-seq data. RPG is defined as $\text{RPG} = N/10^4$ where N is the number of reads mapped to the gene of interest. With this concept, we estimate the expression level of a particular isoform by replacing N with estimated number of long reads originated from isoform i ($\text{RPG}_i = N \cdot \hat{\theta}_i / 10^4$).

Parameter estimation using the EM algorithm

The complete data likelihood is

$$L(\Theta | \mathbf{R}, \mathbf{Z}) = \prod_{r \in R} \prod_{i \in I} (q(r, i) \cdot f(L_r) \cdot \theta_i)^{Z_{R,I}(r, i)}$$

and the update procedure of the EM algorithm is as follows:

E-step: We calculate function

$$\begin{aligned} Q(\Theta | \Theta^{(t)}) &= E_{Z_{R,I} | \Theta^{(t)}} [\log L(\Theta | R)] \\ &= \sum_{r \in R} \sum_{i \in I} E_{Z_{R,I} | \Theta^{(t)}} [Z_{R,I}(r, i)] \cdot \log(q(r, i) f(L_r) \theta_i) \end{aligned}$$

where $E_{Z_{R,I} | \Theta^{(t)}} [Z_{R,I}(r, i)] = \frac{q(r, i) f(L_r) \theta_i^{(t)}}{\sum_{u \in I} q(r, u) f(L_r) \theta_u^{(t)}}$.

M-step: We maximize function $Q(\Theta | \Theta^{(t)})$ and have

$$\theta_i^{(t+1)} = \frac{\sum_{r \in \mathcal{R}} E_{Z_{R,I}|\Theta^{(t)}} [Z_{R,I}(r, i)]}{|\mathcal{R}|}$$

The EM algorithm consists of alternating between the E- and M-steps until convergence. We start the algorithm with $\Theta^{(0)}$ assuming all isoforms are equally expressed and stop when the log likelihood is no longer increasing significantly.

Detection of differential alternative splicing (DAS) with LIQA

The relative abundance of an isoform takes values between 0 and 1. Therefore, we assume it follows a beta distribution, which is well known as a flexible distribution in modeling proportion because its density can have different shapes depending on the values of the two parameters that characterize the distribution, i.e., $\theta_i \sim \text{Beta}(\mu_i, \phi_i)$. The expected value and variance of θ_i are

$$E(\theta_i) = \mu_i$$

$$\text{Var}(\theta_i) = \frac{\mu_i(1-\mu_i)}{1 + \phi_i}$$

To detect splicing difference of isoform i between two groups of samples, we utilized beta regression model with ϕ_i as precision parameter. We apply logit link function and have the model

$$\text{logit}(\theta_i) = \beta_0 + \beta_1 Z$$

where Z is the condition indicator (1 for case; 0 for control), β_0 and β_1 are coefficient parameters.

Since the isoform relative abundances of isoforms within the same gene are correlated, a robust and flexible model is needed when comparing them between conditions at a gene level. To account for this, we utilize Gaussian copula regression model to test splicing difference significance between conditions of correlated isoform relative abundances. The separation of marginal distributions and correlation structure makes Gaussian copula regression versatile in modeling non-normal dependent observations. Therefore, the joint distribution of isoform relative abundances from the same gene is given by

$$\Phi_{I-1}(\Phi^{-1}(F(\theta_1|\beta_0, \beta_1, \phi_1)), \dots, \Phi^{-1}(F(\theta_{I-1}|\beta_0, \beta_{I-1}, \phi_{I-1})))|\mathbf{I}$$

where ϕ_i is the dispersion parameter of the marginal generalized linear model for isoform i . $\Phi_{I-1}(|\mathbf{I})$ is the cumulative distribution function of multivariate normal random variables with $I-1$ dimensions and correlation matrix \mathbf{I} . We choose to use exchangeable correlation structure for \mathbf{I} . Given regression models above, we can detect DAS both for at the isoform level and at gene level. For isoform i , we test $H_0: \beta_1 = 0$ vs $H_1: \beta_1 \neq 0$ to determine splicing change between conditions. For gene g , we test $H_0: \beta_1^1 = \dots = \beta_1^{I-1} = 0$ vs $H_1: \beta_1^i \neq 0$ for any $1 \leq i \leq I-1$.

Nanopore direct mRNA sequencing of Universal Human Reference RNA-seq data

Universal human reference (UHR) RNA comprises of mixed RNA molecules by a diverse set of 10 cancer cell lines with equal quantities of DNase-treated RNA from adenocarcinoma in mammary gland, hepatoblastoma in liver, adenocarcinoma in

cervix, embryonal carcinoma in testis, glioblastoma in brain, melanoma, liposarcoma, histocytic lymphoma in histocyte macrophage, lymphoblastic leukemia, and plasmacytoma in B lymphocyte. This reference sample from MicroArray Quality Control (MAQC) [49–51] project has been utilized in many studies. For example, Gao et al. [52] sequenced this UHR RNA sample and treated it as reference to measure the technical variations of scRNA-seq data. Also, the qRT-PCR measurements of gene or isoform expressions from this sample were used to benchmark and optimize computational tools [17, 53–56]. In this study, we used GridION Nanopore technique to sequence mRNA directly and used Guppy for base calling. In total, we generated 476,000 long reads with 557 MB bases. We aligned the UHR RNA-seq data against a reference genome (hg38) using minimap2 [35], and 95% long reads (89% of total bases) are mapped, demonstrating very high sequencing and basecalling quality. qRT-PCR measurements were downloaded and treated as ground truth to compare the performance between LIQA, FLAIR, Mandalorion, CEM, Cufflinks, and RD.

Chi-squared goodness of fit statistics of read coverage uniformity

Given an isoform of interest, let l denote the length and O_i denote observed read coverage count at base pair position i . Total sequencing depth of this isoform $S = \sum_{1 \leq i \leq l} O_i$.

Under the uniform read coverage assumption, the expected read coverage count at each base pair position $E_i = S/l$. We apply chi-squared goodness of fit statistics to measure the difference between observed read coverage and uniform read distribution. The test statistics is

$$\chi^2 = \sum_{1 \leq i \leq l} \frac{(O_i - E_i)^2}{E_i}$$

The degree of freedom is isoform length $l_i - 1$. The higher value of χ^2 indicates that observed read coverage deviates more from uniform read distribution. We calculated χ^2 for each isoform, then divided them into two categories based on median of the corresponding measure (less uniform: $\chi^2 > \text{median}$, more uniform $\chi^2 < \text{median}$) to evaluate the impact of read coverage distribution on isoform expression quantification.

Statistical test to compare performance of different methods

We simulated ONT RNA-seq data 20 times to assess the statistical significance when comparing the performance of different methods. Each dataset includes 5 million (5 M) reads. We also down-sampled 3 million (3 M), 1 million (1 M), and 0.5 million (0.5 M) reads for the simulated data to evaluate the impact of sequencing depth on performance improvement of LIQA. We ran all methods with the same set of simulated aligned data in BAM format as input and calculated Spearman's correlation of TPM and relative abundance between true and estimated values. Based on this metric from 20 simulated datasets, we conducted pairwise comparison of performance difference between all methods using paired Z-test. Mean difference, standard deviations, test statistics, and P values were calculated. Moreover, we conducted likelihood ratio test to compare different models of LIQA (full model, read length model, read quality model). Likelihood ratio test statistic $Q = -2(\log L_B - \log L_A)$, where L is the optimized likelihood function based on different models.

Availability of data and materials

LIQA is freely available at <https://github.com/WGLab/LIQA> under GPLv3 license [57]. The direct mRNA sequencing data on UHR has been deposited and available at Gene Expression Omnibus (PRJNA639366) [58]. The cDNA sequencing data on a patient with cancer has been deposited and available at Gene Expression Omnibus (PRJNA640456) [59]. The simulation data used in our study can be reproduced using code provided in the LIQA software repository and NanoSim version 2.0.0. The esophageal squamous epithelial cell PacBio RNA-seq data we applied is from Gene Expression Omnibus (PRJNA515570) [40]. Source code used in the manuscript is available via Zenodo with DOI 10.5281/zenodo.4795477 [60].

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-021-02399-8>.

Additional file 1. Figures S1-S13 and Tables S1-S7.

Additional file 2. Review history.

Acknowledgements

We thank the Wang lab members for insightful comments and for testing the software tools. We also thank the developers of the NanoSim software tool, and the generators of the short-read and qRT-PCR results on the UHR datasets and short- and long-read data on the ESCC datasets for making the data publicly available for benchmarking studies. We thank three anonymous reviewers for their constructive comments and suggestions on benchmarking studies.

Review history

The review history is available as Additional file 2.

Peer review information

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

Authors' contributions

YH and KW initiated and designed the project. YH, ML, and KW formulated the model. YH developed and implemented the algorithm. LF performed long-read sequencing experiment and data processing. XC and JFZ prepared cancer sample for sequencing. YH, ML, and KW conducted the analysis and wrote the manuscript. All authors read and approved the final manuscript.

Funding

This study is supported by NIH/NIGMS grant GM132713 and the CHOP Research Institute.

Declarations

Ethics approval and consent to participate

The genomic study on patients with acute myeloid leukemia was carried out under IRB-approved protocol (# HS-20-00223). The patient has given written informed consent for publication. The experimental methods comply with the Helsinki Declaration.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children's Hospital of Philadelphia, Philadelphia, PA 19104, USA. ²Department of Otolaryngology, Keck School of Medicine, University of Southern California, Los Angeles, CA 90033, USA. ³Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA. ⁴Department of Pathology and Laboratory Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA.

Received: 18 September 2020 Accepted: 4 June 2021

Published online: 17 June 2021

References

- Han J, Xiong J, Wang D, Fu XD. Pre-mRNA splicing: where and when in the nucleus. *Trends Cell Biol.* 2011;21:336–43. <https://doi.org/10.1016/j.tcb.2011.03.003>
- Scotti MM, Swanson MS. RNA mis-splicing in disease. *Nat Rev Genet.* 2016;17:19–32. <https://doi.org/10.1038/nrg.2015.3>
- Montes M, Sanford BL, Comiskey DF, Chandler DS. RNA splicing and disease: animal models to therapies. *Trends Genet.* 2019;35:68–87. <https://doi.org/10.1016/j.tig.2018.10.002>
- Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary link between genetic variation and disease. *Science.* 2016;352:600–4. <https://doi.org/10.1126/science.aad9417>
- Kim HK, Pham MHC, Ko KS, Rhee BD, Han J. Alternative splicing isoforms in health and disease. *Pflugers Arch.* 2018;470:995–1016. <https://doi.org/10.1007/s00424-018-2136-x>
- Pan Q, Shai O, Lee LJ, Frey BJ, Blencowe BJ. Deep surveying of alternative splicing complexity in the human transcriptome by high-throughput sequencing. *Nat Genet.* 2008;40:1413–5. <https://doi.org/10.1038/ng.259>
- Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, et al. Alternative isoform regulation in human tissue transcriptomes. *Nature.* 2008;456:470–6. <https://doi.org/10.1038/nature07509>
- Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat Rev Genet.* 2009;10:57–63. <https://doi.org/10.1038/nrg2484>
- Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics.* 2011;12:323. <https://doi.org/10.1186/1471-2105-12-323>
- Roberts A, Pachter L. Streaming fragment assignment for real-time analysis of sequencing experiments. *Nat Methods.* 2013;10:71–3. <https://doi.org/10.1038/nmeth.2251>
- Nariai N, Kojima K, Mimori T, Sato Y, Kawai Y, Yamaguchi-Kabata Y, et al. TIGAR2: sensitive and accurate estimation of transcript isoform expression with longer RNA-Seq reads. *BMC Genomics.* 2014;15 Suppl 10:S5. <https://doi.org/10.1186/1471-2164-15-S10-S5>
- Zhang C, Zhang B, Lin LL, Zhao S. Evaluation and comparison of computational tools for RNA-seq isoform quantification. *BMC Genomics.* 2017;18:583. <https://doi.org/10.1186/s12864-017-4002-1>
- Patro R, Mount SM, Kingsford C. Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nat Biotechnol.* 2014;32:462–4. <https://doi.org/10.1038/nbt.2862>
- Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol.* 2016;34:525–7. <https://doi.org/10.1038/nbt.3519>
- Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, et al. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc.* 2012;7:562–78. <https://doi.org/10.1038/nprot.2012.016>
- Li W, Jiang T. Transcriptome assembly and isoform expression level estimation from biased RNA-Seq reads. *Bioinformatics.* 2012;28:2914–21. <https://doi.org/10.1093/bioinformatics/bts559>
- Hu Y, Liu Y, Mao X, Jia C, Ferguson JF, Xue C, et al. PennSeq: accurate isoform-specific gene expression quantification in RNA-Seq by modeling non-uniform read distribution. *Nucleic Acids Res.* 2014;42:e20. <https://doi.org/10.1093/nar/gkt1304>
- Nicolae M, Mangul S, Mandoiu II, Zelikovsky A. Estimation of alternative splicing isoform frequencies from RNA-Seq data. *Algorithms Mol Biol.* 2011;6:9. <https://doi.org/10.1186/1748-7188-6-9>
- Wan L, Yan X, Chen T, Sun F. Modeling RNA degradation for RNA-Seq with applications. *Biostatistics.* 2012;13:734–47. <https://doi.org/10.1093/biostatistics/kxs001>
- Steijger T, Abril JF, Engstrom PG, Kokocinski F, Consortium R, Hubbard TJ, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods.* 2013;10:1177–84. <https://doi.org/10.1038/nmeth.2714>
- Tilgner H, Grubert F, Sharon D, Snyder MP. Defining a personal, allele-specific, and single-molecule long-read transcriptome. *Proc Natl Acad Sci U S A.* 2014;111:9869–74. <https://doi.org/10.1073/pnas.1400447111>
- Burgess DJ. Genomics: Next regeneration sequencing for reference genomes. *Nat Rev Genet.* 2018;19:125. <https://doi.org/10.1038/nrg.2018.5>
- Pollard MO, Gurdasani D, Mentzer AJ, Porter T, Sandhu MS. Long reads: their purpose and place. *Hum Mol Genet.* 2018;27:R234–41. <https://doi.org/10.1093/hmg/ddy177>
- Sharon D, Tilgner H, Grubert F, Snyder M. A single-molecule long-read survey of the human transcriptome. *Nat Biotechnol.* 2013;31:1009–14. <https://doi.org/10.1038/nbt.2705>
- Tilgner H, Jahanbani F, Blauwkamp T, Moshrefi A, Jaeger E, Chen F, et al. Comprehensive transcriptome analysis using synthetic long-read sequencing reveals molecular co-association of distant splicing events. *Nat Biotechnol.* 2015;33:736–42. <https://doi.org/10.1038/nbt.3242>
- Treutlein B, Gokce O, Quake SR, Sudhof TC. Cartography of neurexin alternative splicing mapped by single-molecule long-read mRNA sequencing. *Proc Natl Acad Sci U S A.* 2014;111:E1291–9. <https://doi.org/10.1073/pnas.1403244111>
- Vollmers C, Penland L, Kanbar JN, Quake SR. Novel exons and splice variants in the human antibody heavy chain identified by single cell and single molecule sequencing. *PLoS One.* 2015;10:e0117050. <https://doi.org/10.1371/journal.pone.0117050>
- Oikonomopoulos S, Wang YC, Djambazian H, Badescu D, Ragoussis J. Benchmarking of the Oxford Nanopore MinION sequencing for quantitative and qualitative assessment of cDNA populations. *Sci Rep.* 2016;6:31602. <https://doi.org/10.1038/srep31602>
- Roberts A, Trapnell C, Donaghey J, Rinn JL, Pachter L. Improving RNA-Seq expression estimates by correcting for fragment bias. *Genome Biol.* 2011;12:R22. <https://doi.org/10.1186/gb-2011-12-3-r22>
- Byrne A, Beaudin AE, Olsen HE, Jain M, Cole C, Palmer T, et al. Nanopore long-read RNAseq reveals widespread transcriptional variation among the surface receptors of individual B cells. *Nat Commun.* 2017;8:16027. <https://doi.org/10.1038/ncomms16027>
- Tang AD, Soulette CM, van Baren MJ, Hart K, Hrabeta-Robinson E, Wu CJ, et al. Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals downregulation of retained introns. *Nat Commun.* 2020;11:1438. <https://doi.org/10.1038/s41467-020-15171-6>

32. Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. Opportunities and challenges in long-read sequencing data analysis. *Genome Biol.* 2020;21:30. <https://doi.org/10.1186/s13059-020-1935-5>
33. Wyman D, Balderrama-Gutierrez G, Reese F, Jiang S, Rahmanian S, Forner S, et al. A technology-agnostic long-read analysis pipeline for transcriptome discovery and quantification. *BioRxiv.* 2020. <https://doi.org/10.1101/672931>
34. Hafezqorani S, Yang C, Lo T, Nip KM, Warren RL, Birol I. Trans-NanoSim characterizes and simulates nanopore RNA-sequencing data. *Gigascience.* 2020;9:gjaa061. <https://doi.org/10.1093/gigascience/gjaa061>
35. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34:3094–100. <https://doi.org/10.1093/bioinformatics/bty191>
36. Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods.* 2017;14:417–9. <https://doi.org/10.1038/nmeth.4197>
37. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29:15–21. <https://doi.org/10.1093/bioinformatics/bts635>
38. Sequel II system data release: universal human reference (UHR) iso seq. [https://github.com/PacificBiosciences/DevNet/wiki/Sequel-II-System-Data-Release-Universal-Human-Reference-\(UHR\)-Iso-Seq](https://github.com/PacificBiosciences/DevNet/wiki/Sequel-II-System-Data-Release-Universal-Human-Reference-(UHR)-Iso-Seq).
39. De Kouchkovsky I, Abdul-Hay M. Acute myeloid leukemia: a comprehensive review and 2016 update. *Blood Cancer J.* 2016;6:e441. <https://doi.org/10.1038/bcj.2016.50>
40. Cheng YW, Chen YM, Zhao QQ, Zhao X, Wu YR, Chen DZ, et al. Long read single-molecule real-time sequencing elucidates transcriptome-wide heterogeneity and complexity in esophageal squamous cells. *Front Genet.* 2019;10:915. <https://doi.org/10.3389/fgene.2019.00915>
41. Hu Y, Lin J, Hu J, Hu G, Wang K, Zhang H, et al. PennDiff: detecting differential alternative splicing and transcription by RNA sequencing. *Bioinformatics.* 2018;34:2384–91. <https://doi.org/10.1093/bioinformatics/bty097>
42. Xie ZC, Wu HY, Ma FC, Dang YW, Peng ZG, Zhou HF, et al. Prognostic alternative splicing signatures and underlying regulatory network in esophageal carcinoma. *Am J Transl Res.* 2019;11:4010–28.
43. Ueno N, Shimizu A, Kanai M, Iwaya Y, Ueda S, Nakayama J, et al. Enhanced expression of fibroblast growth factor receptor 3 Il1c promotes human esophageal carcinoma cell proliferation. *J Histochem Cytochem.* 2016;64:7–17. <https://doi.org/10.1369/0022155415616161>
44. Wu TD, Watanabe CK. GMAP: a genomic mapping and alignment program for mRNA and EST sequences. *Bioinformatics.* 2005;21:1859–75. <https://doi.org/10.1093/bioinformatics/bti310>
45. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res.* 2002;12:656–64. <https://doi.org/10.1101/gr.229202>
46. Maric J, Sovic I, Krizanovic K, Nagarajan N, Sikic M. Graphmap2-splice-aware RNA-seq mapper for long reads. *bioRxiv.* 2019. <https://doi.org/10.1101/720458>.
47. Kellner S, Burhenne J, Helm M. Detection of RNA modifications. *RNA Biol.* 2010;7:237–47. <https://doi.org/10.4161/rna.7.2.11468>
48. Kaplan EL, Meier P. Nonparametric estimation from incomplete observations. *J Am Stat Assoc.* 1958;53:457–81.
49. Consortium M, Shi L, Reid LH, Jones WD, Shippy R, Warrington JA, et al. The MicroArray Quality Control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements. *Nat Biotechnol.* 2006;24:1151–61. <https://doi.org/10.1038/nbt1239>
50. Shi L, Campbell G, Jones WD, Campagne F, Wen Z, Walker SJ, et al. The MicroArray Quality Control (MAQC)-II study of common practices for the development and validation of microarray-based predictive models. *Nat Biotechnol.* 2010;28:827–38. <https://doi.org/10.1038/nbt.1665>
51. Sun P, Sehouli J, Denkert C, Mustea A, Könsgen D, Koch I, et al. Expression of estrogen receptor-related receptors, a subfamily of orphan nuclear receptors, as new tumor biomarkers in ovarian cancer cells. *J Mol Med.* 2005;83:457–67.
52. Gao F, Kim JM, Kim J, Lin M-Y, Liu CY, Russin JJ, et al. Evaluation of biological and technical variations in low-input RNA-Seq and single-cell RNA-Seq. *Int J Comput Biol Drug Des.* 2018;11:5–22.
53. Xu J, Su Z, Hong H, Thierry-Mieg J, Thierry-Mieg D, Kreil DP, et al. Cross-platform ultradeep transcriptomic profiling of human reference RNA samples by RNA-Seq. *Sci Data.* 2014;1:140020. <https://doi.org/10.1038/sdata.2014.20>
54. Garcia-Alonso L, Holland CH, Ibrahim MM, Turei D, Saez-Rodriguez J. Benchmark and integration of resources for the estimation of human transcription factor activities. *Genome Res.* 2019;29:1363–75. <https://doi.org/10.1101/gr.240663.118>
55. Teng M, Love MI, Davis CA, Djebali S, Dobin A, Graveley BR, et al. A benchmark for RNA-seq quantification pipelines. *Genome Biol.* 2016;17:74. <https://doi.org/10.1186/s13059-016-0940-1>
56. Hayer KE, Pizarro A, Lahens NF, Hogenesch JB, Grant GR. Benchmark analysis of algorithms for determining and quantifying full-length mRNA splice forms from RNA-seq data. *Bioinformatics.* 2015;31:3938–45. <https://doi.org/10.1093/bioinformatics/btv488>
57. Hu Y, Li M, Wang K. LIQA: long-read isoform quantification and analysis. Github. 2021; <https://github.com/WGLab/LIQA>
58. Hu Y, Fang L, Chen X, Zhong JF, Li M, Wang K. Long-read sequencing of reference RNA samples. *Datasets Gene Expression Omnibus.* 2021. <https://www.ncbi.nlm.nih.gov/bioproject/?term=PRJNA639366>.
59. Hu Y, Fang L, Chen X, Zhong JF, Li M, Wang K. Oxford nanopore sequencing of acute myeloid leukemia samples. *Datasets Gene Expression Omnibus.* <https://www.ncbi.nlm.nih.gov/bioproject/?term=PRJNA640456>.
60. Hu Y, Li M, Wang K. LIQA: long-read isoform quantification and analysis. Zenodo. 2021; <https://doi.org/10.5281/zenodo.4795477>.

Publisher's Note

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