

LIQA: Long-read Isoform Quantification and Analysis

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

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

² Department of Otolaryngology, Keck School of Medicine, University of Southern California, CA 90033, USA

³ Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA

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

* To whom correspondence should be addressed. Email: wangk@email.chop.edu

Abstract

Long-read RNA sequencing (RNA-seq) technologies have made it possible to sequence full-length transcripts, facilitating the exploration of isoform-specific gene expression over conventional short-read RNA-seq. However, long-read RNA-seq suffers from high per-base error rate, presence of chimeric reads and alternative alignments, and other biases, which require different analysis methods than short-read RNA-seq. Here we present LIQA (Long-read Isoform Quantification and Analysis), an Expectation-Maximization based statistical method to quantify isoform expression and detect differential alternative splicing (DAS) events using long-read RNA-seq data. Rather than summarizing isoform-specific read counts directly as done in short-read methods, LIQA incorporates base-pair quality score and isoform-specific read length information to assign different weights across reads, which reflects alignment confidence. Moreover, given isoform usage estimates, LIQA can detect DAS events between conditions. We evaluated LIQA's performance on simulated data and demonstrated that it outperforms other approaches in rare isoform characterization and in detecting DAS events between two groups. We also generated one direct mRNA sequencing dataset and one cDNA sequencing dataset using the Oxford Nanopore long-read platform, both with paired short-read RNA-seq data and qPCR data on selected genes, and we demonstrated that LIQA performs well in isoform discovery and quantification. Finally, we evaluated LIQA on a PacBio dataset on esophageal squamous epithelial cells, and demonstrated that LIQA recovered DAS events on *FGFR3* that failed to be detected in short-read data. In summary, LIQA leverages the power of long-read RNA-seq and achieves higher accuracy in estimating isoform abundance than existing approaches, especially for isoforms with low coverage and biased read distribution.

Introduction

RNA splicing is a major mechanism for generating transcriptomic variations, and mis-regulation of splicing causes a large array of human diseases due to hereditary and somatic mutation[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. Until now, short-read RNA-seq has been the method of choice for transcriptomics studies due to its high coverage and single nucleotide resolution. However, due to limited read length (ranges from 50 to 150 nucleotides), transcripts cannot be fully sequenced and characterized by short-read RNA-seq data. This partial sequencing of the RNA results in biases and has become a barrier for short reads to be correctly mapped to the reference genome, which is crucial for gene/isoform expression estimation and novel/unique isoform detection. As a consequence, transcriptome profiling using short-read RNA-seq is highly biased by read coverage heterogeneity across isoform transcripts. To tackle these challenges, a number of computational tools, including RSEM[6], eXpress[7], TIGAR2[8], Salmon[9], Sailfish[10], Kallisto[11], Cufflinks[12], CEM[13], PennSeq[14], IsoEM[15], Stringtie[16], SLIDE[17], iReckon[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. Therefore, fragmented short reads cannot quantify isoform expression accurately, especially at complex gene loci.

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[20, 21]. Previous studies have utilized both single-molecule long-read PacBio Iso-Seq and synthetic long-read MOECULO methods[22-25]. For Oxford Nanopore sequencing, there are two types of RNA-seq technologies: direct RNA sequencing and cDNA sequencing. Recently, the Oxford Nanopore Technologies (ONT) MinION has been used to analyze both full-length cDNA samples and RNA samples derived from tissue cells[26]. Nanopore sequencing is able to generate reads as long as 2Mbp, 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.

Moreover, recent success of DNA modification detection using Nanopore data further indicates the opportunity to characterize RNA modifications with the use of direct RNA sequencing.

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. Methods designed specifically for isoform expression estimation in long-read RNA-seq have only emerged recently. For example, Byrne *et al*[27] demonstrated the feasibility of quantifying complex isoform expression using Nanopore RNA-seq data. Tang *et al*[28] 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[29], 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[29]. In Nanopore direct RNA sequencing protocol, pore block and fragmentation can result in truncated reads, leading to biased coverage towards the 3' end of a transcript. 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 (Figure 1). To evaluate the performance of LIQA, we simulated long data with known ground truth and also sequenced two real samples using Nanopore sequencing. Our results demonstrate that LIQA outperforms existing methods in isoform expression quantification.

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/SAM format and isoform annotation file as input. For estimation steps, LIQA first feeds read alignment information to a complete likelihood function and correct biases for each long-read by accounting for quality score and read length probability. 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 and Mandalorion. These two methods use long-read RNA-seq data to detect novel isoforms and quantify transcript expressions by counting the number of reads, which give equal weight for each read. To make the comparisons fair, we ran LIQA, FLAIR and Mandalorion in quantification mode only with isoform annotation information provided. 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 annotation. To simulate a realistic dataset with known ground truth, we used NanoSim [30] 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 (5M) Nanopore reads. To evaluate the impact of sequencing depth on isoform expression quantification, we down-sampled 3 (3M), 1 (1M) and 0.5 (0.5M) million reads for the simulated data, respectively. These reads were aligned to the reference genome using minimap2 [31]. Then, we selected genes with 2 or more isoforms to evaluate the performance of LIQA in isoform expression quantification. For each isoform, we

compared it with Mandalorion [27] and FLAIR [32]. 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 Figure 2(A) and Supplementary Figure 1. The median lengths of ONT reads in the 0.5M, 1M, 3M and 5M datasets are 896, 922, 1010 and 993 base pairs, respectively. Among the evaluated genes with multiple isoforms 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 = 0.75, 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. These simulated data thus provide an ideal basis to evaluate the performance of LIQA as the ground truth is known.

Gene/Isoform expression quantification accuracy

For each simulated dataset, we computed different 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 the coefficient of determination (i.e. R squared). 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 2(A) shows the summary statistics between estimated and true values of global isoform expression (global gene expression and isoform relative abundances) at different read coverages. Spearman correlation and RMSE were calculated for all three methods. LIQA and FLAIR clearly have higher Spearman correlation than Mandalorion across all simulated datasets. Compared with FLAIR, LIQA has similar estimation accuracy based on Spearman correlation. Figure 2(B) gives summary statistics of relative abundance estimates for the three methods. For relative abundance estimation, LIQA outperforms FLAIR and Mandalorion with 6.6% and 10.1% higher RMSE, respectively. The improved performance of LIQA is 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. In contrast, FLAIR 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.

Next, we evaluated the impact of the isoform length and 3' bias in isoform expression estimation. We considered two ways to compare the performance between methods. First, we divided the isoforms into 3 categories (length < 33% quantile, $33\% \leq \text{length} < 66\%$ quantile, length $\geq 66\%$) and summarized Spearman correlation coefficient and RMSE for each group of isoforms. The isoform lengths were calculated based on GENCODE annotation. Figure 3(A) shows the bar plots of statistics for three methods of the 0.5M dataset. For isoforms with length less than 66% quantile, Spearman correlation of FLAIR is greater than LIQA. However, despite of reduced Spearman correlation value, LIQA clearly outperforms the other two methods for isoforms longer than 66%. Longer isoforms are more challenging to estimate because the 5' end is less likely to be sequenced compared to shorter reads. We observed similar pattern when accuracy was quantified by RMSE. This is because LIQA models potential truncated reads from longer isoforms when quantifying isoform expression. Second, we compared the accuracy statistics for 5' terminal exon and 3' terminal exon of each isoform (Figure 3(B)). LIQA appears to be much more accurate than the other two methods, especially for 5' terminal exon. For example, the Spearman correlation coefficient of LIQA is 11% higher than the second best performed method FLAIR for 5' terminal exons, while only 6% higher for 3' exons. This superior performance of LIQA in terminal exons quantification is also revealed by RMSE. LIQA has 8%-15% RMSE improvement compared to FLAIR and Mandalorion. These results clearly demonstrate the advantage of LIQA in isoform length bias and 3' bias correction.

Differential alternative splicing (DAS) detection accuracy

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 isoform 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 and Mandalorion, respectively. We compared the performance of DAS detection between these methods using three summary statistics. First, 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 ($F1\ score = 2 \cdot \frac{precision \cdot recall}{precision + recall}$) was summarized to average the precision and recall values. As shown in Figure 3(C), LIQA and FLAIR clearly outperforms Mandalorion for all three evaluation metrics. This is not surprising because Mandalorion has lower accuracy than other two approaches for 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). This indicates that FLAIR has more inflated type-I error rate under significance level = 0.05. Overall, LIQA and FLAIR had similar performance in detecting DAS events based on F1 score.

Application to 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 hard to be full representation of reality. To evaluate the performance of LIQA in a real setting, we sequenced the Universal Human Reference sample with Nanopore Direct RNA-sequencing. Then, the resulting ONT-RNA-seq data were analyzed using all three long-read based methods (LIQA, FLAIR, Mandalorion). As quantitative real time polymerase chain reaction (qRT-PCR) is considered as the most popular technology for measuring true isoform abundance, we downloaded the qRT-PCR measurements from 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 Genome Analyzer platform. This dataset was analyzed using Cufflinks and CEM 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 and STAR, respectively, and applied each quantification method to the RNA-seq data. qRT-PCR measurements were treated as golden 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. Because estimating isoform-specific expression for these single-transcript genes is trivial, to better assess the performance of different methods, we only considered those transcripts that are derived from genes with two or more isoforms.

To quantify the similarity between estimates and qRT-PCR measurements, we calculated spearman correlation of the isoform abundance values in log-scale. As shown in Figure 4(A)(B), the estimation accuracy of all methods is lower than simulated data because the qRT-PCR measures may not be accurate, especially for those transcripts with qRT-PCR measures close to 0. Nevertheless, we observed consistent relative performance of different methods with simulation results. LIQA clearly outperforms 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 FPKM values being compacted toward 0. For ONT data, the spearman correlation of LIQA is 0.68, whereas the corresponding values from FLAIR and Mandalorian are only 0.45 and 0.48 only. For Illumina data, Cufflinks seems to correlate with the qRT-PCR measurements better than CEM. Comparison of different methods using RMSE reveals a similar pattern. The major reason of the better performance for LIQA is 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 Figure 4(C) to show the read coverage difference between direct RNA 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 3' degradation in Illumina data, but with full length and even coverage across the transcripts for long-read data. Terminal exons at 3' end in red square are crucial informative regions for splicing analysis, which enable us to differentiate read origin from different isoforms. As shown in Figure 4(C), these exonic regions were captured by Nanopore reads but missed by Illumina, 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.

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[33]. In this study, we next sequenced peripheral blood from an AML patient using GridION Nanopore sequencer with Guppy basecalling (<https://denbi-nanopore-training-course.readthedocs.io/en/latest/basecalling/basecalling.html#references>). In total, we generated 8,061,683 long-reads with 6.6 GB bases. We aligned the data against a reference genome (hg38) using minimap2[31], and 63% long-reads (73% bases) are mapped, indicating

high sequencing and basecalling quality. 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 a short-read data sequenced based on the same AML sample and treated the estimates as gold standard. This dataset included 440M short read with 150bp in length. Figure 5(A) shows the scatter plots of isoform relative abundance estimates between long- and short-read data. Pearson 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 Figure 5(B), long-read and short-read shared consistent estimates for the major isoforms. This is not surprised 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* (Figure 5(C)). It is clear that 3' read coverage biases exist in this read data. 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' 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 clearly 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)[34]. 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 difference 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 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, 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 richest event among detected DAS between normal-like and cancer esophageal cells, followed by RI, A5SS and A3SS. MXE is the most infrequent splicing types. As shown in Figure 5(D), 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 are greater than 98%. Compared to short-read, long-read data shows preference in detecting more differential splicing events at both exon and gene level. This is not surprising to us because read coverage heterogeneity, which might bias DAS detection, is significantly alleviated in long-read data by capturing full-length transcript in each read.

The expression of alternatively spliced isoforms from gene *FGFR3* shows preference in cancerous ESCC cells compared to non-cancerous[35]. Figure 5(E) provides the sashimi plots of two DAS exons at gene *FGFR3* detected by LIQA, but were missed by PennDiff using short-read data. From long-read data, it is clear that the relative expression of exon in the green square is higher in cancerous cells than normal-like ESCC. However, this event is missed by short-read data. Similarly, differential usage of the exon in the red square is detected in long-read data but missed in short-read data. The read coverage difference between normal-like and cancerous ESCC in sashimi plots indicates the flip-flop expression pattern of two exons between samples, 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 effective in bias correction than the limited existing approaches.

However, we note that there is still room to improve LIQA. LIQA is computationally intensive because the approximation of non-parametric 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 the running time of LIQA is slower than FLAIR and Mandalorion. 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 is 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[36], GMAP[37], BLAT[38], BMap (<https://sourceforge.net/projects/bbmap/>), and GraphMap2[39]. LIQA can take SAM/BAM file 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 covered (less than 5), the likelihood function of LIQA will be flattened, then optimal points are harder to be reached by EM-algorithm and estimates maybe 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-based error rate

compared to Nanopore sequencing, but potentially shorter reads. Smaller read length may introduce much larger biases in 5'/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 amplified and requires less amount 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 R denote the set of reads that are mapped to the gene of interest, and I denote the set of known isoforms. For a specific isoform $i \in I$, let θ_i denote its relative abundance, with $0 \leq \theta_i \leq 1$ and $\sum_{i \in I} \theta_i = 1$ 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 $P(\text{iso.} = i) = \frac{\theta_i l_i}{\sum_{i \in I} \theta_i l_i} = \tilde{\theta}_i$. We define $Z_{R,I}$ as a $|R| \times |I|$ matrix with $Z_{R,I}(r, i) = 1$ if long-read r is generated from a molecule that is originated from isoform i , and $Z_{R,I}(r, i) = 0$ otherwise. For isoform quantification, our goal is to estimate $\Theta = \{\theta_i, i \in I\}$ 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

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

This formula is based 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

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

where $q(r, i)$ is isoform-specific read quality score and $f(L_r \mid \text{iso.} = i)$ 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 $f(L_r \mid \text{iso.} = i)$

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 read, 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 [22, 29, 40]. To offer flexibility in modeling read length distribution, we employ a nonparametric approach. For all long-reads mapped to the genome, we categorize them into two groups: complete reads and truncated reads. The read is treated as complete when the distance between its ending alignment position and any known isoform 5' end is less than 20 bp (Figure 1(B)). 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 is 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. Function $\hat{F}(l) = P(\text{read len.} < l)$ can be estimated based on Kaplan-Meier estimator[41], 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

$$f(L_r | \text{iso.} = i) = \frac{f(L_r) \cdot P(\text{iso.} = i | L_r)}{P(\text{iso.} = i)} = \frac{f(L_r) \cdot \tilde{\theta}_i / \sum_{l_j > L_r} \tilde{\theta}_j}{\tilde{\theta}_i} = \frac{f(L_r)}{\sum_{l_j > L_r} \tilde{\theta}_j}$$

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, $\tilde{\theta}$ are estimated using EM algorithm. Then, we have isoform relative abundance $\theta_i = \frac{\tilde{\theta}_i / l_i}{\sum_{i \in I} \tilde{\theta}_i / l_i}$. 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 10K reads (RPG 10K) 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 ($\text{RPG}_i = N \cdot \tilde{\theta}_i / 10^4$).

Parameter estimation using EM algorithm

Our interested parameter θ will be estimated by inferring $\tilde{\theta}$ with the fact that $\theta_i = \frac{\tilde{\theta}_i / l_i}{\sum_{i \in I} \tilde{\theta}_i / l_i}$. The complete data likelihood is

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

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

E-step: We calculate function

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

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

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

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

The EM algorithm consists of alternating between the E- and M-steps until convergence. We start the algorithm with $\tilde{\theta}^{(0)}$ 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 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

$$\Phi_{I-1}(\Phi^{-1}(F(\theta_1|\beta_0, \beta_1, \varphi_1)), \dots, \Phi^{-1}(F(\theta_{I-1}|\beta_0, \beta_{I-1}, \varphi_{I-1}))|\Gamma)$$

where φ_i is the dispersion parameter of the marginal generalized linear model for isoform i . $\Phi_{I-1}(\Gamma)$ is the cumulative distribution function of multivariate normal random variables with $I - 1$ dimensions and correlation matrix Γ . We choose to use exchangeable correlation structure for Γ . Given regression models above, we can detect DAS both for at isoform level and at gene level. For isoform i , we test $H_0: \beta_0 = 0$ vs $H_1: \beta_0 \neq 0$ to determine splicing change between conditions. For gene g , we test $H_0: \beta_0 = \dots = \beta_{I-1} = 0$ vs $H_1: \beta_i \neq 0$ 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)[42-44] project has been utilized in many studies. For example, Gao et al[45] 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/isoform expressions from this sample were used to benchmark and optimize computational tools[14, 46-49]. 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[31], 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, Mandalorian, CEM, Cufflinks and RD.

Data availability

The direct mRNA sequencing data on UHR are available at BioProject data base (PRJNA639366). The cDNA sequencing data on a patient with cancer are available at BioProject data base (PRJNA639366). The simulation data used in our study can be reproduced using code provided in the LIQA software repository and NanoSim version 2.0.0.

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 UHR datasets for making the data publicly available for benchmarking studies.

Funding

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

References

1. Han, J., et al., *Pre-mRNA splicing: where and when in the nucleus*. Trends Cell Biol, 2011. **21**(6): p. 336-43.
2. Scotti, M.M. and M.S. Swanson, *RNA mis-splicing in disease*. Nat Rev Genet, 2016. **17**(1): p. 19-32.
3. Montes, M., et al., *RNA Splicing and Disease: Animal Models to Therapies*. Trends Genet, 2019. **35**(1): p. 68-87.
4. Li, Y.L., et al., *RNA splicing is a primary link between genetic variation and disease*. Science, 2016. **352**(6285): p. 600-4.
5. Kim, H.K., et al., *Alternative splicing isoforms in health and disease*. Pflugers Arch, 2018. **470**(7): p. 995-1016.
6. Li, B. and C.N. Dewey, *RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome*. BMC Bioinformatics, 2011. **12**: p. 323.
7. Roberts, A. and L. Pachter, *Streaming fragment assignment for real-time analysis of sequencing experiments*. Nat Methods, 2013. **10**(1): p. 71-3.
8. Nariai, N., et al., *TIGAR2: sensitive and accurate estimation of transcript isoform expression with longer RNA-Seq reads*. BMC Genomics, 2014. **15 Suppl 10**: p. S5.
9. Zhang, C., et al., *Evaluation and comparison of computational tools for RNA-seq isoform quantification*. BMC Genomics, 2017. **18**(1): p. 583.
10. Patro, R., S.M. Mount, and C. Kingsford, *Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms*. Nat Biotechnol, 2014. **32**(5): p. 462-4.
11. Bray, N.L., et al., *Near-optimal probabilistic RNA-seq quantification*. Nat Biotechnol, 2016. **34**(5): p. 525-7.
12. Trapnell, C., et al., *Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks*. Nat Protoc, 2012. **7**(3): p. 562-78.
13. Li, W. and T. Jiang, *Transcriptome assembly and isoform expression level estimation from biased RNA-Seq reads*. Bioinformatics, 2012. **28**(22): p. 2914-21.
14. Hu, Y., et al., *PennSeq: accurate isoform-specific gene expression quantification in RNA-Seq by modeling non-uniform read distribution*. Nucleic Acids Res, 2014. **42**(3): p. e20.
15. Nicolae, M., et al., *Estimation of alternative splicing isoform frequencies from RNA-Seq data*. Algorithms Mol Biol, 2011. **6**(1): p. 9.
16. Pertea, M., et al., *StringTie enables improved reconstruction of a transcriptome from RNA-seq reads*. Nat Biotechnol, 2015. **33**(3): p. 290-5.
17. Li, J.J., et al., *Sparse linear modeling of next-generation mRNA sequencing (RNA-Seq) data for isoform discovery and abundance estimation*. Proc Natl Acad Sci U S A, 2011. **108**(50): p. 19867-72.
18. Mezlini, A.M., et al., *iReckon: simultaneous isoform discovery and abundance estimation from RNA-seq data*. Genome Res, 2013. **23**(3): p. 519-29.
19. Wan, L., et al., *Modeling RNA degradation for RNA-Seq with applications*. Biostatistics, 2012. **13**(4): p. 734-47.
20. Burgess, D.J., *Genomics: Next generation sequencing for reference genomes*. Nat Rev Genet, 2018. **19**(3): p. 125.
21. Pollard, M.O., et al., *Long reads: their purpose and place*. Hum Mol Genet, 2018. **27**(R2): p. R234-R241.
22. Sharon, D., et al., *A single-molecule long-read survey of the human transcriptome*. Nat Biotechnol, 2013. **31**(11): p. 1009-14.

23. Tilgner, H., et al., *Comprehensive transcriptome analysis using synthetic long-read sequencing reveals molecular co-association of distant splicing events*. Nat Biotechnol, 2015. **33**(7): p. 736-42.
24. Treutlein, B., et al., *Cartography of neurexin alternative splicing mapped by single-molecule long-read mRNA sequencing*. Proc Natl Acad Sci U S A, 2014. **111**(13): p. E1291-9.
25. Vollmers, C., et al., *Novel exons and splice variants in the human antibody heavy chain identified by single cell and single molecule sequencing*. PLoS One, 2015. **10**(1): p. e0117050.
26. Oikonomopoulos, S., et al., *Benchmarking of the Oxford Nanopore MinION sequencing for quantitative and qualitative assessment of cDNA populations*. Sci Rep, 2016. **6**: p. 31602.
27. Byrne, A., et al., *Nanopore long-read RNAseq reveals widespread transcriptional variation among the surface receptors of individual B cells*. Nat Commun, 2017. **8**: p. 16027.
28. Tang, A.D., et al., *Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals downregulation of retained introns*. Nat Commun, 2020. **11**(1): p. 1438.
29. Amarasinghe, S.L., et al., *Opportunities and challenges in long-read sequencing data analysis*. Genome Biol, 2020. **21**(1): p. 30.
30. Yang, C., et al., *NanoSim: nanopore sequence read simulator based on statistical characterization*. Gigascience, 2017. **6**(4): p. 1-6.
31. Li, H., *Minimap2: pairwise alignment for nucleotide sequences*. Bioinformatics, 2018. **34**(18): p. 3094-3100.
32. Tang, A.D., et al., *Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals downregulation of retained introns*. bioRxiv, 2018: p. 410183.
33. De Kouchkovsky, I. and M. Abdul-Hay, 'Acute myeloid leukemia: a comprehensive review and 2016 update'. Blood Cancer J, 2016. **6**(7): p. e441.
34. Cheng, Y.W., et al., *Long Read Single-Molecule Real-Time Sequencing Elucidates Transcriptome-Wide Heterogeneity and Complexity in Esophageal Squamous Cells*. Front Genet, 2019. **10**: p. 915.
35. Ueno, N., et al., *Enhanced Expression of Fibroblast Growth Factor Receptor 3 Il1c Promotes Human Esophageal Carcinoma Cell Proliferation*. J Histochem Cytochem, 2016. **64**(1): p. 7-17.
36. Dobin, A., et al., *STAR: ultrafast universal RNA-seq aligner*. Bioinformatics, 2013. **29**(1): p. 15-21.
37. Wu, T.D. and C.K. Watanabe, *GMAP: a genomic mapping and alignment program for mRNA and EST sequences*. Bioinformatics, 2005. **21**(9): p. 1859-75.
38. Kent, W.J., *BLAT--the BLAST-like alignment tool*. Genome Res, 2002. **12**(4): p. 656-64.
39. Maric, J., et al., *Graphmap2-splice-aware RNA-seq mapper for long reads*. bioRxiv, 2019: p. 720458.
40. Kellner, S., J. Burhenne, and M. Helm, *Detection of RNA modifications*. RNA Biol, 2010. **7**(2): p. 237-47.
41. Kaplan, E.L. and P. Meier, *Nonparametric estimation from incomplete observations*. Journal of the American statistical association, 1958. **53**(282): p. 457-481.
42. Consortium, M., et al., *The MicroArray Quality Control (MAQC) project shows inter- and intraplatform reproducibility of gene expression measurements*. Nat Biotechnol, 2006. **24**(9): p. 1151-61.
43. Shi, L., 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**(8): p. 827-38.
44. Sun, P., et al., *Expression of estrogen receptor-related receptors, a subfamily of orphan nuclear receptors, as new tumor biomarkers in ovarian cancer cells*. Journal of Molecular Medicine, 2005. **83**(6): p. 457-467.

45. Gao, F., et al., *Evaluation of biological and technical variations in low-input RNA-Seq and single-cell RNA-Seq*. International Journal of Computational Biology and Drug Design, 2018. **11**(1-2): p. 5-22.
46. Xu, J., et al., *Cross-platform ultradeep transcriptomic profiling of human reference RNA samples by RNA-Seq*. Sci Data, 2014. **1**: p. 140020.
47. Garcia-Alonso, L., et al., *Benchmark and integration of resources for the estimation of human transcription factor activities*. Genome Res, 2019. **29**(8): p. 1363-1375.
48. Teng, M., et al., *A benchmark for RNA-seq quantification pipelines*. Genome Biol, 2016. **17**: p. 74.
49. Hayer, K.E., et al., *Benchmark analysis of algorithms for determining and quantifying full-length mRNA splice forms from RNA-seq data*. Bioinformatics, 2015. **31**(24): p. 3938-45.

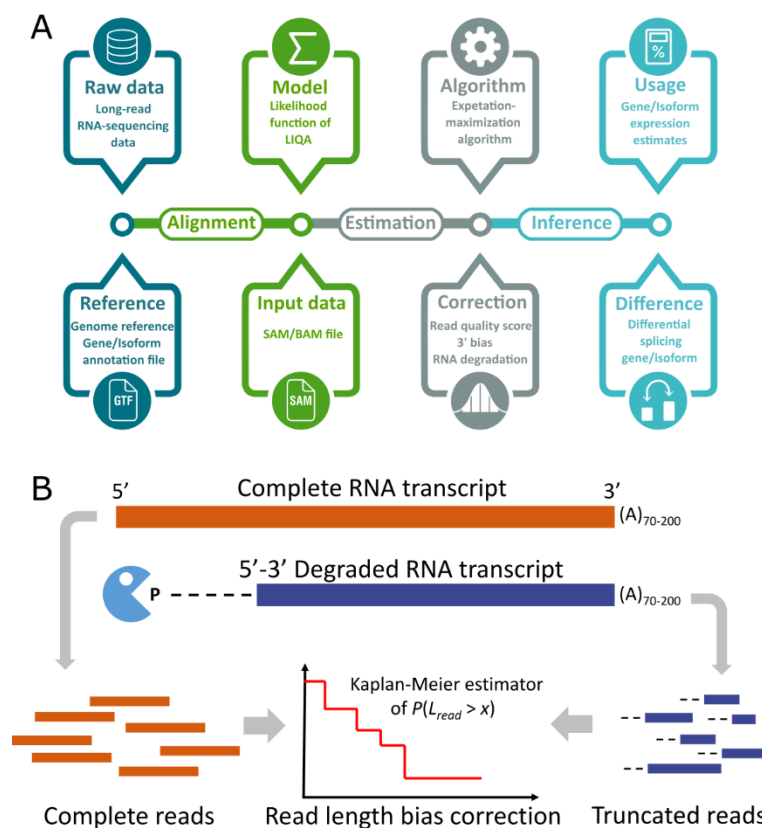


Figure 1. Framework of LIQA. (A) The flowchart to illustrate how LIQA works. The input of LIQA are long-read RNA-seq data and isoform annotation file. LIQA models observed splicing information, high error rate of data and read length bias. The output of LIQA are isoform expression estimates and detected DAS event (B) Quantification of potential 3' bias of long-read RNA-seq data. Complete and degraded RNA transcript are indicated by orange and blue. Complete (orange) and truncated (blue) long reads are jointly modeled to correct read length bias by estimating read length distribution.

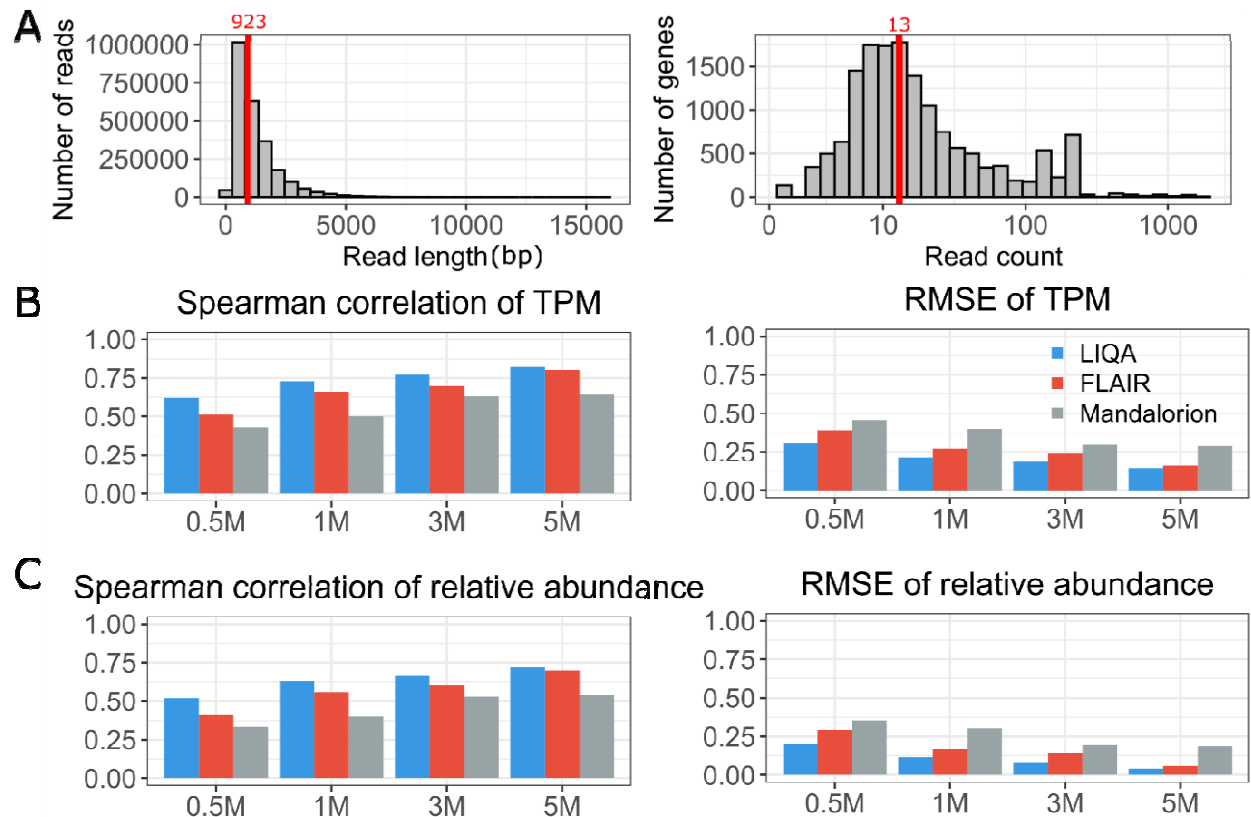


Figure 2. Simulation study results. (A) Characteristics of simulated data with 5M. Read length distribution (left) and read count distribution by genes in log-scale (right). (B-C) Summary statistics between true and estimated isoform expressions using LIQA (blue), FLAIR (red) and Mandalorion (gray) at different read coverages. (B) Spearman correlations (left) and RMSE (right) between estimated and true TPM. (C) Spearman correlations (left) and RMSE (right) between estimated and true relative abundance.

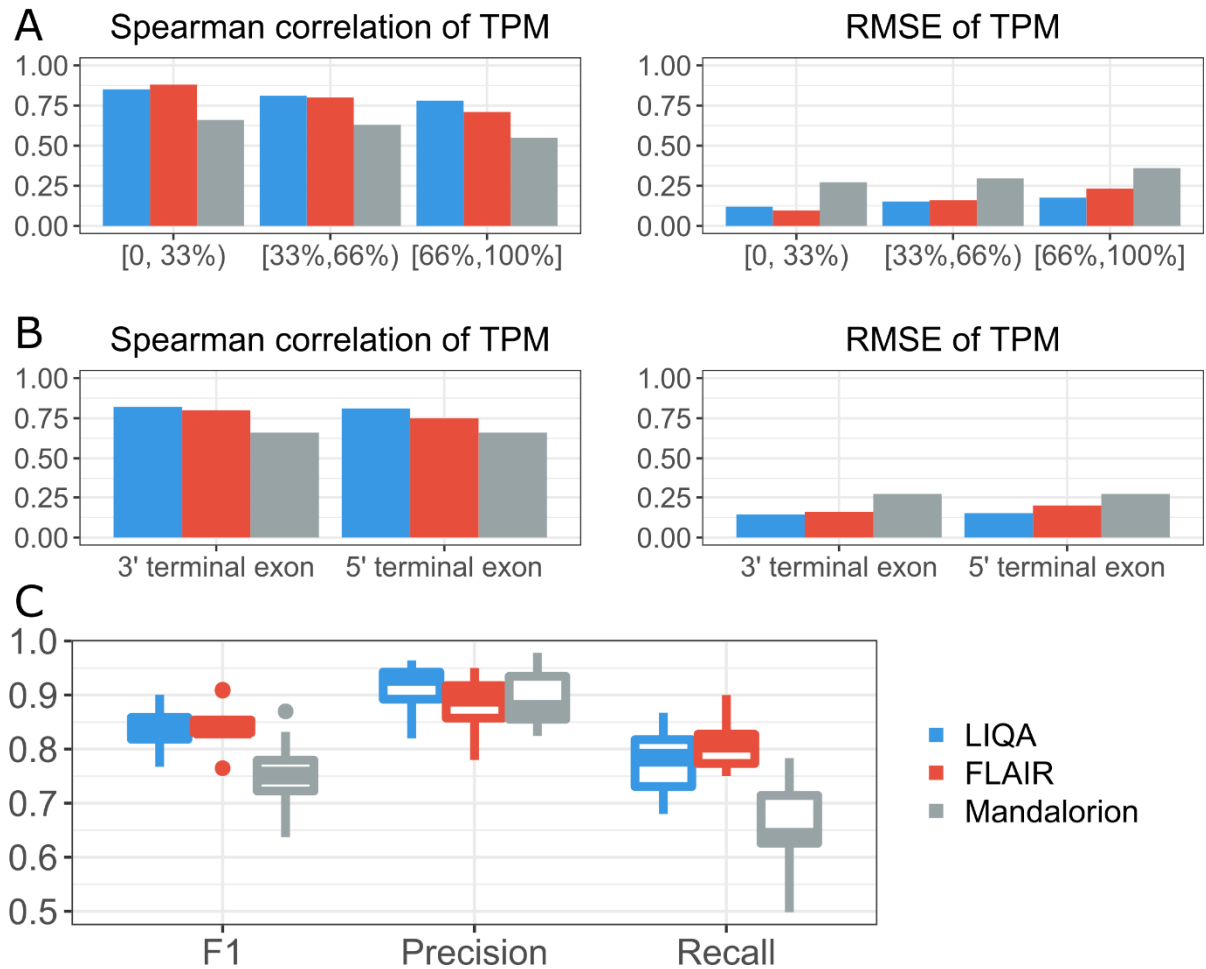


Figure 3. Evaluation of robustness to read coverage bias (A-B) and DAS events detection (C) using LIQA (blue), FLAIR (red) and Mandalorion (gray). (A) Spearman correlations (left) and RMSE (right) between estimated and true TPM for isoforms with different lengths (length < 33% quantile, 33% ≤ length < 66% quantile, length ≥ 66%). (B) Spearman correlations (left) and RMSE (right) between estimated and true TPM for 3' and 5' terminal exons. (C) Summary statistics (recall, precision and F1 score) of DAS event detection analysis for 10 simulations.

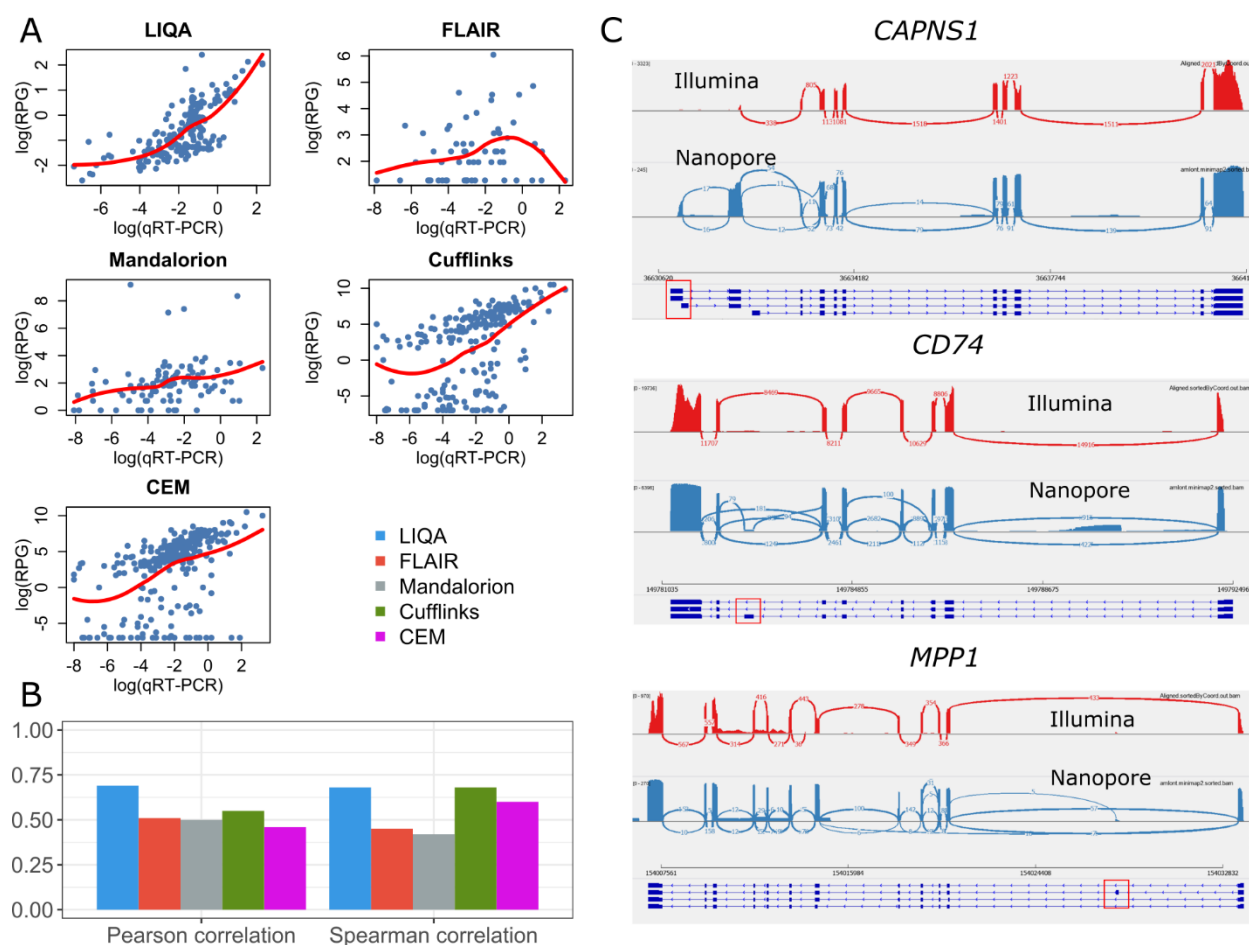


Figure 4. UHR data analysis results. (A-B) Comparison of different methods. (A) Scatter plots of estimated isoform-specific expression versus qRT-PCR measurements in log scale. (B) Pearson correlation coefficients (left) and spearman correlation coefficient (right) between estimated isoform-specific expression versus qRT-PCR measurements in log scale. (C) Examination of read coverage difference between Illumina and Nanopore data at randomly selected 3 genes. Informative exonic regions were in red square.

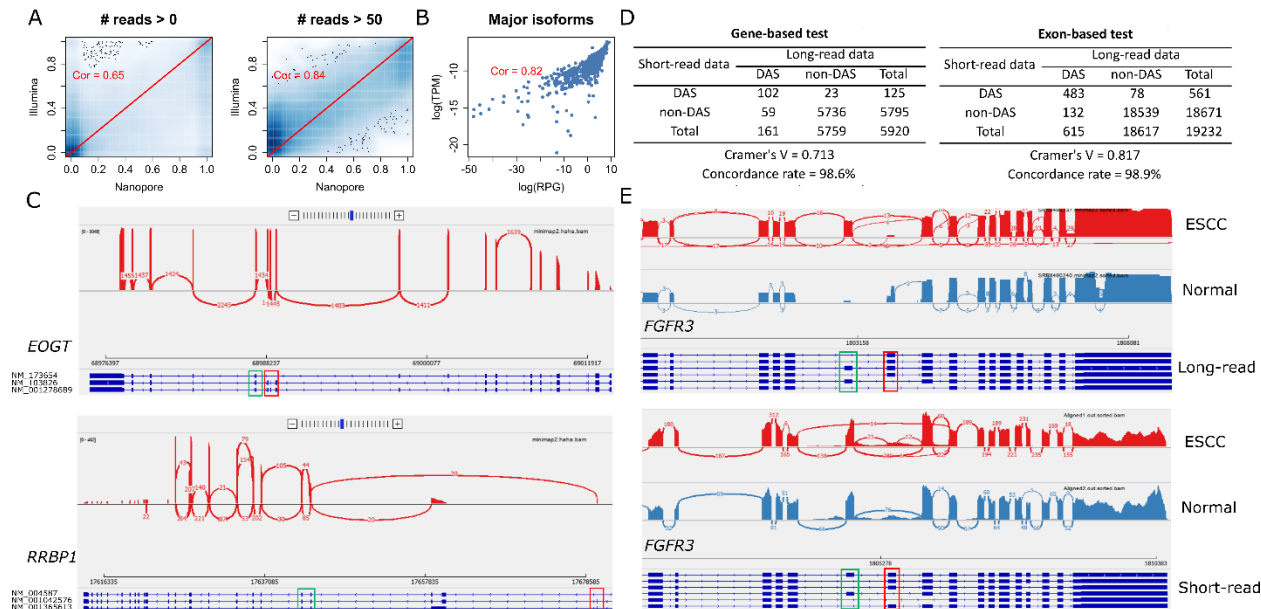


Figure 5. Performance of LIQA using AML data (A-C) and ESCC data (D, E). (A) Scatter plots of estimated isoform relative abundances using long-read data (LIQA) versus short-read data (PennSeq) for all genes (left) and genes with at least 50 read coverage (right). (B) Scatter plot of estimated isoform-specific expression using long-read data (LIQA) versus short-read data (PennSeq) in log-scale for all major isoforms. (C) Examination of isoform usage inferred by LIQA. Sashimi plots of gene *EOGT* and *RRBP1*. Informative exonic regions were in green and red squares. (D) DAS detections between long- and short-read data. Consistency of detected DAS events between long- and short-read data were quantified using Cramer's V and concordance rate. (E) Examination of AS exon usage inferred by LIQA (long-read) but missed by PennDiff (short-read). Sashimi plots of gene *FGFR3*. Informative exonic regions were in green and red squares.