

# IsoQuant: a tool for accurate novel isoform discovery with long reads

**Andrey Prjibelski**

Saint Petersburg University

**Alla Mikheenko**

Center for Algorithmic Biotechnology, Institute of Translational Biomedicine, Saint Petersburg State University

**Anoushka Joglekar**

Weill Cornell Medicine <https://orcid.org/0000-0002-7818-6867>

**Alexander Smetanin**

Bioinformatics Institute

**Julien Jarroux**

Weill Cornell Medicine

**Alla Lapidus**

St.Petersburg State University <https://orcid.org/0000-0003-0427-8731>

**Hagen Tilgner** (✉ [hagen.u.tilgner@gmail.com](mailto:hagen.u.tilgner@gmail.com))

Weill Cornell Medicine <https://orcid.org/0000-0002-7058-3606>

---

## Brief Communication

### Keywords:

**Posted Date:** April 26th, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1571850/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

# 1 IsoQuant: a tool for accurate novel isoform discovery 2 with long reads

3 *Andrey D. Prjibelski<sup>1,+</sup>, Alla Mikheenko<sup>1</sup>, Anoushka Joglekar<sup>2,3,4</sup>, Alexander Smetanin<sup>5</sup>, Julien Jarroux<sup>3,4</sup>,  
4 Alla L. Lapidus<sup>1</sup>, Hagen U. Tilgner<sup>3,4+</sup>*

5 <sup>1</sup> *Center for Algorithmic Biotechnology, Institute of Translational Biomedicine, St. Petersburg State  
6 University*

7 <sup>2</sup> *Tri-Institutional Computational Biology & Medicine, Weill Cornell Medicine, New York, NY, USA*

8 <sup>3</sup> *Brain and Mind Research Institute, Weill Cornell Medicine, New York, NY, USA*

9 <sup>4</sup> *Center for Neurogenetics, Weill Cornell Medicine, New York, NY, USA*

10 <sup>5</sup> *Bioinformatics Institute, St. Petersburg, Russia*

11 + *Corresponding authors: [hui2006@med.cornell.edu](mailto:hui2006@med.cornell.edu), [andrewprzh@gmail.com](mailto:andrewprzh@gmail.com)*

12

13 **Long reads are reshaping RNA biology. However, determining alternative isoforms from**  
14 **long-read RNA data is a complex and incompletely solved problem even when the reference**  
15 **genome is known. Here we present IsoQuant — a reference-based tool that accurately**  
16 **discovers novel transcripts with at least 3-fold lower false positive rate and 1.8-fold increase**  
17 **in F1-score compared to other tools for Oxford Nanopore data. IsoQuant also increases**  
18 **performance for Pacific Biosciences data.**

19

20 Long-read RNA sequencing is now widely used in bulk, sorted cells, single cells and spatial  
21 approaches. This wide field of applications has led to the development of multiple spliced  
22 alignment programs<sup>1-4</sup>, transcript discovery methods<sup>5-11</sup>, tools for transcript classification<sup>12</sup>,  
23 annotation<sup>13</sup> and visualization<sup>14,15</sup>. Additionally, several reference-free tools for RNA long-read  
24 correction and assembly have been developed<sup>16,17</sup>. Current community efforts address the problem  
25 of understanding performance, weaknesses and advantages of each approach for various  
26 applications<sup>18</sup>.

27

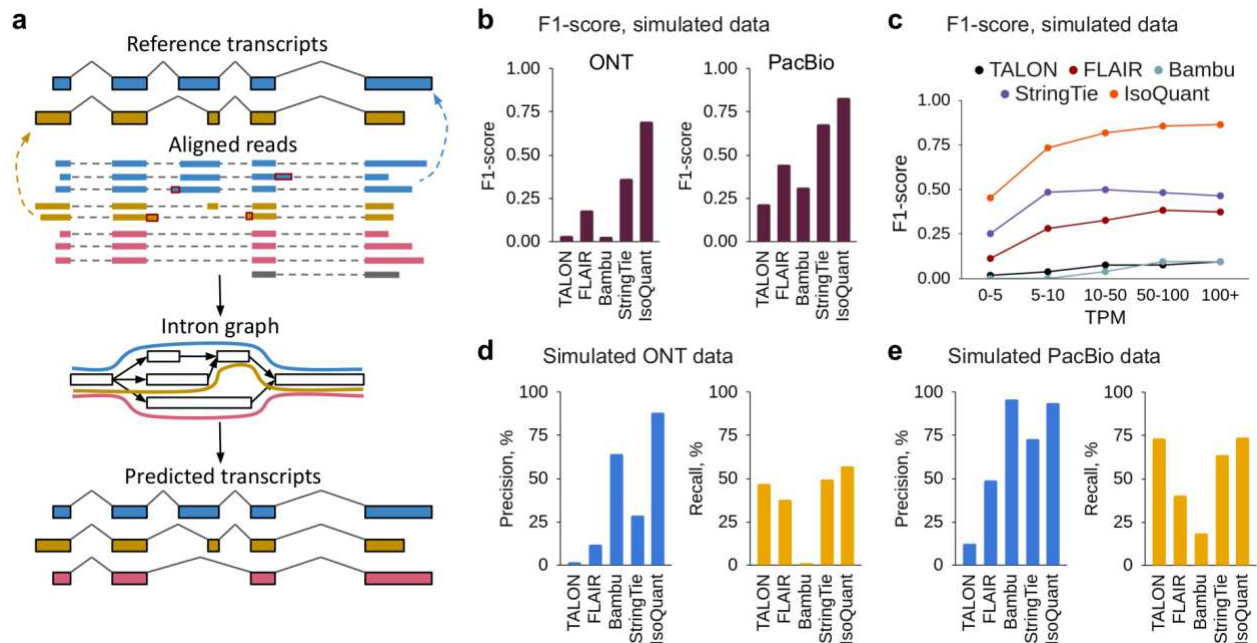
28 Here we present IsoQuant — a tool for transcript discovery and quantification with long RNA  
29 reads. IsoQuant takes as input a reference genome, a gene annotation, and a dataset containing  
30 PacBio or ONT RNA reads. By default, IsoQuant maps input reads to the genome via minimap2  
31 in splice mode<sup>2</sup>. Alternatively, a user may provide BAM files generated with a spliced aligner of  
32 their choice, e.g. STARlong<sup>1</sup> for PacBio and uLTRA<sup>4</sup> or deSALT<sup>3</sup> for ONT reads.

33

34 IsoQuant first assigns reads to known isoforms via an inexact intron chain matching algorithm that  
35 takes into account splice site shifts, which are typical for alignment of error-prone reads<sup>19</sup>.  
36 Uniquely assigned reads with consistent intron chains are then used for isoform and gene  
37 quantification (Methods). For inconsistent reads that are likely to originate from unannotated  
38 isoforms, IsoQuant also reports newly detected alternative splicing events. IsoQuant exploits read-

1 to-isoform assignments for transcripts quantification and correction of inaccurately detected splice  
2 junctions and misalignments, such as skipped microexons. Corrected alignments are used to  
3 construct an intron graph, in which vertices are introns, and two vertices are connected with a  
4 directed edge if the corresponding introns are consecutive in at least one read (Methods). Finally,  
5 this graph is exploited for constructing paths that correspond to full-length transcripts (Figure 1a).  
6  
7 To compare IsoQuant performance against existing transcript discovery tools, we first simulated  
8 mouse PacBio and ONT data using realistic gene expression profiles with IsoSeqSim  
9 (<https://github.com/yunhaowang/IsoSeqSim>) and Trans-NanoSim<sup>20</sup> respectively. To mimic real-  
10 life datasets containing unannotated transcripts, we arbitrarily removed 5,311 (15%) of 35,684  
11 expressed isoforms (the ones contributing to at least one read during the simulation) from the  
12 GENCODE gene annotation. These 5,311 hidden transcripts were further used as a ground truth  
13 for novel transcript discovery. The reduced GENCODE annotation was used as an input for all  
14 tools. Each output annotation was then separated into a set of known and a set of novel transcripts,  
15 which were compared against the respective baselines using gffcompare<sup>21</sup> (Methods).  
16  
17 For known transcripts, IsoQuant has the highest F1-score (the harmonic mean of precision and  
18 recall) compared to TALON<sup>7</sup>, FLAIR<sup>8</sup>, Bambu<sup>11</sup> and StringTie<sup>5</sup>, but these advances are not  
19 dramatic (Supplemental Tables 1-3). However, IsoQuant produces novel transcripts with a 1.8-  
20 fold higher F1-score on ONT data compared to the second best tool, StringTie. In comparison to  
21 TALON, FLAIR and Bambu, the improvement in F1-score is even more noticeable (Figure 1b  
22 left). On PacBio data, IsoQuant again shows the best F1-score, but the difference from other tools  
23 is smaller than for ONT data (Figure 1b, right). While IsoQuant's F1-score is clearly higher than  
24 that of other tools for highly-expressed transcripts, IsoQuant is able to maintain these advances for  
25 novel transcript discovery regardless of the expression levels (Figure 1c, Supplementary Figure  
26 1). Thus, IsoQuant is likely to be highly useful across many genes, including but not limited to  
27 lowly expressed long-non-coding RNAs and marker genes of cell types.  
28  
29 Compared to most tools, IsoQuant's improvements in F1-score is primarily caused by its very high  
30 precision of novel transcripts. As compared to TALON, FLAIR and StringTie, IsoQuant shows a  
31 minimum of 6-fold drop in false-positive rate on ONT data, while still maintaining slight gains in  
32 recall (Figure 1d). Importantly, the situation is of a different nature for Bambu. IsoQuant has higher  
33 precision (88.2% vs. 64.2%), but substantially higher recall: while Bambu only reconstructs 70 out  
34 of 5,311 novel isoforms (1.3% recall), IsoQuant reconstructs 3,031 (57.2%).  
35  
36 For novel transcripts generated with PacBio data, similar trends can be observed, although with a  
37 less drastic difference in specificity. Bambu shows slightly higher precision (95.8%) compared to  
38 IsoQuant (94.0%), but again has the lowest recall (18.7% for Bambu vs. 74.2% for IsoQuant).  
39 StringTie, TALON and FLAIR again predict transcripts with comparable recall, but have at least  
40 4.5-fold higher false-positive rate compared to IsoQuant (Figure 1e).

1



2

3 **Figure 1. IsoQuant pipeline outline and characteristics of novel transcripts generated from mouse**  
 4 **simulated data.** **a.** Outline of the IsoQuant pipeline. Reads are assigned to annotated isoforms and  
 5 alignment artifacts are corrected (top). The intron graph is constructed from corrected read alignments  
 6 (middle) and transcripts are discovered via path construction (bottom). **b.** F1-score for novel transcripts  
 7 reported by different tools on simulated ONT (left) and PacBio data (right). **c.** F1-score for novel transcripts  
 8 reported by different tools on simulated ONT data broken up by expression levels in transcript-per-million  
 9 (TPM). **d.** Precision (left) and recall (right) for novel transcripts reported by different tools on simulated  
 10 ONT data. **e.** Same as d, but for simulated PacBio data.

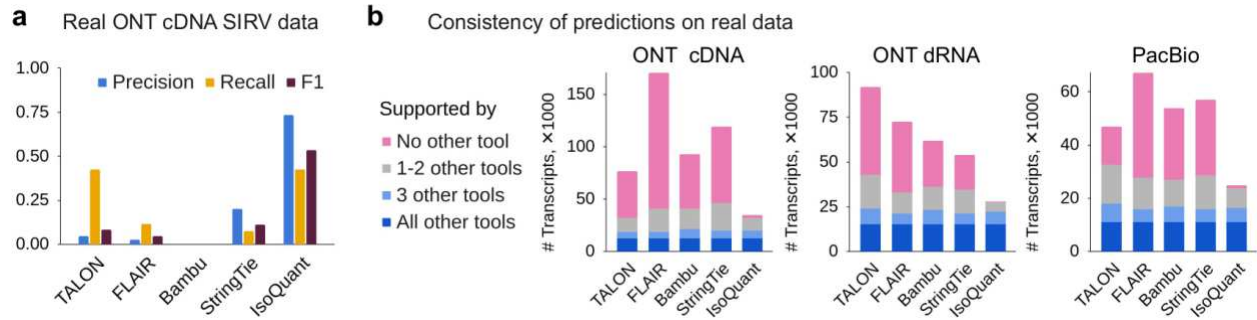
11

12 To complement our benchmarks on simulated data, we also sequenced Lexogen SIRVs synthetic  
 13 molecules on the Oxford Nanopore MinION (Methods). Along with the complete SIRV  
 14 annotation, Lexogen provides an incomplete annotation, missing 26 out of the total 69 SIRV  
 15 isoforms, which allows the evaluation of novel transcript discovery, similar to the one we  
 16 performed for simulated data with the reduced GENCODE annotation.

17

18 Results on SIRV sequencing data resemble the ones obtained on simulated reads. When predicting  
 19 novel isoforms, IsoQuant shows at least 5 times higher F1-score and 3-fold lower false positive  
 20 rate than any other tool. In comparison to most tools, with the exception of TALON, IsoQuant  
 21 shows high gains in both precision and recall. TALON and IsoQuant have identical recall of  
 22 42.3%, but IsoQuant has 16-fold higher precision (Figure 2a). Similar to simulated data, all tools  
 23 are able to accurately predict SIRV transcripts kept in the annotation, with Bambu, StringTie and  
 24 IsoQuant having perfect precision for known isoforms alone (Supplementary Table 4).

25



**Figure 2. Characteristics of transcripts obtained from real sequencing data. a.** Precision, recall, and F1-score for novel transcripts generated on real SIRV ONT cDNA sequencing data. **b.** Consistency of predictions made by different methods on real human ONT cDNA, ONT dRNA and PacBio data.

To support our observations, we also applied all tools to the real human ONT cDNA, ONT dRNA<sup>22</sup>, and PacBio public datasets, for which the ground truth is indeed unknown. We used gffcompare to estimate the consistency of predictions by computing the number of identical transcript models reported by the different tools. On the human ONT cDNA dataset, IsoQuant shows the highest percentage of transcripts confirmed by at least 3 other methods (58.1%), while no other tool surpasses the 25% threshold. This suggests that IsoQuant transcript models are significantly more consistent with other methods (Figure 2b left). In comparison to the other approaches, IsoQuant also reports the lowest number of transcripts that are not predicted by any other method. If one interprets such transcript models as false positives, IsoQuant again stands out in the lowest false-discovery rate (7.7%, 2,655 transcripts). In contrast, other tools output annotations containing more than 50% of unconfirmed transcript models (varying from 44,000 to 130,000). Additionally, for each tool we computed the number of potentially missed transcripts that were reported by all other methods. While Bambu has the lowest number of such transcripts (574), IsoQuant shows the second-best results of 1507 possible false negatives (Supplementary Table 5).

Similar trends can be observed in ONT dRNA and PacBio datasets, although the overall percentage of common transcripts appears to be higher compared to ONT cDNA data (Figure 2b, middle and right). IsoQuant again shows the highest fraction of transcripts predicted by at least 3 other tools (79.8% for ONT dRNA, 65.9% for PacBio), while other programs have 40% at best. Of the others, all four tools produce annotations containing more than 30% of transcripts that are not confirmed by any other method, while IsoQuant's potential false predictions are below 5% on both datasets.

Although these values cannot be explicitly treated as false positives and false negatives, they advocate that unlike other tools, IsoQuant produces highly specific annotations that are strongly consistent with transcripts reported by several alternative approaches. Moreover, because IsoQuant typically misses very few isoforms predicted by all other tools simultaneously, it is likely to also be highly sensitive.

1 Additionally, we used long-read RNA sequencing data from a mouse brain sample, wherein a  
2 previous study reported 76 novel isoforms of high biological importance<sup>23</sup>, which were confirmed  
3 by manual annotation by the GENCODE team. Here, we compared IsoQuant only with StringTie,  
4 which has the second-best F1-score across all simulated datasets. On PacBio data, IsoQuant  
5 correctly reconstructs 71% of the confirmed novel isoforms, while StringTie restores half as many  
6 novel transcripts — only 37% (Supplementary Table 6). Similarly, on two ONT datasets (single-  
7 cell and spatial) from the same brain sample IsoQuant restores up to 50% of these 76 novel  
8 isoforms, whereas StringTie reports 30% at best. Although it is not possible to evaluate specificity  
9 in this kind of experiment, it confirms that IsoQuant is capable of maintaining high recall values  
10 on real sequencing data.

11  
12 Beside transcript discovery, IsoQuant implements additional functionality, such as read-to-isoform  
13 assignment and transcript quantification. Benchmarks of these supplementary features,  
14 information on computational performance, as well as IsoQuant results obtained with different  
15 spliced aligners can be found in the Supplementary Notes 2-5.

16  
17 In summary, IsoQuant accurately predicts transcript models from PacBio or ONT RNA sequencing  
18 data. For known isoforms, IsoQuant has higher F1-score compared to other tested tools, but these  
19 differences are not dramatic. For unannotated isoforms, however, IsoQuant provides very strong  
20 increases in F1-score over other existing approaches. In comparison to most tools, it achieves this  
21 F1-score increase by maintaining higher recall, while substantially increasing precision. Thus,  
22 IsoQuant is a valuable tool for predicting novel alternatively spliced isoforms in the age of long-  
23 read sequencing.

## 24 25 **References**

- 26 1. Dobin, A. et al. *Bioinformatics* 29, 15–21 (2013).
- 27 2. Li, H. *Bioinformatics* 34, 3094–3100 (2018).
- 28 3. Liu, B. et al. *Genome Biol.* 20, 274 (2019).
- 29 4. Sahlin, K. & Mäkinen, V. *Bioinformatics* 37, 4643-4651 (2021).
- 30 5. Kovaka, S. et al. *Genome Biol.* 20, 278 (2019).
- 31 6. Tung, L.H., Shao, M. & Kingsford, C. *Genome Biol.* 20, 287 (2019).
- 32 7. Wyman, D. et al. *BioRxiv*, <https://doi.org/10.1101/672931>
- 33 8. Tang, A.D. et al. *Nat. Commun.* 11, 1438 (2020).
- 34 9. Kuo, R.I. et al. *BMC Genomics* 21, 751 (2020).
- 35 10. Byrne, A. et al. *Nat. Commun.* 8, 16027 (2017).
- 36 11. Chen, Y. et al. *Bioconductor*, <https://doi.org/doi:10.18129/B9.bioc.bambu>
- 37 12. Tardaguila, M. et al. *Genome Res.* 28, 396-411 (2018).
- 38 13. de la Fuente, L. et al. *Genome Biol.* 21, 119 (2020).
- 39 14. Reese, F. & Mortazavi, A. doi:10.1101/2020.06.09.143024
- 40 15. Stein, A.N. et al. *BioRxiv*, <https://doi.org/10.1101/2022.04.14.488347>
- 41 16. Sahlin, K. & Medvedev, P. *Nat. Commun.* 12, 2 (2021).
- 42 17. Nip, K.M. et al. *Genome Res.* 30, 1191–1200 (2020).

- 1 18. Pardo-Palacios, F. et al. *In Review*, <https://doi.org/10.21203/rs.3.rs-777702/v1>
- 2 19. Mikheenko, A. et al. *Genome Res.* 32, 726–737 (2022).
- 3 20. Hafezqorani, S. et al. *Gigascience* 9, (2020).
- 4 21. Pertea, G. & Pertea, M. *F1000Res.* 9, (2020).
- 5 22. Workman, R.E. et al. *Nat. Methods* 16, 1297–1305 (2019).
- 6 23. Joglekar, A. et al. *Nat. Commun.* 12, 463 (2021).

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

# 1 **Methods**

2 **Data simulation.** To simulate PacBio CCS reads we used IsoSeqSim  
3 (<https://github.com/yunhaowang/IsoSeqSim>), which generates a read by truncating a transcript  
4 sequence according to given probabilities and randomly inserts sequencing errors at a specified  
5 rate with uniform distribution. As reported in previous studies<sup>24</sup>, a uniform error distribution is a  
6 realistic model for PacBio CCS reads. Here we used 5' and 3' truncation probabilities typical for  
7 PacBio Sequel II (provided within the package) and an overall error rate of 1.6%: 0.6% deletions,  
8 0.6% insertions and 0.4% substitutions. While these discrepancies do not necessarily represent  
9 sequencing errors, they nevertheless have to be modeled, as they can confuse transcript  
10 reconstruction. The above values were obtained by mapping real PacBio CCS reads to the  
11 reference genome<sup>18</sup>.

12  
13 ONT reads were simulated with the NanoSim software in the transcriptome mode<sup>20</sup>. NanoSim is  
14 designed specifically for simulating ONT-specific sequencing errors and biases. It first constructs  
15 error-profile and length-distribution models, which are further used to mutate reference transcript  
16 sequences. Here we used a pre-trained model provided within the NanoSim package, which was  
17 obtained using publicly available ONT cDNA data<sup>22</sup> from NA12878 human cell line and has an  
18 average error rate of 15.9%: 6% deletions, 5.1% insertions, 4.8% substitutions. In addition, we  
19 turned off the simulation of intron retention events and random unaligned reads representing the  
20 background noise.

21  
22 However, additional analysis of the simulated ONT data and NanoSim code revealed that  
23 NanoSim randomly selects a start position of a read in a transcript sequence with a uniform  
24 distribution, thus introducing no 5' or 3' bias. To simulate more realistic ONT reads, we aligned  
25 real ONT cDNA data obtained from the mouse brain sample to the reference transcriptome using  
26 minimap2 and derived empirical truncation probability distributions on both 5' and 3' ends.  
27 Further, we changed the NanoSim source code to enable sequence truncation with respect to  
28 obtained probabilities (Supplementary Figure 2). The modified version is available at  
29 <https://github.com/andrewprzh/lrgasp-simulation>.

30  
31 For both ONT and PacBio simulation we used Mouse GENCODE v26 and Human GENCODE  
32 v36 basic annotations<sup>25</sup>. Prior to simulation, we also attached a 30 bp polyA tail to every transcript  
33 sequence. To simulate realistic mouse data, a transcript expression profile was obtained using  
34 PacBio data from a mouse brain sample<sup>23</sup>. For human data, a gene expression profile was computed  
35 with PacBio GM12878 data. Complete description of every dataset used in this study is provided  
36 in the Supplementary Table 7.

37  
38 **Sequencing Lexogen SIRV transcripts.** First, total RNA from HeLa cells was extracted using  
39 the miRNeasy Tissue/Cells Advanced Mini Kit (Qiagen, 217604), and polyA transcripts were  
40 pulled-down using the NEBNext Poly(A) mRNA Magnetic Isolation Module (NEB, E7490S).



1 Next, the SIRV-Set 4 (Iso Mix E0 / ERCC / Long SIRVs) (Lexogen, 141.01) was spiked-in to the  
2 RNA and reverse transcribed using the Maxima H Minus RT (Thermo Scientific, EP0752). The  
3 RT reaction final concentrations are as follows: 1.25 ng/μl polyA HeLa RNA, 0.33 ng/μl SIRV-  
4 Set 4, 0.5 mM dNTP, 5 μM dT-VN oligo, 5 μM TSO, 1X RT Buffer, 2 U/μl RiboLock RNase  
5 Inhibitor (Thermo Scientific, EO0382) and 20 U/μl Maxima H Minus Reverse Transcriptase. The  
6 reaction was incubated for 30 minutes at 50°C and 5 minutes at 85°C. Then, 5 μl of RT reaction  
7 were amplified using the Platinum Superfi II Mastermix (ThermoFisher, 12368010) for 12 cycles,  
8 according to manufacturer instructions and using Forward- and Reverse-Amplification primers.  
9 Finally, the cDNA was cleaned up using SPRIselect beads at a 0.8x ratio (Beckman Coulter,  
10 B23318) and used as input for Amplicons by Ligation library preparation (Oxford Nanopore  
11 Technology, SQK-LSK110). Sequencing was run for 72 hours on a R9.4.1 MinION flowcell  
12 (ONT, FLO-MIN106D).

13

14 **Quality evaluation of predicted novel transcripts.** To mimic real-life situations and assess the  
15 ability of an algorithm to predict novel transcripts, we created reduced gene annotations by  
16 removing a fraction of expressed isoforms. First, we define a subset of true expressed transcripts  
17 that contributed to at least one read during the simulation. Among this set, we select a fraction of  
18 transcripts to be excluded from the annotation. These transcripts are denoted as the true novel  
19 isoforms. The remaining transcripts (among the expressed) are defined as true known isoforms.  
20 To create a reduced gene annotation we remove all true novel isoforms from the comprehensive  
21 GENCODE annotation. Here we created a reduced mouse annotation with 15% of expressed  
22 transcripts removed, and four human reduced annotations with 10%, 15%, 20% and 25% of  
23 excluded expressed isoforms (Supplementary Note 1).

24

25 To evaluate a transcript prediction tool, we provided the entire set of simulated reads and the  
26 reduced annotation as an input. Thus, true novel isoforms are hidden from the annotation, but  
27 present in the reads. We then compute precision and recall by running gffcompare<sup>21</sup> for (i) the  
28 entire output annotation vs. the complete set of expressed transcripts, (ii) reported known isoforms  
29 vs. the set of true known isoforms, and (iii) predicted novel transcript models vs. the true novel  
30 set. The information on whether a transcript is known or novel is obtained from the output GTF  
31 file. The script for computing these metrics can be found in the IsoQuant repository in  
32 misc/reduced\_db\_gffcompare.py.

33

34 To estimate how recall and precision of novel transcripts depend on the expression levels,  
35 predicted transcripts are grouped into bins by their TPM values. For computing recall, TPM values  
36 used during the simulation are used as the number of false negative calls is required (undetected  
37 transcripts). However, computing precision requires the number of false positive predictions  
38 within each bin and thus only reported TPM values can be used (the true TPM for a false prediction  
39 is 0). Thus, it may happen that the same transcript may fall into different bins when benchmarking  
40 different tools. Although it is not possible to compute precision, recall and F1-score exactly for an

1 arbitrary TPM range, the bias has a minor effect as only a small number of bins was used in this  
2 experiment (5). Therefore, despite being imperfect, these estimations can provide additional  
3 insights on whether a transcript discovery method has any bias towards highly- or lowly-expressed  
4 isoforms.

5  
6 To evaluate SIRV transcripts we used an incomplete SIRV annotation containing only 43 out of  
7 69 SIRV transcripts. The output annotations were again split into known and novel transcripts and  
8 compared against the respective reference set using gffcompare. The SIRV Set 4 annotations are  
9 available at <https://www.lexogen.com/sirvs/download/>.

10  
11 **Estimating consistency between annotations.** Consistency between transcripts generated on real  
12 data was estimated using gffcompare (without providing a reference annotation). Based on  
13 gffcompare output, for each tool we computed how many of its transcripts are supported by (i) all  
14 4 other tools, (ii) exactly 3 other tools, (iii) 1-2 other tools and (iv) no other tool (possible false  
15 predictions). We also counted the number of potentially missed transcripts that were reported by  
16 all methods except the one being evaluated (possible false negative). This approach is implemented  
17 in misc/denovo\_model\_stats.py.

18  
19 **Command line options.** For PacBio data minimap2 was launched with “splice:hq” preset; for  
20 ONT data we used k-mer size 14 with the usual “splice” preset. We also provided annotated splice  
21 junctions in BED format as an input. In each experiment, all tools were provided with the same  
22 BAM file and the same reference annotation. IsoQuant was launched with the default parameters  
23 setting the appropriate data type via “--data\_type” option. StringTie2 was launched with the “-L”  
24 option. All other tools were run with the default parameters in 20 threads. In contrast to all other  
25 tools, Bambu outputs all reference transcripts, including unexpressed ones. Thus, we filtered out  
26 all transcripts with read count values < 1 from the Bambu output. As recommended in the user  
27 manual, we also ran TALON using preliminary alignment correction with TranscriptClean<sup>30</sup>  
28 (<https://github.com/mortazavilab/TALON>). However, as the results with and without correction  
29 were almost identical, we decided to use the annotations obtained from raw data for a fair  
30 comparison. Complete information on all options and software versions are provided in the  
31 Supplementary Table 8.

32  
33 **IsoQuant algorithm.** To process long RNA reads, IsoQuant requires a reference genome and a  
34 corresponding gene annotation. If the reads are provided in the FASTQ format, IsoQuant maps  
35 them to the reference with minimap2 in splice mode<sup>2</sup>. Alternatively, a user may provide a sorted  
36 and indexed BAM file generated with a spliced aligner of their choice. The IsoQuant algorithm  
37 consists of 4 main steps: (i) assigning mapped reads to known isoforms, (ii) transcript  
38 quantification, (iii) alignment correction and (iv) transcript model construction. Below we describe  
39 the key aspects of all 4 procedures.

1 **IsoQuant algorithm: assigning long reads to known isoforms.** The algorithm for assigning long  
2 reads to annotated isoforms is based on intron chain matching and detecting exonic overlaps. To  
3 assign reads, IsoQuant processes each gene individually by extracting reads that map to the  
4 respective region from the sorted BAM file.

5  
6 IsoQuant first processes the annotation to construct intron and exon profiles of all known isoforms.  
7 A set of annotated introns in the gene is sorted according to their coordinates in the genome and  
8 enumerated from 1 to N. Thus, an annotated isoform can be represented as a vector of length N, in  
9 which the element at position  $i$  is set to 1 if this isoform includes the  $i$ -th intron and -1 otherwise  
10 (Supplementary Figure 3a). This vector is henceforth referred to as an isoform intron profile. The  
11 exon profile is constructed in a similar manner: all annotated exons are first split into a minimal  
12 set of M non-overlapping fragments, such that every exon can be represented as their combination,  
13 and these exonic fragments are sorted and enumerated. The exon profile for an annotated isoform  
14 is similarly denoted as a vector of length M, where the  $i$ -th element is set to 1 if this isoform  
15 contains the  $i$ -th exon fragment and -1 otherwise (Supplementary Figure 3b).

16  
17 To assign a read to an annotated isoform, each intron from the alignment is matched against  
18 annotated introns from the current gene and a read intron profile is constructed (also a vector of  
19 length N). In this vector the  $i$ -th element is set to 1 if the annotated intron with index  $i$  matches to  
20 an intron from the read, -1 if it is overlapped or spanned by the read, but no match is detected, and  
21 0 otherwise. A zero value indicates that the intron is located outside of the alignment region and  
22 therefore no information can be derived, e.g. due to read truncation. Similarly, the exon profile of  
23 the read is constructed based on M exonic fragments described above: 1 indicates that the  
24 respective exonic fragment is overlapped, -1 means it is spanned, and 0 is set for exonic fragments  
25 outside of the alignment region (Supplementary Figure 4).

26  
27 Due to sequencing errors, an aligner may detect splice site positions inaccurately<sup>19</sup>. To avoid  
28 considering them as alternative or novel, the algorithm allows a small difference  $\Delta$  between  
29 annotated and alignment splice site coordinates when matching introns. Formally speaking, an  
30 annotated intron  $(x_1, x_2)$  matches a read intron  $(y_1, y_2)$  if  $|x_1 - y_1| \leq \Delta$  and  $|x_2 - y_2| \leq \Delta$ . The default  $\Delta$   
31 value varies for different types of input data: 4 bp used for PacBio CCS reads and 6 bp for ONT  
32 reads (can be set manually). Although an aligned read can be assigned to an isoform by simply  
33 comparing its intron chain and exonic coordinates to the annotation, vectorizing the alignment as  
34 described above allows one to easily implement inexact splice site comparison with a delta, and  
35 quickly detect candidate isoforms for read assignment.

36  
37 Further, to assign a read to an isoform, its exon and intron profiles are matched against the  
38 respective profiles of the annotated isoforms. The distance between two profiles is computed  
39 simply as the number of distinct elements in which the read profile has non-zero values. A read is  
40 said to be consistent with an isoform if the distances between their exon and intron profiles are 0,

1 and the read has no unannotated introns/exons (Supplementary Figure 4). When a read is consistent  
2 with a single isoform, it is reported as a unique match. When a read is consistent with multiple  
3 isoforms simultaneously, it is classified as ambiguous, which may happen, for example, due to  
4 read truncation. If a read contains unannotated introns/exons, or its profiles are not consistent with  
5 any isoform, it is marked as inconsistent. For such alignments IsoQuant reports the most similar  
6 reference transcript and detected alternative splicing events.

7  
8 Some inconsistencies can be, however, caused by misalignments, rather than by real alternative  
9 splicing events<sup>17</sup>: (i) skipped short exons, (ii) intron shifts exceeding  $\Delta$  bp and (iii) short  
10 unannotated exons at transcript ends (Supplementary Figure 5). If an inconsistent alignment  
11 contains only these types of discrepancies, the read is reclassified as conditionally consistent.

12  
13 **IsoQuant algorithm: transcript quantification.** Once long reads are assigned to annotated  
14 isoforms, quantification becomes rather trivial. Uniquely assigned reads are counted as a single  
15 detected transcript, while ambiguous reads are treated as multi-mappers and contribute to multiple  
16 assigned isoforms with lower weight. A transcript is reported as expressed only if it has at least  
17 one uniquely assigned read. Inconsistent reads are considered as potential novel isoforms and  
18 ignored during the quantification step. Beside genes and transcripts, IsoQuant can also count  
19 inclusion and exclusion abundances for separate exons and introns, which can be useful for  
20 computing percent-spliced-in (PSI) values.

21  
22 IsoQuant implements additional functionality for barcoded long RNA reads, e.g. barcoded by  
23 single-cell or spatial location<sup>23,26</sup>. A user can provide information on how the reads are grouped,  
24 for example, as a TSV file that indicates a barcode or a cell type of origin for every read. Isoform  
25 and gene abundances are then calculated for every read group separately, which can facilitate an  
26 expression comparison between different groups or cell types.

27  
28 **IsoQuant algorithm: spliced alignment correction.** IsoQuant corrects each uniquely assigned  
29 read individually. If a read contains misalignments described above (Supplementary Figure 5) or  
30 its intron chain is not identical to the intron chain of the assigned isoform, the alignment is  
31 corrected as follows. Short skipped exons are restored according to the annotation and minor intron  
32 shifts are replaced with the respective introns from the assigned transcript. Unannotated terminal  
33 micro-exons are simply removed from the alignment. Finally, any unannotated splice site is  
34 substituted with the nearest site from the assigned transcript if (i) these splice sites are located  
35 within  $\Delta$  bp and (ii) read alignment contains sequencing errors near this splice site. Coordinates of  
36 corrected alignments are then saved in BED12 format.

37  
38 **IsoQuant algorithm: transcript model construction.** The transcript reconstruction procedure  
39 implemented in IsoQuant includes 4 steps: (i) intron graph construction from read alignments, (ii)  
40 intron graph simplification, (iii) attaching terminal vertices, (iv) construction of paths representing

1 full-length transcripts. Below we provide a detailed description of all algorithms and intuition  
2 behind them.

3

4 *Intron graph construction.* To construct transcript models IsoQuant implements a concept of intron  
5 graph, which was influenced by the previously designed splice graph approach<sup>27</sup>, used, for  
6 example, in StringTie<sup>5</sup>. For a given set of transcripts, an intron graph is constructed as follows.  
7 First, we define internal vertices as a set of all introns from all transcripts. Formally, each internal  
8 vertex is denoted by an ordered pair of coordinates in the genome. Two vertices are connected  
9 with a directed edge if the respective introns are consecutive in any transcript. Finally, for every  
10 first/last intron in a transcript, the corresponding vertex is connected with a terminating vertex that  
11 represents the transcript start/end position (formally, a single integer). Intron graph is a directed  
12 acyclic graph since every edge connects only consecutive elements. Each transcript can now be  
13 represented as a path in the graph that traverses from the initial to terminal vertex (Supplementary  
14 Figure 6a).

15

16 The described approach can be used to construct an intron graph from read alignments. Similarly  
17 to the read-to-isoform assignment procedure, the genes are processed by IsoQuant individually.  
18 First, the algorithm constructs a set of internal vertices corresponding to introns from the selected  
19 alignments. Two vertices are likewise connected when the respective introns are consecutive in  
20 any read alignment. Due to the presence of inexact detected splice sites, which may remain even  
21 after the alignment correction, such a graph may contain false vertices and connections. These  
22 false nodes typically form topological patterns, such as tips and bulges. A tip is defined as a dead  
23 end (dead start) edge that has a starting (ending) vertex with outdegree (indegree) at least 2. A  
24 bulge consists of two alternative paths having the same start and end vertices (Supplementary  
25 Figure 6b). Similar patterns are also typical for de Bruijn graphs, which are used for short read  
26 assembly, where bulges and tips are caused by sequencing errors. To remove tips and bulges  
27 assemblers exploit various techniques broadly called graph simplification<sup>28,29</sup>.

28

29 *Intron graph simplification.* Here we implement a graph simplification procedure based on the  
30 following observations: (i) a false intron is typically unannotated, (ii) splice site shifts that cause a  
31 false intron are short, and (iii) the number of reads supporting the correct intron often exceeds read  
32 support of a false intron. Formally, a bulge/tip is removed from the graph if it represents an  
33 unannotated intron that has at least twice lower read support compared to the alternative path and  
34 there is a known intron with splice sites within 20 bp (10 bp for PacBio). In other cases, when an  
35 unannotated intron has a high read support or no similar known intron exists, a bulge or a tip is  
36 likely to represent a part of a novel isoform and thus should be preserved (Supplementary Figure  
37 6b). Although intron graph simplification strongly resembles naive intron clustering, it has an  
38 important difference: an intron is removed not only based on its properties, such as splice site  
39 positions and read support, but based on the graph topology as well, thus taking into account  
40 adjacent introns. Such a method allows one to, for example, preserve similar introns from distinct

1 isoforms. It is worth noting, that the simplification procedure keeps track of all collapsed tips and  
2 bulges, thus preserving the possibility to later traverse alignment containing removed introns  
3 though the graph.

4  
5 *Collecting terminal positions.* After the graph is simplified, the algorithm proceeds to attach  
6 starting and terminal vertices. In contrast to annotated transcripts, read alignments do not provide  
7 the exact terminal positions, as their sequences can be truncated. Thus, to avoid having an extreme  
8 number of terminal vertices, terminal positions are detected using the heuristics presented below.  
9 Without loss of generality here we assume that the gene of interest is on the forward strand and  
10 polyA tails are on the right.

11  
12 For every intron  $V$  in the graph the algorithm selects only read alignments that contain  $V$  as a  
13 terminal intron and processes them as follows. First, the polyA sites are collected and clustered.  
14 Clustered polyA positions  $\{p_1, \dots, p_k\}$  are added to the graph as terminal vertices and connected  
15 to vertex  $V$  (Supplementary Figure 7a). Further, the algorithm adds the rightmost non-polyA  
16 terminal position  $P$  as a terminal vertex if one of the conditions is satisfied: (i)  $V$  has no outgoing  
17 edges, (ii)  $V$  has an outgoing edge to an intron  $(u_1, u_2)$  and  $P > u_1 + \Delta$ , or (iii)  $V$  has adjacent polyA  
18 vertices  $\{p_1, \dots, p_k\}$  and  $P > \max(p_1, \dots, p_k) + \Delta$  (where  $\Delta$  is the parameter defined above). Thus,  
19 a non-polyA terminal position can only be attached if it is located to the right of adjacent exons or  
20 polyA vertices. Starting positions are collected in a similar manner, but, indeed, without looking  
21 for polyA sites (Supplementary Figure 7b). The described approach, however, may lose  
22 information when several isoforms share the same starting intron but have distinct TSS/TES. Thus,  
23 we also apply an additional transcripts correction, which is described below.

24  
25 *Transcript discovery via path construction.* Once the intron graph is constructed and simplified,  
26 IsoQuant detects full-length paths that connect starting and terminal vertices. Paths entirely  
27 supported by at least a single read alignment (i.e., full splice match) are marked as transcript  
28 prediction candidates (Supplementary Figure 7c). To filter out unreliable novel transcripts  
29 IsoQuant applies read support cutoffs: at least 5 full-splice match reads (3 for PacBio) and at least  
30 2% from the maximum graph coverage. Since some isoforms may not have a full-splice matching  
31 alignment, IsoQuant also reports known transcripts that (i) have at least one uniquely assigned read  
32 and (ii) can be traversed through the intron graph. It also reports known mono-exonic transcripts  
33 that have (i) a uniquely assigned read and (ii) a confirmed polyA site.

34  
35 To correct terminal positions of a novel transcript, the algorithm selects all alignments consistent  
36 with this transcript and uses them to extract terminal positions using the approach described above  
37 (Supplementary Figure 7d). In contrast to detecting terminal vertices for the entire graph, where  
38 all alignments are used, the subset of consistent reads likely belongs specifically to this isoform  
39 and thus provides correct start and end positions. The resulting transcripts are saved in GTF format,  
40 providing additional information about transcript types and their reference genes.

## 1 **Data availability**

2 Nanopore sequencing data obtained from the human NA12878 cell line is available at  
3 <https://github.com/nanopore-wgs-consortium/NA12878/blob/master/RNA.md>. PacBio human  
4 GM12878 data is available at ENCODE (<https://www.encodeproject.org/search>) under the  
5 accession numbers ENCFF450VAU and ENCFF694DIE. Sequencing data obtained from mouse  
6 brain samples is available at NCBI Gene Expression Omnibus (GEO;  
7 <https://www.ncbi.nlm.nih.gov/geo/>) under accession numbers GSE158450 and GSE178175. ONT  
8 SIRV data, simulated data and reduced gene annotations are available at  
9 <https://data.cab.spbu.ru/index.php/s/dgcaSaGME2xF7ed?path=%2FIsoQuant>. Complete  
10 information on all datasets used in this work is listed in the Supplementary Table 7.

## 11 **Code availability**

12 IsoQuant and the supplementary scripts used for the evaluation are available at  
13 <https://github.com/ablab/IsoQuant>.

## 14 **Acknowledgements**

15 We thank Nanopore WGS consortium and Ali Mortazavi's lab at UCI for making the ONT and  
16 PacBio data publicly available. This work was supported by St. Petersburg State University, Russia  
17 (grant ID PURE 93023437 to A.M., A.S., A.L.L. and A.D.P.). Scientific research was performed  
18 at the Research park of St.Petersburg State University «Computing Center».

## 19 **Author contributions**

20 A.D.P., A.M. and A.S. designed and implemented the software. A.D.P., A.M. and A.J. performed  
21 the benchmarks. J.J. performed the sequencing experiments. A.D.P., A.L.L., and H.U.T. acquired  
22 funding and supervised the project. A.D.P., A.M., A.J. and H.U.T. wrote the manuscript.

## 23 **Competing interests**

24 The authors declare no competing interests.

## 25 **References**

- 26 24. Ono, Y. et al. *Bioinformatics*. 29, 119-221 (2013).  
27 25. Frankish, A. et al. *Nucleic acids research* 49, D916-D923 (2021).  
28 26. Gupta, I. et al., *Nature biotechnology* 36, 1197-1202 (2018).  
29 27. Heber, S. et al. *Bioinformatics* 18, S181-S188 (2002).  
30 28. Zerbino, D.R. & Birney, E. *Genome research* 18, 821-829 (2008).  
31 29. Bankevich, A. et al. *Journal of computational biology* 19, 455-477 (2012).  
32 30. Wyman, D. & Mortazavi, A. *Bioinformatics* 35, 340 (2019).

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [IsoQuantsupplementarymaterial.pdf](#)