



Quality assessment of long read data in multisample lrrNA-seq experiments with SQANTI-reads

Netanya Keil, Carolina Monzó, Lauren McIntyre, et al.

Genome Res. published online March 3, 2025

Access the most recent version at doi:[10.1101/gr.280021.124](https://doi.org/10.1101/gr.280021.124)

P<P	Published online March 3, 2025 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Quality assessment of long read data in multisample lrrRNA-seq experiments using SQANTI-reads

Netanya Keil^{1,2}, Carolina Monzó⁴, Lauren McIntyre^{1,2,3,*}, Ana Conesa^{4,*}

¹Department of Molecular Genetics and Microbiology, University of Florida, Gainesville, FL, USA, 32610

²University of Florida Genetics Institute, University of Florida, Gainesville, FL, USA, 32610

³UF Health Cancer Center, University of Florida, Gainesville, FL, USA, 32610

⁴Institute for Integrative Systems Biology (I2SysBio), Spanish National Research Council (CSIC), Paterna, 46980, Spain

*corresponding authors: ana.conesa@csic.es; mcintyre@ufl.edu

ABSTRACT

SQANTI-reads leverages SQANTI3, a tool for the analysis of the quality of transcript models, to develop a read-level quality control framework for replicated long-read RNA-seq experiments. The number and distribution of reads, as well as the number and distribution of unique junction chains (transcript splicing patterns), in SQANTI3 structural categories are informative of raw data quality. Multisample visualizations of QC metrics are presented by experimental design factors to identify outliers. We introduce new metrics for 1) the identification of potentially under-annotated genes and putative novel transcripts and for 2) quantifying variation in junction donors and acceptors. We applied SQANTI-reads to two different datasets, a *Drosophila* developmental experiment and a multiplatform dataset from the LRGASP project and demonstrate that the tool effectively reveals the impact of read coverage on data quality, and readily identifies strong and weak splicing sites. SQANTI-reads is open source and is available in versions $\geq 5.3.0$ in the SQANTI3 GitHub repository.

30 INTRODUCTION

31 Short-read RNA sequencing (srRNA-seq) is the most common and cost-effective approach for
32 studying the transcriptome. In srRNA-seq, transcripts must be inferred computationally, which
33 can lead to inaccuracies in transcript identification (Liu et al. 2016; Newman et al. 2018). Recent
34 advances in single-molecule long-read sequencing technologies have opened new avenues for
35 transcriptome analysis (reviewed in (Marx 2023; van Dijk et al. 2023)). In long-read RNA
36 sequencing (lrRNA-seq), full-length transcripts can be observed as single sequencing reads,
37 allowing for direct transcript detection without the need for an assembly step. As with any
38 technology, lrRNA-seq is not without errors, and factors such as mRNA degradation, library
39 preparation failures, and sequencing inaccuracies can result in bias in the data.

40

41 A database tracking bioinformatic tools for long-read sequencing (Amarasinghe et al. 2021)
42 identifies numerous tools for the initial processing of lrRNA-seq data primarily assessing the
43 accuracy of basecalling and the length of the reads. These include pycoQC (Leger and Leonardi
44 2019), longQC (Fukasawa et al. 2020) and nanoQC (De Coster et al. 2018), which offer a first-
45 pass analysis of lrRNA-seq data. Other tools, such as SQANTI3 (Pardo-Palacios et al. 2024a),
46 TALON (Wyman et al. 2020), FLAMES (Holmqvist et al. 2021), Iso-Seq (<https://isoseq.how/>), and
47 IsoTools (Lienhard et al. 2023), focus on evaluating transcript models inferred from the data.
48 However, most current tools for lrRNA-seq read quality control were developed during the
49 early stages of these technologies and are generally limited in the number of evaluated features
50 and/or samples. As long-read sequencing technologies rapidly improve in quality and decrease
51 in cost, the experimental scope possible with these technologies has expanded. The need for a

52 comprehensive and comparative read quality assessment tool capable of analyzing millions of
53 reads and dozens (or more) samples is critical.

54

55 The rapid decline in costs implies that the use of lrRNA-seq will continue to expand, with
56 experimental designs involving multiple samples becoming more common (Glinos et al. 2022;
57 Joglekar et al. 2024; Mahmoud et al. 2024; Patowary et al. 2024). From a quality control
58 perspective, this necessitates that datasets are homogeneous, without systematic bias
59 associated with experimental groups, and free of outliers. Moreover, the generated data must
60 be sufficient to address the research questions that motivated the experiment. The increase in
61 throughput now makes it possible to design experiments that include barcoding and
62 multiplexing to balance library preparation and sequencing across experimental groups (Auer
63 and Doerge 2010). This approach helps avoid confounding technical variation with the
64 treatments of interest and facilitates discriminating between failed technical replicates and
65 failed samples. Finally, technological advancements such as more accurate basecallers
66 (<https://github.com/nanoporetech/dorado>) and the availability of novel library preparation
67 methods such as MAS-Iso-Seq (Al'Khafaji et al. 2024), CapTrap (Carbonell-Sala et al. 2024), R2C2
68 (Volden et al. 2018), Nano3P-seq (Begik et al. 2023), or FLAM-seq (Legnini et al. 2019)
69 motivates the need for tools that can easily evaluate how these improvements impact data
70 quality.

71

72 In this context, we present SQANTI-reads, an extension of SQANTI3 (Pardo-Palacios et al.
73 2024a), a tool originally designed for transcript model quality control, to jointly provide quality

74 control metrics for long-read data and to analyze multiple samples for consistency and bias. We
75 demonstrate that SQANTI3's structural categories and other quality control metrics, repurposed
76 in SQANTI-reads, are highly effective for assessing the homogeneity in a lrrNA-seq multisample
77 experiment, identifying read quality control failures, and detecting outliers. Additionally, we
78 have added new metrics that provide insights into the potential utility and discovery power of
79 the data, including variation at donor/acceptor sites and identification of potentially under-
80 annotated genes and mis-annotated transcripts.

81

82 **RESULTS**

83 **SQANTI-reads can be used to evaluate long-read technology improvements**

84

85 Long-read methods are rapidly improving, and both Nanopore and PacBio are updating their
86 instruments, molecular protocols, and algorithms. Tools that can readily evaluate the broad
87 impact on data quality of these technical improvements are highly needed to make informed
88 decisions for downstream analyses. Our *Drosophila* dataset contained the same set of samples
89 processed on the MinION and PromethION instruments through several runs. In addition, the
90 development of basecalling algorithms for ONT sequencing is a highly active and rapidly
91 evolving field (Pagès-Gallego and de Ridder 2023; Diensthuber et al. 2024). We expect that as
92 new basecallers emerge and algorithms are improved upon, users will want to evaluate the
93 impact of different basecallers and algorithms on the quality of their experiments. This
94 comparison can be done using SQANTI-reads. While Dorado is now the default basecaller for
95 ONT, we demonstrate the utility of SQANTI-reads by comparing real-time Guppy basecalled

96 reads and Dorado basecalled reads. We first evaluated whether Dorado effectively improved
97 data quality without introducing biases and if data from several sequencing experiments could
98 be merged. We compared the Guppy and Dorado basecallers using SQANTI-reads metrics. As
99 anticipated, Dorado resulted in more reads with assignable barcodes, a higher number of
100 mapped reads, more reads aligning to annotated genes, more reads aligning to annotated
101 transcripts, and longer reads, without an increase in the proportion of reads with technical
102 artifacts (Supplementary Figure 1, Supplementary File 1). This confirms that Dorado improves
103 basecalling accuracy without introducing unwanted biases and motivated the selection of
104 Dorado basecalled reads for subsequent analyses.

105
106 In the *Drosophila* experiment, libraries were barcoded, pooled, and multiplexed across different
107 MinION and PromethION runs, with a re-pooling step between the two machines. We used
108 SQANTI-reads to compare the quality of the MinION and the PromethION runs, and to evaluate
109 the consistency of the MinION and PromethION technology across technical replicates. The first
110 MinION run (TR1) had higher percentages of reads classified as novel in catalog (NIC) and novel
111 not in catalog (NNC) and with non-canonical junctions compared to the second and third
112 technical replicates MinION runs (TR2, TR3) and compared to the PromethION runs for the
113 same libraries (Supplementary Figure 2). NIC refers to reads with novel combinations of
114 annotated junctions and NNC refers to reads with at least one unannotated junction. MinION
115 runs TR2 and TR3 were similar in their quality metrics (described below) to the PromethION run
116 of the same samples, and the technical replicates on the PromethION were similar. These
117 results indicate that the technology performs consistently across instruments and runs. Based

118 on these SQANTI-reads QC results, we aggregated data across technical replicates to further
119 evaluate the quality of the lrRNA-seq experiment. Although TR1 had lower quality than the
120 other technical replicates, the overall read numbers were low, and we decided to keep this
121 technical replicate in our evaluation of the samples.

122

123 **SQANTI-reads metrics can be used to evaluate the global quality of the lrRNA-seq experiment**

124 In a multisample lrRNA-seq experiment, all samples should be of similar quality. SQANTI3 uses
125 the full-splice match (FSM) structural category to identify long-read sequences whose junctions
126 are consistent with an annotated transcript model. However, for an lrRNA-seq experiment to
127 accurately reflect the analyzed transcriptome, the reads should also capture the distribution of
128 transcript lengths of the expressed transcriptome. The distribution of transcript lengths
129 depends on many factors, including the species and tissue used as input. Consistent with the
130 species transcript model annotations, we observe that *Drosophila* has less complex and shorter
131 transcripts than humans (Supplementary Figure 3). A dataset with reads substantially shorter
132 than the annotated, transcriptome but a high proportion of FSM indicates capture of short
133 transcripts, while combining a high proportion of incomplete-splice match (ISM) may indicate
134 RNA degradation. ISM refers to a read whose junction matches an annotated transcript but is
135 missing junctions on the 5' end, 3' end or both ends. We looked at these values for an initial
136 assessment of the quality of the *Drosophila* experiment.

137

138 First, we compared the number of reads and distributions of read lengths for all samples. The
139 difference in sequencing depth (number of reads) between the two developmental stages was

140 evident, as expected because of the additional PromethION run on the 3-8 day samples (Figure
141 1A). For all samples, most reads were shorter than 1 kb, with less than 20% of them above the 1
142 kb threshold (Figure 1B, Supplementary Figure 4A, Supplementary Figure 5). When we
143 evaluated all reads, 53% to 67% of the reads across samples were classified as FSM, 20% to 38%
144 were labelled as ISM. For all samples the proportion of NIC/NNC was less than 10% of the reads
145 (Figure 1C). For the reads that were above the 1kb threshold, between 73% and 82% of the
146 reads were FSM, with only 10% to 18% ISM (Figure 1D). The decrease in proportion of ISM and
147 increase in proportion of FSM when evaluating only the reads >1kb suggests that the shorter
148 reads represent incomplete transcript sequences, potentially due to mRNA degradation.

149

150 Reads that share all internal junctions are grouped together and annotated using a string made
151 up of the junction locations, the unique junction chain (UJC) (Nanni et al. 2024). We examined
152 both gene-level and UJC-level metrics. We found that, despite the large sequencing depth
153 differences between developmental stages, the number of detected genes was only slightly
154 lower in the 0-1 h samples (Figure 1E). However, these genes were quantified with fewer reads
155 (80% genes with < 50 reads) than the 3-8 d samples, which had between 30% and 50% of genes
156 with more than 100 reads (Figure 1E). When evaluating UJC, we found that, while the number
157 of UJC mirrored the sequencing depth pattern (Figure 1F), with 3-8 d samples showing five
158 times more UJC than 0-1 h samples, and a larger number of FSM and ISM UJC, there were many
159 additional UJC detected by fewer than 10 reads, and usually by a single read (Figure 1F & 1G).
160 These UJC were most frequently NIC/NNC (Supplementary Figure 4B & 4C). Downstream
161 analyses would therefore need to address whether this represents novel low-expressed

162 transcripts or technology errors. In contrast, the percentage of FSM reads between the two
163 time points differed by less than 1x in all replicates (Figure 1H). These results indicate that the
164 higher sequencing depth of the 3-8 d samples does not change the number of detected genes
165 or annotated transcripts (FSM). The higher read depth per gene/UJC suggests that more genes
166 and transcripts will be able to be quantitatively evaluated in the 3-8 day samples compared to
167 the 0-1 hour samples.

168

169 In the *Drosophila* data, we noticed two samples (Sample 4 Rep 1 (RIL 12279,0-1 hour) - red
170 arrow); Sample 6 Rep 3 (RIL 11255,3-8 day) - teal arrow) that had the lowest percentage of FSM
171 and highest percentage of ISM in the 0-1 hour and 3-8 day groups respectively (Figure 1C). To
172 determine whether these two samples were of overall lower quality than the rest, we examined
173 their SQANTI-reads metrics. We found that Sample 4 Rep 1 had a lower proportion of FSM
174 across all genes (Figure 1H) and a higher proportion of genes quantified with only one read
175 (Figure 1F), while Sample 6 Rep 3 had a similar gene (Figure 1H), UJC (Figure 1G), and % FSM in
176 genes (Figure 1H) than other 3-8 day samples. We concluded that RIL 12279 rep 1 0-1 hour is a
177 lower-quality sample.

178

179 Altogether, this example shows that SQANTI-reads metrics can be used to compare samples
180 and experimental conditions in a multisample experiment, detect outliers, and suggest points
181 of attention for downstream data processing.

182

183

184 **SQANTI-reads metrics can be used to identify systematic differences among samples**

185 The previous example demonstrated that SQANTI-reads metrics are effective in assessing
186 dataset consistency. However, SQANTI-reads evaluates over 35 quality metrics, making it
187 challenging to determine which features contribute to potential differences among samples.
188 We include Principal Component Analysis (PCA) analysis to identify which metrics are the most
189 relevant for quality variability when there are differences among samples or between groups.
190 The percentage of reads and UJs in each structural category, percentage of artifact reads (RT-
191 switching, non-canonical junctions, and intrapriming), percentage of junctions in each category,
192 as well as length metrics, are included in the PCA.

193

194 We applied SQANTI-reads PCA analysis of quality features to investigate differences in read
195 quality among various long-read sequencing methods used in the Long-read RNA-seq Genome
196 Annotation Assessment Project (LRGASP) challenge (Pardo-Palacios et al. 2024b), focusing on
197 the WTC11 dataset. We evaluated triplicate transcriptome measurements of the WTC11 human
198 cell line, analyzed using three technologies: cDNA PacBio Sequel II (cDNA PacBio), cDNA Oxford
199 Nanopore MinION (cDNA ONT), and direct RNA Oxford Nanopore MinION (dRNA ONT). The
200 analysis revealed that WTC11 samples clustered based on the long-read technology applied
201 (Figure 2A). Specifically, PC1, which explains 56% of the variance, distinguished cDNA ONT
202 samples from those generated by the other two technologies, while PC2, accounting for 35% of
203 the variance, highlighting differences between dRNA ONT and cDNA PacBio. To further explore
204 these differences, we examined the loadings for each principal component. Quality features
205 with the highest positive loadings in PC1 included the number of reads, the percentage of

206 reads, and the proportion of UJCs in the NNC category, while features with high negative
207 loadings included Intergenic and Genic Genomic reads. Several junction-related variables also
208 exhibited high absolute loadings on PC1 (Figure 2B). SQANTI-reads plots confirmed these
209 structural category differences between cDNA ONT samples and other library preparations.
210 cDNA ONT had both the highest proportion of NNC reads and UJCs (Figure 2C and 2D) and also
211 had the lowest proportion of intergenic reads (Figure 2C). Other differences in sequencing
212 throughput and junction characteristics were also confirmed (Supplementary Figure 6).

213
214 Upon examining the feature loadings for PC2, we found that variables with high contributions
215 included several metrics related to read length (Figure 2B). Consequently, we evaluated the
216 SQANTI-reads “Lengths of All Mapped Reads” plot for this experiment. Indeed, we observed
217 that the cDNA PacBio method produced a significantly higher proportion of reads between 1-2
218 kb, 2-3 kb, and greater than 3 kb, as suggested by their negative loadings, compared to the
219 dRNA ONT method, which predominantly generated reads shorter than 1 kb (Figure 2E).
220 Similarly, the percentage of reads assigned as ISM with high positive values was higher in dRNA
221 ONT samples (Figure 2C).

222
223 In conclusion, we showed that the SQANTI-reads PCA analysis is an effective tool for uncovering
224 significant read quality differences between long-read sequencing methods. The technological
225 differences explored here are based on a benchmarking experiment that has been superseded
226 by new technology and the results here should not be interpreted as an evaluation of current
227 technological capabilities. Nonetheless, we have shown that SQANTI-reads can be used to

228 evaluate long-read sequencing methods and for revealing technological biases.

229

230 **SQANTI-reads identifies potentially under-annotated genes**

231 Long-read data often contain a large number of sequences that cannot be exactly matched to
232 existing annotations. In many cases, these UJCs belong to annotated genes and are identified by
233 only a few reads, as illustrated by the SQANTI-reads analysis in Figures 1 and 2. This suggests
234 they could be either low-expressed transcripts or technological artifacts. However, in some
235 instances, a high proportion of reads in a gene may correspond to the same novel UJC,
236 indicating the possibility of a previously unannotated transcript that warrants closer
237 examination. SQANTI-reads includes a customizable decision tree to identify such cases (see
238 Methods, Figure 3A). Basically, the tool identifies genes with a high ($R > \text{threshold}$) number of
239 reads, with novel UJC containing a large fraction ($Q > \text{threshold}$) of the gene's splice sites and
240 capturing a high proportion ($P > \text{threshold}$) of the reads. We applied this approach to the
241 WTC11 PacBio data using default parameters ($R=100$, $P=20$, and $Q=80$). In addition, we classify
242 expressed genes as well annotated or underannotated using the parameters P and R . The logic
243 for defining expressed genes as well annotated or underannotated is described in Figure 3A.
244 The annotation category for all expressed genes (default: number of reads in gene $[R] > 100$) is
245 provided in the `gene_classification.csv` file.

246

247 From the set of expressed genes, we identified 8,556 well-annotated genes, 88% of which have
248 a well-covered annotated transcript ($>20\%$ of total gene coverage) (Figure 3B, Supplementary
249 Figure 7). We also identified 101 genes for which there are no reads with an FSM match to an

250 annotated transcript. Of these, 54% have a well-covered UJC. For all expressed genes with a
251 well-covered unannotated UJC, we identified 424 that contained most of the observed
252 junctions in that gene (>80%), and we label these putative novel transcripts (Figure 3C,
253 Supplementary Figure 7). Of these, 316 were NIC and 108 were NNC (Supplementary Figure 7).
254 The SQANTI-reads output for putative novel transcripts is included in the
255 putative_underannotation.csv file (Table 1).

256

257 For genes with at least one putative novel transcript ($R > 100$, $P > 20$, $Q > 80$) and an annotated
258 transcript (FSM), we selected the FSM with the highest proportion of reads. We then compared
259 the structure of the putative novel transcripts to the most expressed annotated transcript using
260 TranD (Nanni et al. 2024). For genes with an annotated transcript that is relatively highly
261 expressed (>20% of the reads for that gene), 103 putative candidate transcripts differed from
262 the annotated transcript by donor/acceptor variation, suggesting a possible alternative splice
263 site. In addition, 10 putative candidate transcripts had an extra exon, 15 a skipped exon, and 9
264 with both missing and skipped exons relative to the most expressed annotated UJC (Figure 3D).
265 For the genes where the annotated transcript represented less than 20% of the reads in that
266 gene, the putative novel transcript differed from the annotated transcript by an alternative
267 exon in 147 cases (33 extra exons, 86 skipped exons, and 43 with both an extra and skipped
268 exon) (Figure 3E). Details for this analysis are provided in the Supplementary Methods.

269

270 This analysis shows that SQANTI-reads can readily identify under-annotated genes and flag
271 putative novel transcript models that contain interpretable alternative exonic patterns that

272 deserve further attention.

273

274 **SQANTI-reads metrics for donors/acceptors identify noisy splicing and potentially novel splice**
275 **sites**

276 SQANTI-reads calculates the mean, standard deviation, and coefficient of variation (CV) for all
277 expressed annotated donors/acceptors (Figure 4A). A $CV > 0$ indicates variability in the
278 donor/acceptor, with higher CV values indicating more variability. Variability around a splice
279 junction may be due to weak splicing (Wang and Marín 2006) or to technology errors and
280 mapping accuracy, for example, due to junction ambiguity (Li 2018). We evaluated these
281 metrics on the WTC11 dataset for reference junctions with at least 10 reads. We found similar
282 patterns in the variability ($CV > 0$) in donors and acceptors (Figure 4B). All three technologies
283 identify donors and acceptors with variability around the splice site ($CV > 0$). Donors/acceptors
284 with $CV > 0$ consistently across the three technologies are highly suspicious of ‘noisy’ splicing or
285 a weak splice site and may warrant follow-up (Supplementary Figure 8A and 8B). The SQANTI-
286 reads output file `cv.csv` identifies the donors/acceptors with $CV > 0$, making it straightforward
287 to follow up on particular locations with tools such as Integrative Genomics Viewer (IGV)
288 (Robinson et al. 2011).

289

290 A reference match junction indicates that the splice signal is strong (Wang and Marín 2006;
291 Dent et al. 2021). Results for cDNA PacBio and dRNA-seq were similar, with both showing a
292 higher number and proportion of reference match donors/acceptors compared to cDNA ONT
293 (Figure 4B and C), despite these technologies detecting similar numbers of FSM UJs (Figure

294 2D). We compared the FSM UJCs identified by the three methods (Supplementary Figure 9A).
295 Most of the FSM UJCs were detected by all three technologies ($n = 19,690$), with a similar
296 number detected by cDNA PacBio only ($n = 8,110$) and cDNA ONT ($n = 9,360$) only. We
297 hypothesized that the difference in the number of reference match donors/acceptors was
298 potentially due to longer transcripts with more junctions being detected in cDNA PacBio
299 compared to shorter transcripts with fewer junctions in dRNA ONT and cDNA ONT. For the
300 FSMs detected only in one technology, we plotted the distribution of the number of junctions
301 and confirmed that the cDNA PacBio FSM transcripts had a larger number of junctions
302 compared to dRNA ONT and cDNA ONT (Supplementary Figure 10). This agrees with cDNA
303 PacBio showing longer reads than both ONT technologies in the WTC11 dataset (Figure 2E).
304
305 Splice junctions may differ from annotated sites due to the presence of novel donors/acceptors.
306 The category $CV = 0$ identifies donors/acceptors with no variability but differing from the
307 annotated donor/acceptor, representing strong candidates for bona fide alternative splice sites.
308 We evaluated the donors and acceptors with $CV = 0$ in all the technologies (Supplementary
309 Figures 8E and 8F). We identified 51 donors and 34 acceptors with $CV = 0$ in all three
310 technologies. Of these, 76% of donors and 65% of acceptors are within 12 nt of an annotated
311 donor/acceptor, indicating a potential misannotation of the splice site (Supplementary Figure
312 11). These donors and acceptors with $CV = 0$ across all technologies were detected by a
313 minimum of 10 reads and, in some cases, could be detected with >100 reads (Supplementary
314 Figure 12). Supplementary Figure 13 shows an example of one of the reference donors with CV
315 $= 0$ detected across all three technologies. This donor is in the H2AZ1 gene at position

316 99,948,814 on Chromosome 4. All the reads associated with this donor map to position
317 99,948,811, which is 3 nts away from the annotated donor position (Supplementary Figure 13).
318 Donors and acceptors with CV = 0 detected consistently across samples and/or technologies
319 indicate potential robust detection of alternative splice sites.

320

321 **Memory usage and runtime of SQANTI-reads**

322 Computational efficiency is a critical factor for users when choosing a quality control tool. To
323 evaluate memory usage and running time of SQANTI-reads, datasets containing 5 different
324 number of reads were generated, as well as 5 replicates for each dataset. After independently
325 processing these datasets through SQANTI3, they were tested 5 times in 15 combinations to
326 analyze memory usage and runtime for: a single sample of each size, 3 samples of each size,
327 and 5 samples of each size.

328 The performance evaluation of SQANTI-reads revealed that memory consumption remained
329 consistent at 2 GB of RAM across all tested conditions, regardless of the number of reads or the
330 number of samples analyzed. Runtime, however, increased with both the number of reads per
331 sample and the total number of samples processed simultaneously (Supplementary Figure 14)
332 When processing a single sample, runtimes ranged from 100 seconds for datasets with 10,000
333 reads to approximately 3 minutes for datasets with 1,000,000 reads. For multisample runs,
334 processing 3 samples of 1,000,000 reads each required approximately 4 minutes, while
335 processing 5 samples of the same size took about 5 minutes. These findings highlight the
336 scalability of SQANTI-reads, with a predictable increase in runtime proportional to the dataset
337 size and sample number. Overall, SQANTI-reads demonstrated computational efficiency with

338 minimal memory requirements, making it suitable for large-scale transcriptomic analyses.

339 Table 1. SQANTI-reads specific output files.

Output file	Description	Default output	Output type
gene_counts.csv	Number of reads in each structural category, per gene and per sample	Yes	Multiple samples file
ujc_counts.csv	List of junction hashes in each sample and the number of reads in each sample associated with each junction string. Flags the most expressed UJC per gene	Yes	Multiple samples file
length_summary.csv	Number and percentage of reads in length categories per sample	Yes	Multiple samples file
cv.csv	Metrics on the coefficient of variance of reference junctions for each sample	Yes	Multiple samples file
jxn_counts.csv	Number of known canonical, novel canonical, known non-canonical, and novel non-canonical junctions in reads of each sample	--all-tables	Multiple samples file
cv_acc_counts.csv cv_don_counts.csv	Number of detected annotated donors and acceptors in each junction variation category	--all-tables	Multiple samples file
FSM_counts.csv ISM_counts.csv NIC_NNC_counts.csv	Number of reads in each subcategory for FSMs, ISMs, NICs and NNCs	--all-tables	Multiple samples file
err_counts.csv	Number and percentage of reads with evidence of intrapriming, RT-switching, and non-canonical junctions per sample	--all-tables	Multiple samples file
pca_loadings.csv	PC1 and PC2 loadings from PCA	--pca-tables	Summary file
pca_variance.csv	Variance explained by each PC	--pca-tables	Summary file
sample_quality_flags.csv	Binary quality indicators to flag potential sample issues	Yes	Summary File
gene_classification.csv	For genes with coverage over a user-defined threshold, gives the annotation category of each gene	Yes	Summary file
putative_underannotation.c	Metrics on NIC and NNC UJCs	Yes	Summary file

sv	and flags putative novel transcripts		
----	--------------------------------------	--	--

340

341

342 Table 2. Comparison between SQANTI3 and SQANTI-reads.

Feature	SQANTI3	SQANTI-reads
Sequences analyzed	transcript models	reads
Annotation with SQANTI3 categories	yes, for transcript models	yes, for reads and UJCs
Computation of SQANTI3 quality metrics	yes	yes
Samples processed	One	Multiple
Visualizations across samples	no	yes
PCA analysis between samples	no	yes
Summary of read counts (per gene, per UJC)	no	yes
Donor/acceptor variation metrics	no	yes
Identification of putative under annotated genes	no	yes
Identification of putative novel transcripts	no	yes
Machine learning validation of transcript models	yes	no
IsoAnnot annotation	yes	no

343

344

345

346 **Figure 1: SQANTI-reads analysis of *Drosophila* samples.** A) Number of mapped reads by
347 experimental group labeled with read length. B) Percentage of mapped reads >1kb vs
348 percentage of reads that are FSM. Dots represent early stage (0-1 hours after eclosion) and
349 crosses indicate adult stage (3 to 8 days old). The four genotypes are indicated with four
350 different colors. C) Percentage of reads mapping to genes in each SQANTI3 structural category
351 D) Percentage of reads mapping to genes in each SQANTI3-QC structural category for reads
352 >1kb E) Number of genes detected with breakdown by the number of reads mapped to each
353 gene. F) Number of UJCs detected with breakdown by the number of reads associated with
354 each UJC. G) Proportion of UJCs detected with breakdown by the number of reads associated
355 with each UJC. H) Distribution of the percentage of FSM reads by gene across samples.

356

357 **Figure 2: SQANTI-reads PCA analysis of LRGASP WTC11 samples.** A) PCA using SQANTI-reads
358 quality features. The percentage of variance explained by each principal component (PC) is
359 labelled on each axis in brackets. B) Top 10 Loadings for PC1 and PC2. C) Distribution of reads in
360 SQANTI3 structural categories. D) Distribution of UJCs in structural categories E) Distribution of
361 read lengths for all mapped reads.

362
363 **Figure 3: Evaluation of WTC11 PacBio samples for under-annotated genes.** A) Decision tree for
364 classifying genes as well-annotated or under-annotated and classifying transcripts as putative
365 novel transcripts. B) Number of genes by their annotation status according to SQANTI-reads
366 parameters. C) The coverage (percent of total reads) vs length (percentage of maximum
367 junctions) for all UJCs in under-annotated genes with well covered candidate transcripts. UJCs
368 that meet the thresholds for putative novel transcripts are colored in orange. D&E) Upset plot
369 of putative novel transcripts with the most expressed annotated transcript in that gene for
370 well-annotated genes with well-covered FSMs (D) and with an FSM detected but without < 20%
371 of the total reads in the gene(E). Upper and middle panels show the distribution of the number
372 and proportion of nucleotides different between the putative novel transcript and most
373 expressed FSM.

374
375 **Figure 4. Variation in donors/acceptors.** Metrics are only calculated for annotated
376 donor/acceptors with a minimum threshold of reads (10 by default). A) Metrics for the
377 classification of donor/acceptor variation. When all reads align to the annotated
378 donor/acceptor this is classified as a reference match (Ref Match). When all reads align to the
379 same donor/acceptor location, but this is not the annotated position this is classified as $CV = 0$.
380 When reads align in multiple positions in proximity to an annotated donor/acceptor this is
381 classified as $CV > 0$. B) Classification of the number (left) and proportion (right) of detected
382 acceptors faceted by technology. C) Classification of the number (left) and proportion (right) of
383 detected donors faceted by technology.

384

385

386 **DISCUSSION**

387 The increase in throughput and generalization of lrRNA-seq has led to the growing popularity of
388 experiments containing many samples. Multiple lrRNA-seq studies have shown that, despite
389 technological improvements, biases associated with read length and accuracy are still present
390 (Amarasinghe et al. 2020; Delahaye and Nicolas 2021). These biases vary between the major
391 long-read sequencing (LRS) platforms and with the introduction of new instruments and
392 chemistries. Other types of systematic biases, such as those introduced by batches utilized in

393 large experiments, may also occur. Benchmarking studies have also revealed that transcript
394 reconstruction methods can yield highly different results due to the varying strategies used to
395 resolve data inaccuracies (Pardo-Palacios et al. 2024b). This highlights the need for tools that
396 can comprehensively evaluate the characteristics of raw long-read data before making
397 decisions on downstream analyses. Direct examination of read quality enables the researcher
398 to evaluate the experiment for consistency and identify any outlier samples and any systematic
399 differences in read quality between sample groups.

400

401 We have developed SQANTI-reads as a tool that enables a comprehensive assessment of reads
402 obtained from an LRS experiment. We leverage the widely adopted SQANTI3 framework and
403 add several metrics for the assessment of reads. SQANTI-reads is a flexible tool that can be
404 used to evaluate the quality of lrrNA-seq multisample experiments. We present metrics that
405 evaluate reads in terms of whether they correspond to annotated transcript models using
406 SQANTI3 categories and subcategories. If metadata are included in the design file, alternate
407 sample groupings can be used without needing to rebuild the classification and junction files.

408

409 For example, in Oxford Nanopore Technology (ONT), the existence of multiple platforms at
410 different price points for different numbers of pores (Flongle, MinION, GridION, PromethION)
411 but with the same library protocols means that, in a large experiment, samples can be initially
412 evaluated at low cost on one of the lower-throughput platforms (Flongle MinION, GridION). If
413 samples are of sufficient quality, they can then be run on higher-throughput platforms
414 (PromethION). Sample multiplexing and running on multiple "lanes" is good experimental

415 design practice (Auer and Doerge 2010). Our SQANTI-reads analysis of the *Drosophila* dataset
416 illustrates such a scenario. In this case, the pool for the 0-1 hour samples was initially evaluated
417 on the MinION. The resulting data from TR1 were unbalanced and had relatively short reads
418 with a high proportion of ISM. These observations enabled adjustment of the library
419 concentrations and run parameters, and a second MinION run resulted in fewer ISM and more
420 balance in read numbers across libraries. A PromethION run based on the rebalanced libraries
421 efficiently used the sequencing resources across samples. The same procedure was deployed
422 with the 3-8 day samples, and rebalancing helped ensure efficiency of the subsequent two
423 PromethION runs. This example illustrates both good experimental design practices for large
424 experiments and the utility of SQANTI-reads in assessing data quality at early stages of data
425 acquisition, enabling corrections that lead to a successful sequencing experiment.

426

427 Another important aspect of lrrna-seq multisample experiments is the ability to quickly assess
428 whether the data can address the biological questions that motivated the study. This includes,
429 among other things, whether genes and transcripts are sufficiently quantified and if potential
430 novel transcripts are adequately supported. While these questions may ultimately be answered
431 after full data processing with transcript reconstruction algorithms, users may find it helpful to
432 evaluate support directly from the raw data.

433

434 SQANTI-reads provides information on the distribution of reads across genes and UJC (a raw-
435 data proxy for transcripts) and introduces new metrics for identifying variations in
436 donors/acceptors, under-annotated genes, and putative novel transcripts for further

437 evaluation. These metrics enable the researcher to quickly determine if more reads are needed
438 and whether there are highly expressed putative novel transcripts potentially worth detailed
439 experimentation.

440

441 The examples presented in this work demonstrate that SQANTI-reads is flexible and
442 customizable, allowing users to explore the impact of various experimental design factors on
443 read, UJC, and donor/acceptor properties, as well as identifying potential novel transcripts. The
444 output from SQANTI-reads can be easily mined for additional insights and used to direct
445 attention and resources toward interesting and novel features of lrrNA-seq experiments. We
446 expect SQANTI-reads to become an essential tool for the QC of multisample lrrNA-seq datasets.

447

448

449 **METHODS**

450

451 **SQANTI-reads basics**

452 SQANTI-reads is an adaptation of SQANTI3 designed to evaluate individual reads rather than
453 transcript models and is available in SQANTI3 version 5.3.0. It allows for the comparison of
454 multiple samples, providing quality control results across the entire experiment. Several new
455 features have been introduced to address the specific needs of QC in multisample experiments,
456 while some functionalities of SQANTI3 have been removed as they are not applicable to read-
457 level processing. Table 2 highlights the major differences between SQANTI3 and SQANTI-reads.
458 The input files for SQANTI-reads include: 1) a GTF file of read alignments, 2) a reference genome

459 FASTA file, 3) a GTF file of the reference transcript model annotation, and 4) a design file
460 containing metadata for multiple samples. The first step of SQANTI-reads involves using the
461 SQANTI3 QC module to generate SQANTI3-like classification and junction files, with the
462 classification file containing one row for each mapped read. Reads are classified according to
463 the SQANTI categories (Tardaguila et al. 2018) as full-splice match (FSM), incomplete-splice
464 match (ISM), novel-in-catalog (NIC), novel-not-in-catalog (NNC), antisense, fusion, genic
465 genomic, and intergenic. SQANTI3 subcategories are also included, based on 5' and 3' end
466 positions relative to the annotated transcription start sites (TSS) and transcription termination
467 sites (TTS) (Pardo-Palacios et al. 2024a). Additionally, the reverse transcriptase (RT) switching
468 algorithm of SQANTI3 identifies reads with evidence of RT switching events, while reads with
469 more than 60% adenines in the 20 bp downstream of the reported TTS at the genomic level are
470 flagged as potential intrapriming events. The length of each read and the number of exons in
471 each read are also recorded in the classification file.

472

473 **Junction Metrics**

474 The SQANTI-reads junction file follows the same format as the SQANTI3 junction file, with each
475 row representing a junction in a read, including the start and end positions of the junction. The
476 distance from the junction start and end to the nearest annotated junction start and end in the
477 reference GTF is calculated. It's important to note that the nearest annotated start and end
478 positions may not belong to the same annotated junction. SQANTI classifies junctions as known
479 or novel, and as canonical or non-canonical, based on the dinucleotide pairs at the junction's
480 start and end. By default, dinucleotide combinations of GT-AG, GC-AG, and AT-AC are

481 considered canonical, while any other combinations are classified as non-canonical, although
482 the user can specify additional canonical sites.

483
484 SQANTI-reads introduces new metrics to evaluate the relationship between the junctions in
485 mapped reads and the annotated donors and acceptors. In the SQANTI3 junction file, the
486 distance from each donor/acceptor in each read to the nearest annotated donor/acceptor is
487 recorded. In SQANTI-reads, the mean absolute distance in nucleotides from the annotated
488 donor/acceptor site, the standard deviation, and the coefficient of variation (CV = standard
489 deviation/mean) are calculated and included in the cv.csv file. Each detected junction is
490 classified as: 1) Reference Match junction if the mean distance and the standard deviation to an
491 annotated junction are both equal to 0; 2) CV = 0 junction when the mean distance is greater
492 than 0 and the standard deviation equals 0, and 3) CV > 0 junction when the CV is greater than
493 0.

494

495 **Unique Junction Chain and Gene-Level Information**

496 SQANTI-reads groups mapped reads based on their full junction pattern and refers to them as
497 Unique Junction Chains (UJCs). Each UJC is labeled with a string that includes the chromosome
498 and junction coordinates (Nanni et al. 2024). To enhance computational efficiency, UJC strings
499 are encoded as an index in a hash table (JxnHash). The read count for each JxnHash is
500 calculated and included in the ujc_counts.csv file. Additionally, the number of known canonical,
501 known non-canonical, novel canonical, and novel non-canonical junctions within each UJC is
502 annotated, along with the SQANTI structural category of the UJC. The number of reads within

503 each structural category for each gene, as well as the total number of reads per gene, is stored
504 in the summary file `gene_counts.csv`.

505

506 **Identifying Genes That May Be Under-Annotated and Transcripts That May Be Mis-Annotated**

507 For expressed genes, a high proportion of reads from a UJC classified as NIC/NNC may indicate
508 the presence of a potentially novel transcript. SQANTI-reads includes a customizable pipeline to
509 identify genes with such potential under-annotation events. The procedure identifies NIC/NNC
510 UJCs present in genes with a minimum number of reads (R) and representing a minimum
511 proportion (P) of reads in the gene, with default values set at 100 reads and 20%, respectively.

512 To mitigate the risk that the NIC/NNC UJC is merely a degradation product, an additional
513 condition is applied: the candidate UJC must include at least 80% of the gene's junctions (Q).

514 The R , P , and Q thresholds are pipeline parameters that can be adjusted by the user.

515 Furthermore, SQANTI-reads allows for the evaluation of under-annotated genes and novel
516 transcripts within a specific subset of samples associated with a particular experimental factor
517 (e.g., developmental stage or technology) using the `--factor-level` option.

518

519 **Multisample Processing**

520 SQANTI-reads processes multiple samples to generate classification and junction files when a
521 design file (e.g. Supplementary File 2) is provided to the `sqanti-reads.py` command. If individual
522 samples have already been pre-processed with SQANTI3, SQANTI-reads can be run in `--fast`
523 mode, where the design file links the individual classification and junction files to sample IDs for
524 the calculation of SQANTI-reads metrics, summaries, and a series of visualizations. If pre-

525 processing has not been done, SQANTI-reads is run in --simple mode, where SQANTI3 is run on
526 each sample, followed by the calculation of SQANTI3 metrics and summaries. The output also
527 includes a summary for each sample, reporting the mean, median, upper quartile, and lower
528 quartile of mapped read length, as well as the number and proportion of reads that are shorter
529 than 1 kb, between 1 and 2 kb, between 2 and 3 kb, and greater than 3 kb in length, all of which
530 are included in the length_summary.csv file.

531

532 ***Drosophila melanogaster* data**

533 A total of 24 female *D. melanogaster* abdomen samples corresponding to 2 developmental
534 stages (0-1 hours and 3-8 days post-hatching), four genotypes (dmel 11037, 11255, 12272 and
535 12279) and 3 replicates (2 time points * 4 genotypes * 3 replicates = 24 samples) were
536 sequenced using Oxford Nanopore Technology (ONT). For each sample, mRNA was isolated
537 (DynaBeads mRNA direct kit) from a pool of ~20 abdomens. ONT libraries were constructed
538 using the ONT PCR-cDNA Barcoding Kit (SQK-PCB109) starting with poly(A) mRNA according to
539 the manufacturer's protocol. Libraries were pooled to a total of 100fmol and run on a MinION
540 Mk1c with real-time basecalling and demultiplexing (Guppy v6.1.5, MinKNOW v22.05.8). Read
541 length and quality was evaluated with all samples passing the PycoQC metrics (v2.5.2, Leger
542 and Leonardi 2019). Based on the MinION read counts, libraries were repooled prior to
543 obtaining additional sequencing data on the ONT PromethION (Guppy v5.1.13, MinKNOW
544 v23.04.5) at the University of Florida Interdisciplinary Center for Biotechnology Research (ICBR).
545 Technical replicates (TRs) are defined as the same library run on different ONT flow cells
546 (MinION or PromethION). TRs 1-3 were run on the MinION and TRs 4-6 were run on the

547 PromethION. Detailed metadata for these samples and technical replicates are provided in
548 Supplementary Table 1.

549
550 Dorado basecalled reads were generated from the FAST5 files by converting to pod5 formats
551 (pod5 v 0.3.6) prior to basecalling by Dorado (v 0.5.2)
552 (<https://github.com/nanoporetech/dorado>) using options --recursive --device "cuda:0,1" --kit-
553 name SQK-PCB109 --trim none. Reads were demultiplexed using the demux mode of Dorado (v
554 0.5.2) with options --no-classify --emit-fastq.

555
556 Both Guppy and Dorado basecalled reads were processed using pypochopper (v 2.7.1). Re-
557 oriented FASTQ files were aligned to *D. melanogaster* 6.50 using minimap2 (v 2.17) (Li 2018)
558 and the resulting SAM files converted to GTF format using SAMtools (v 1.10) (Li et al. 2009) and
559 BEDtools (v 2.29.2) (Quinlan and Hall 2010).

560
561 The resulting 67 GTF files ((5 samples * 2 TRs) + (19 samples * 3 TRs) = 67) were used as input
562 into SQANTI-reads along with the *D. melanogaster* 6.50 FASTA and GTF reference files
563 (https://ftp.flybase.net/releases/FB2023_01/dmel_r6.50/, (Öztürk-Çolak et al. 2024)) and a
564 design file (Supplementary File 3). The design file included all experimental factors for
565 evaluation including time, genotype, sequencing platform (MinION and PromethION) and
566 basecaller (Dorado vs Guppy).

567

568 **Human Cell Line WTC11**

569 We used publicly available lrrRNA-seq data from the Long-read RNA-seq Genome Annotation
570 Assessment Project (Pardo-Palacios et al. 2024b) to illustrate the utility of SQANTI-reads.
571 Specifically, we used triplicate measurements of the transcriptome of the WTC11 human cell
572 line that were profiled by cDNA PacBio Sequel II, cDNA Oxford Nanopore MinION, and direct
573 RNA Oxford Nanopore MinION methods. Data were downloaded from the ENCODE website
574 (https://www.encodeproject.org/search/?type=Experiment&internal_tags=LRGASP). Accession
575 numbers for these samples are provided in Supplementary Table 2. The FASTQ files were pre-
576 processed by LRGASP researchers as described in (Pardo-Palacios et al. 2024b). We used the
577 GTF files of read alignments, GENCODE's GRCh38.p13 reference genome GTF and FASTA for
578 release 38 (https://www.gencodegenes.org/human/release_38.html), and a design file
579 (Supplementary File 4) to run SQANTI-reads on the WTC11 samples. The SQANTI-reads output
580 is provided in Supplementary File 5.

581

582 **Computational efficiency evaluation**

583 SQANTI-reads can operate in two modes: slow and fast. The slow mode uses SQANTI3 to
584 generate the standard classification file, which serves as the input for further analysis. In
585 contrast, the fast mode begins with preprocessed samples that have already been run through
586 SQANTI3. This mode generates a unique junction chain and hash dictionary for each sample to
587 facilitate comparisons and calculate summary statistics.

588 To evaluate the performance of SQANTI-reads, we generated 25 datasets representing GTF files
589 with varying numbers of reads. To do so, we used the shuf command line tool, with the -n flag,
590 to extract a random set of 10,000, 50,000, 100,000, 500,000, and 1,000,000 reads from the

591 ENCF003QZT sample of the LRGASP study (Pardo-Palacios et al. 2024b). For each number of
592 reads, five different random subsets of reads were extracted, representing five replicates.
593 These datasets were first processed independently using SQANTI3 to produce the necessary
594 classification files required for fast-mode analysis, as the SQANTI3 performance does not reflect
595 the actual runtime or memory usage of SQANTI-reads itself.

596 Following preprocessing, the datasets were tested under 15 experimental conditions designed
597 to evaluate runtime and memory usage across different scales. The conditions included
598 processing (1) a single sample of each dataset size, (2) three samples of each size, and (3) five
599 samples of each size. All experiments were conducted on a Linux High Performance Computing
600 Cluster requesting 10 GB of RAM and a single CPU.

601 We measured runtime as the wall-clock time from the start to the completion of SQANTI-reads
602 execution and memory usage using sacct command from Slurm workload manager. These
603 metrics were recorded for each experimental condition to assess scalability and resource
604 efficiency under varying dataset sizes and sample counts.

605

606 **Software availability**

607 SQANTI-reads is available for download on GitHub and is integrated into SQANTI3 from version
608 5.3.0 (<https://github.com/ConesaLab/SQANTI3>). Instructions for running SQANTI-reads are
609 provided on the SQANTI3 wiki (<https://github.com/ConesaLab/SQANTI3/wiki/Running-SQANTI%E2%80%90reads>).

611 Code for generating subsampled datasets to evaluate computational efficiency, along with
612 scripts for analysis, is available on GitHub

613 (https://github.com/ConesaLab/SQANTI_reads_ComputationalEfficiency_plots).

614

615 **Data access**

616 *Drosophila* long-read data have been submitted to the NCBI BioProject database

617 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under the accession number PRJNA1134728. The

618 WTC11 long-read data from the LRGASP project, are available from the ENCODE Portal

619 (<https://www.encodeproject.org/>) and accession numbers are provided in Supplementary

620 Table 2.

621 Subsampled datasets to evaluate computational efficiency, are available on GitHub

622 (https://github.com/ConesaLab/SQANTI_reads_ComputationalEfficiency_plots).

623

624 **Acknowledgements**

625 N.K. conducted research, developed software implementations, created figures, and drafted

626 the manuscript. C.M. developed software implementations, integrated new methods into

627 SQANTI3 GitHub, corrected figures and helped with data analysis. L.M. conceived and

628 supervised the study and drafted the manuscript. A.C. supervised the study, contributed to

629 conceptualizations, and drafted the manuscript.

630

631 This work was supported in part by a grant from the National Institutes of Health

632 (1R21HG011280-01), the Spanish MICIN (PID2020-119537RB-I00), the European Union's

633 programme Horizon Europe under the Marie Skłodowska-Curie Actions postdoctoral

634 fellowship to C.M. (101149931). This work is also supported in part by a grant from the

635 National Cancer Institute (NCI P01 CA214091), the National Institute of General Medical
636 Sciences (NIGMS GM137430), the University of Florida Department of Molecular Genetics and
637 Microbiology, the University of Florida Genetics Institute, the University of Florida Cancer
638 Center, and the University of Florida Research Computing Center (www.rc.ufl.edu) and the
639 Latin American and Caribbean Scholars award to N.K. Part of the computations were
640 performed on the high-performance computing cluster Garnatxa at the Institute for
641 Integrative Systems Biology (I2SysBio). I2SysBio is a joint research center formed by University
642 of Valencia (UV) and Spanish National Research Council (CSIC). We acknowledge Ashley Myrick
643 for help with some of the initial CV plots and Knife Bankole for the initial coding of the junction
644 hash. Alison Morse prepared all samples and libraries for the *Drosophila* experiment, ran all
645 initial QC analyses, and recalled all the bases for ONT data. We acknowledge Rolf Renne for his
646 support.

647

648 **Competing interest statement**

649 A.C. has received in-kind funding from Pacific Biosciences for library preparation and
650 sequencing. A.C. collaborates with Oxford Nanopore in the Marie Skłodowska-Curie Actions
651 Doctoral Network project LongTREC (Grant Agreement # 101072892).

652

653 **References**

- 654 Al'Khafaji AM, Smith JT, Garimella KV, Babadi M, Popic V, Sade-Feldman M, Gatzem M, Sarkizova S,
655 Schwartz MA, Blaum EM et al. 2024. High-throughput RNA isoform sequencing using
656 programmed cDNA concatenation. *Nature Biotechnology* **42**: 582-586.
657 Amarasinghe SL, Ritchie ME, Gouil Q. 2021. long-read-tools.org: an interactive catalogue of analysis
658 methods for long-read sequencing data. *GigaScience* **10**.
659 Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. 2020. Opportunities and challenges in long-

- 660 read sequencing data analysis. *Genome Biol* **21**: 30.
- 661 Auer PL, Doerge RW. 2010. Statistical Design and Analysis of RNA Sequencing Data. *Genetics* **185**: 405-
662 416.
- 663 Begik O, Diensthuber G, Liu H, Delgado-Tejedor A, Kontur C, Niazi AM, Valen E, Giraldez AJ, Beaudoin JD,
664 Mattick JS et al. 2023. Nano3P-seq: transcriptome-wide analysis of gene expression and tail
665 dynamics using end-capture nanopore cDNA sequencing. *Nat Methods* **20**: 75-85.
- 666 Carbonell-Sala S, Perteghella T, Lagarde J, Nishiyori H, Palumbo E, Arnan C, Takahashi H, Carninci P,
667 Uszczynska-Ratajczak B, Guigó R. 2024. CapTrap-seq: a platform-agnostic and quantitative
668 approach for high-fidelity full-length RNA sequencing. *Nature Communications* **15**: 5278.
- 669 De Coster W, D'hert S, Schultz DT, Cruts M, Van Broeckhoven C. 2018. NanoPack: visualizing and
670 processing long-read sequencing data. *Bioinformatics* **34**: 2666-2669.
- 671 Delahaye C, Nicolas J. 2021. Sequencing DNA with nanopores: Troubles and biases. *PLoS one* **16**:
672 e0257521.
- 673 Dent CI, Singh S, Mukherjee S, Mishra S, Sarwade RD, Shamaya N, Loo KP, Harrison P, Sureshkumar S,
674 Powell D et al. 2021. Quantifying splice-site usage: a simple yet powerful approach to analyze
675 splicing. *NAR Genomics and Bioinformatics* **3**.
- 676 Diensthuber G, Prysycz LP, Llovera L, Lucas MC, Delgado-Tejedor A, Cruciani S, Roignant J-Y, Begik O,
677 Novoa EM. 2024. Enhanced detection of RNA modifications and read mapping with high-
678 accuracy nanopore RNA basecalling models. *Genome Research* **34**: 1865-1877.
- 679 Fukasawa Y, Ermini L, Wang H, Carty K, Cheung MS. 2020. LongQC: A Quality Control Tool for Third
680 Generation Sequencing Long Read Data. *G3 (Bethesda)* **10**: 1193-1196.
- 681 Glinos DA, Garborcauskas G, Hoffman P, Ehsan N, Jiang L, Gokden A, Dai X, Aguet F, Brown KL, Garimella
682 K et al. 2022. Transcriptome variation in human tissues revealed by long-read sequencing.
683 *Nature* **608**: 353-359.
- 684 Holmqvist I, Bäckholm A, Tian Y, Xie G, Thorell K, Tang KW. 2021. FLAME: long-read bioinformatics tool
685 for comprehensive spliceome characterization. *Rna* **27**: 1127-1139.
- 686 Joglekar A, Hu W, Zhang B, Narykov O, Diekhans M, Marrocco J, Balacco J, Ndhlovu LC, Milner TA,
687 Fedrigo O et al. 2024. Single-cell long-read sequencing-based mapping reveals specialized
688 splicing patterns in developing and adult mouse and human brain. *Nature Neuroscience* **27**:
689 1051-1063.
- 690 Leger A, Leonardi T. 2019. pycoQC, interactive quality control for Oxford Nanopore Sequencing. *Journal*
691 *of Open Source Software* **4**: 1236.
- 692 Legnini I, Alles J, Karaiskos N, Ayoub S, Rajewsky N. 2019. FLAM-seq: full-length mRNA sequencing
693 reveals principles of poly(A) tail length control. *Nature Methods* **16**: 879-886.
- 694 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094-3100.
- 695 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The
696 Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078-2079.
- 697 Lienhard M, van den Beucken T, Timmermann B, Hochradel M, Börno S, Caiment F, Vingron M, Herwig R.
698 2023. IsoTools: a flexible workflow for long-read transcriptome sequencing analysis.
699 *Bioinformatics* **39**.
- 700 Liu P, Sanalkumar R, Bresnick EH, Keleş S, Dewey CN. 2016. Integrative analysis with ChIP-seq advances
701 the limits of transcript quantification from RNA-seq. *Genome Res* **26**: 1124-1133.
- 702 Mahmoud M, Huang Y, Garimella K, Audano PA, Wan W, Prasad N, Handsaker RE, Hall S, Pionzio A,
703 Schatz MC et al. 2024. Utility of long-read sequencing for All of Us. *Nature Communications* **15**:
704 837.
- 705 Marx V. 2023. Method of the year: long-read sequencing. *Nature Methods* **20**: 6-11.
- 706 Nanni A, Titus-McQuillan J, Bankole KS, Pardo-Palacios F, Signor S, Vlaho S, Moskalenko O, Morse
707 Alison M, Rogers RL, Conesa A et al. 2024. Nucleotide-level distance metrics to quantify

708 alternative splicing implemented in TranD. *Nucleic Acids Research* **52**: e28-e28.

709 Newman JRB, Concannon P, Tardaguila M, Conesa A, McIntyre LM. 2018. Event Analysis: Using
710 Transcript Events To Improve Estimates of Abundance in RNA-seq Data. *G3*
711 *Genes/Genomes/Genetics* **8**: 2923-2940.

712 Öztürk-Çolak A, Marygold SJ, Antonazzo G, Attrill H, Goutte-Gattat D, Jenkins VK, Matthews BB, Millburn
713 G, dos Santos G, Tabone CJ et al. 2024. FlyBase: updates to the Drosophila genes and genomes
714 database. *Genetics* **227**.

715 Pagès-Gallego M, de Ridder J. 2023. Comprehensive benchmark and architectural analysis of deep
716 learning models for nanopore sequencing basecalling. *Genome Biology* **24**: 71.

717 Pardo-Palacios FJ, Arzalluz-Luque A, Kondratova L, Salguero P, Mestre-Tomás J, Amorín R, Estevan-Morió
718 E, Liu T, Nanni A, McIntyre L et al. 2024a. SQANTI3: curation of long-read transcriptomes for
719 accurate identification of known and novel isoforms. *Nature Methods* **21**: 793-797.

720 Pardo-Palacios FJ, Wang D, Reese F, Diekhans M, Carbonell-Sala S, Williams B, Loveland JE, De María M,
721 Adams MS, Balderrama-Gutierrez G et al. 2024b. Systematic assessment of long-read RNA-seq
722 methods for transcript identification and quantification. *Nature Methods* **21**: 1349-1363.

723 Patowary A, Zhang P, Jops C, Vuong CK, Ge X, Hou K, Kim M, Gong N, Margolis M, Vo D et al. 2024.
724 Developmental isoform diversity in the human neocortex informs neuropsychiatric risk
725 mechanisms. *Science* **384**: eadh7688.

726 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
727 *Bioinformatics* **26**: 841-842.

728 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. 2011.
729 Integrative genomics viewer. *Nat Biotechnol* **29**: 24-26.

730 Tardaguila M, de la Fuente L, Marti C, Pereira C, Pardo-Palacios FJ, Del Risco H, Ferrell M, Mellado M,
731 Macchietto M, Verheggen K et al. 2018. SQANTI: extensive characterization of long-read
732 transcript sequences for quality control in full-length transcriptome identification and
733 quantification. *Genome Res* **28**: 396-411.

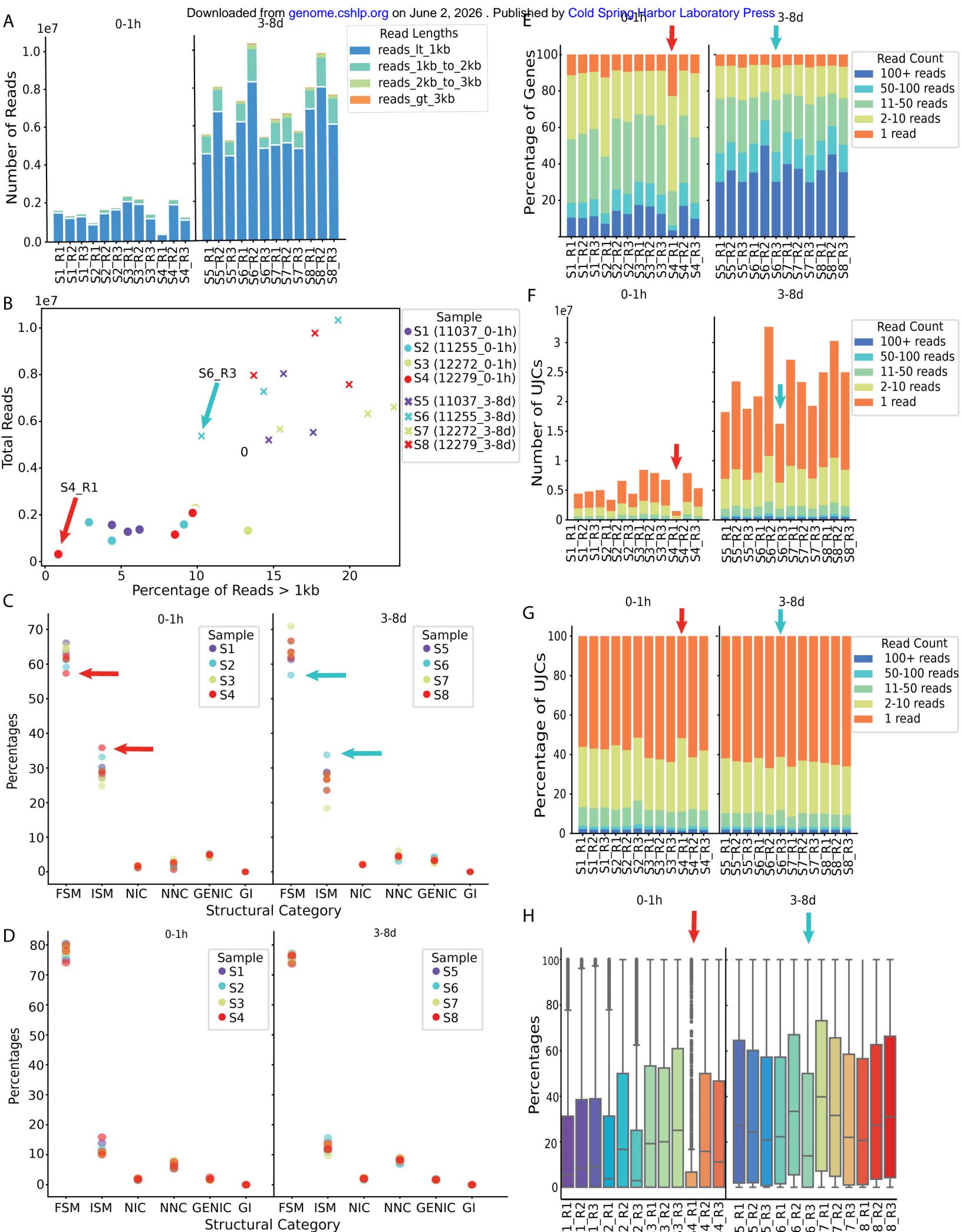
734 van Dijk EL, Naquin D, Gorrichon K, Jaszczyszyn Y, Ouazahrou R, Thermes C, Hernandez C. 2023.
735 Genomics in the long-read sequencing era. *Trends in Genetics* **39**: 649-671.

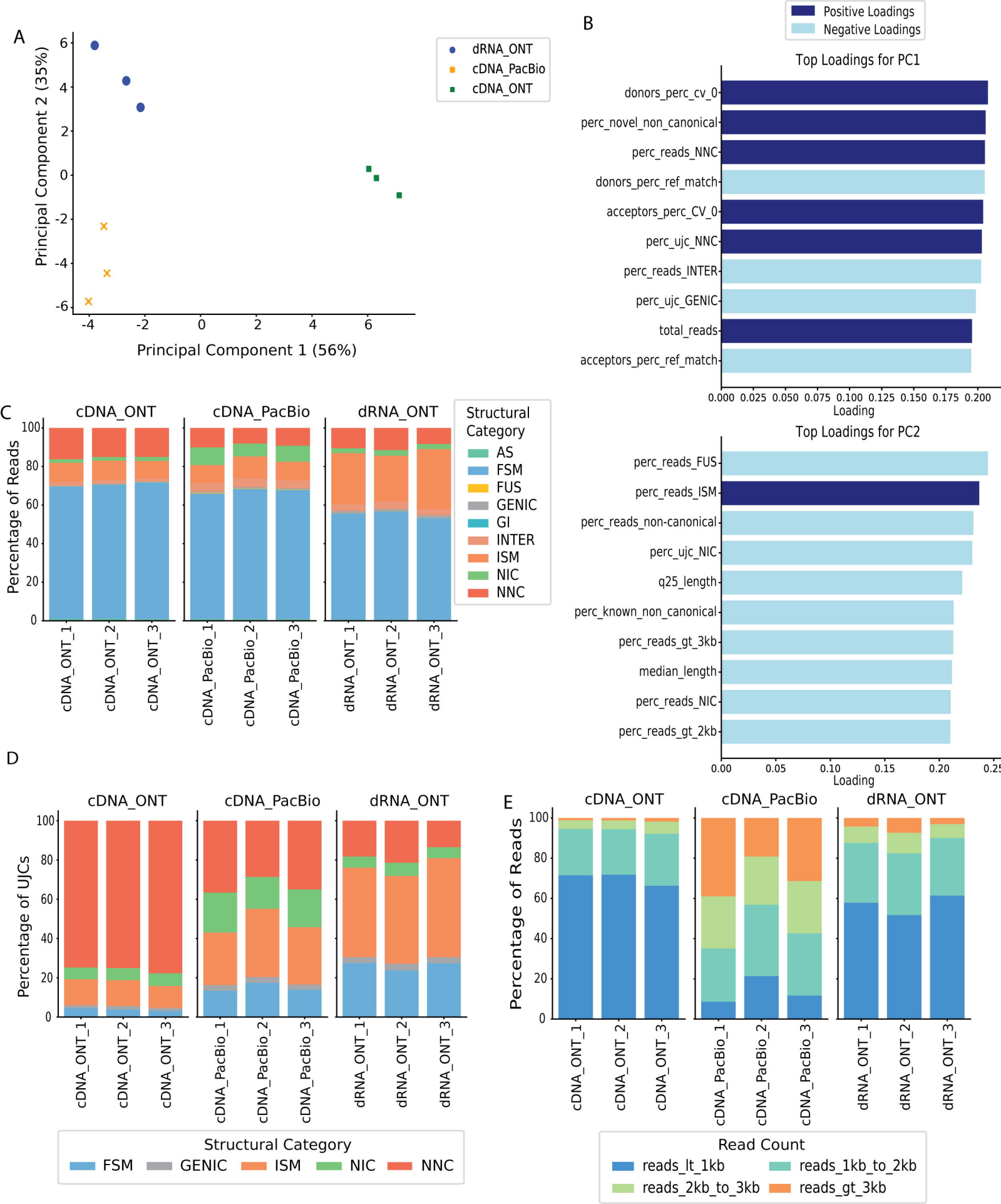
736 Volden R, Palmer T, Byrne A, Cole C, Schmitz RJ, Green RE, Vollmers C. 2018. Improving nanopore read
737 accuracy with the R2C2 method enables the sequencing of highly multiplexed full-length single-
738 cell cDNA. *Proceedings of the National Academy of Sciences* **115**: 9726-9731.

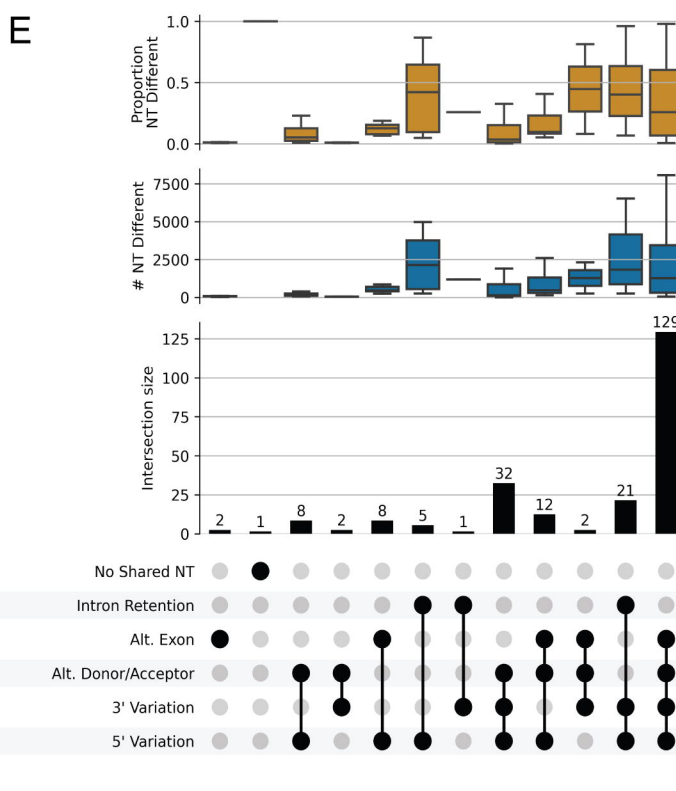
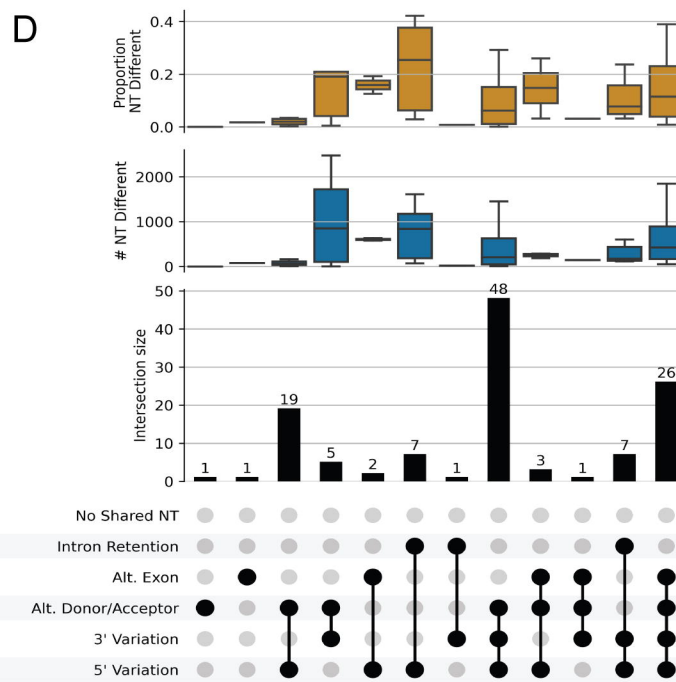
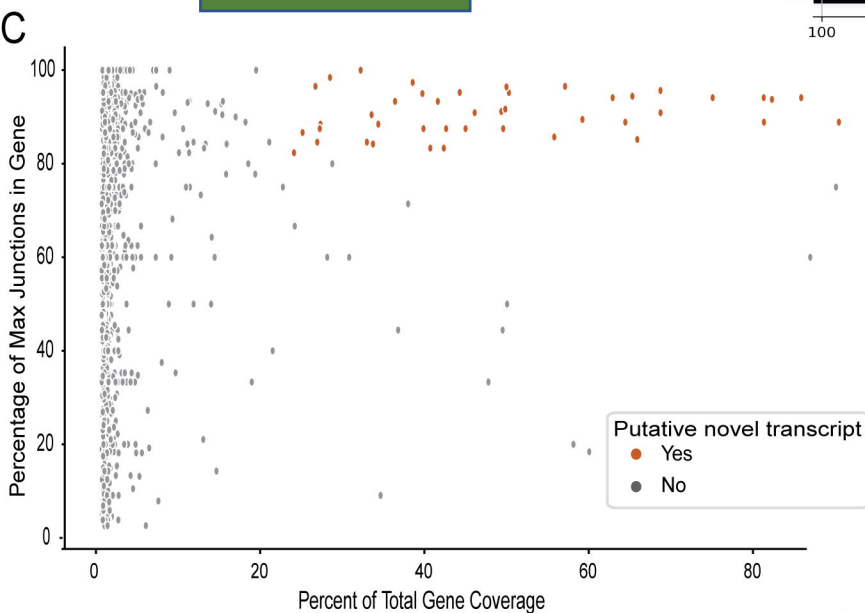
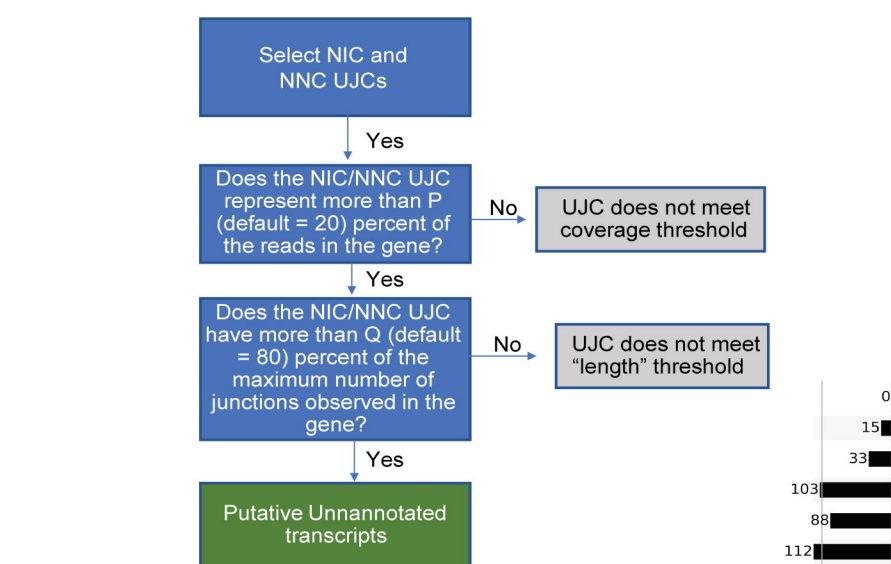
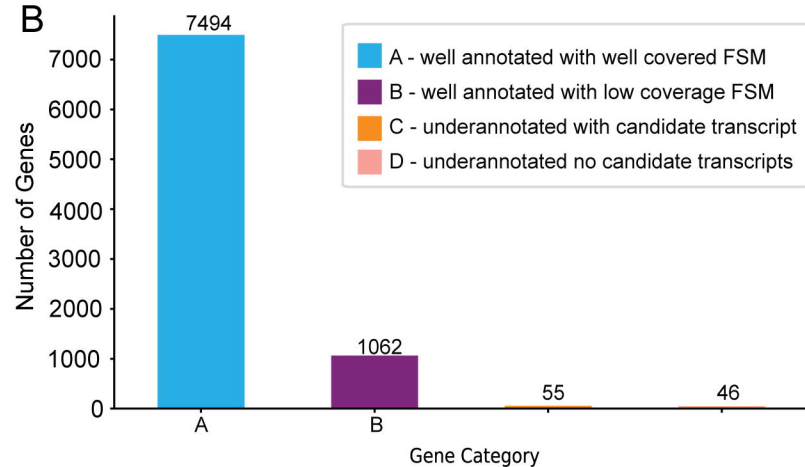
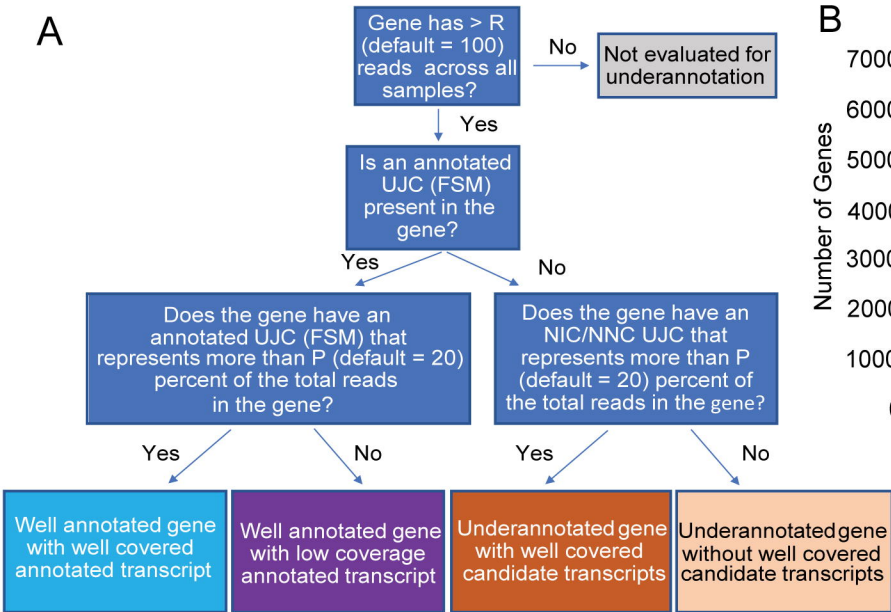
739 Wang M, Marín A. 2006. Characterization and prediction of alternative splice sites. *Gene* **366**: 219-227.

740 Wyman D, Balderrama-Gutierrez G, Reese F, Jiang S, Rahmanian S, Forner S, Matheos D, Zeng W,
741 Williams B, Trout D et al. 2020. A technology-agnostic long-read analysis pipeline for
742 transcriptome discovery and quantification. *bioRxiv* doi:10.1101/672931: 672931.

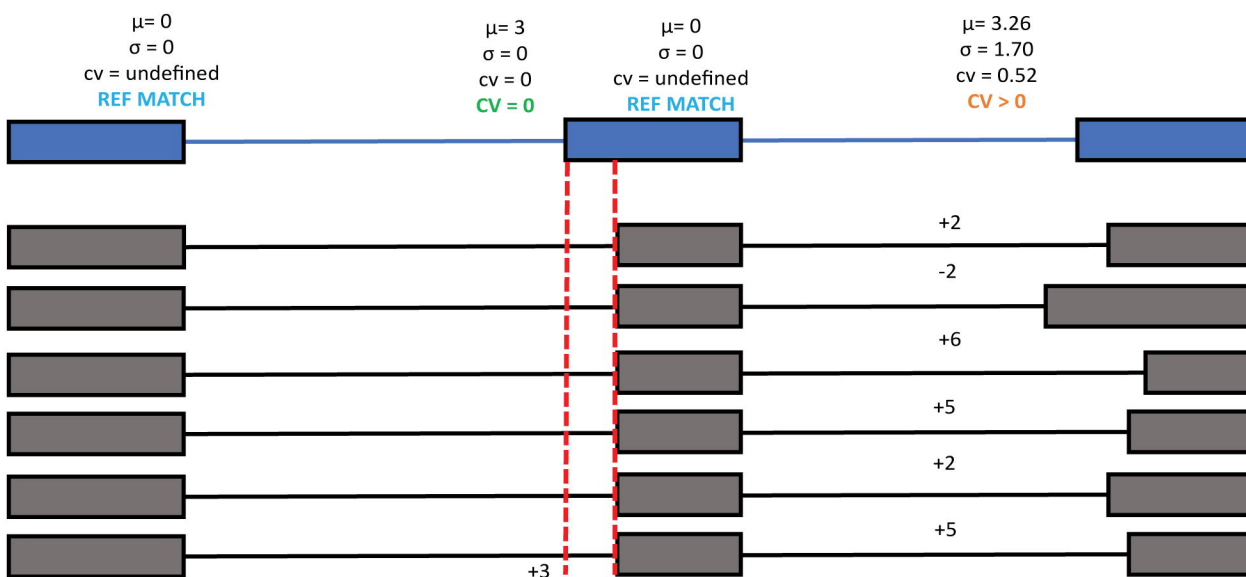
743



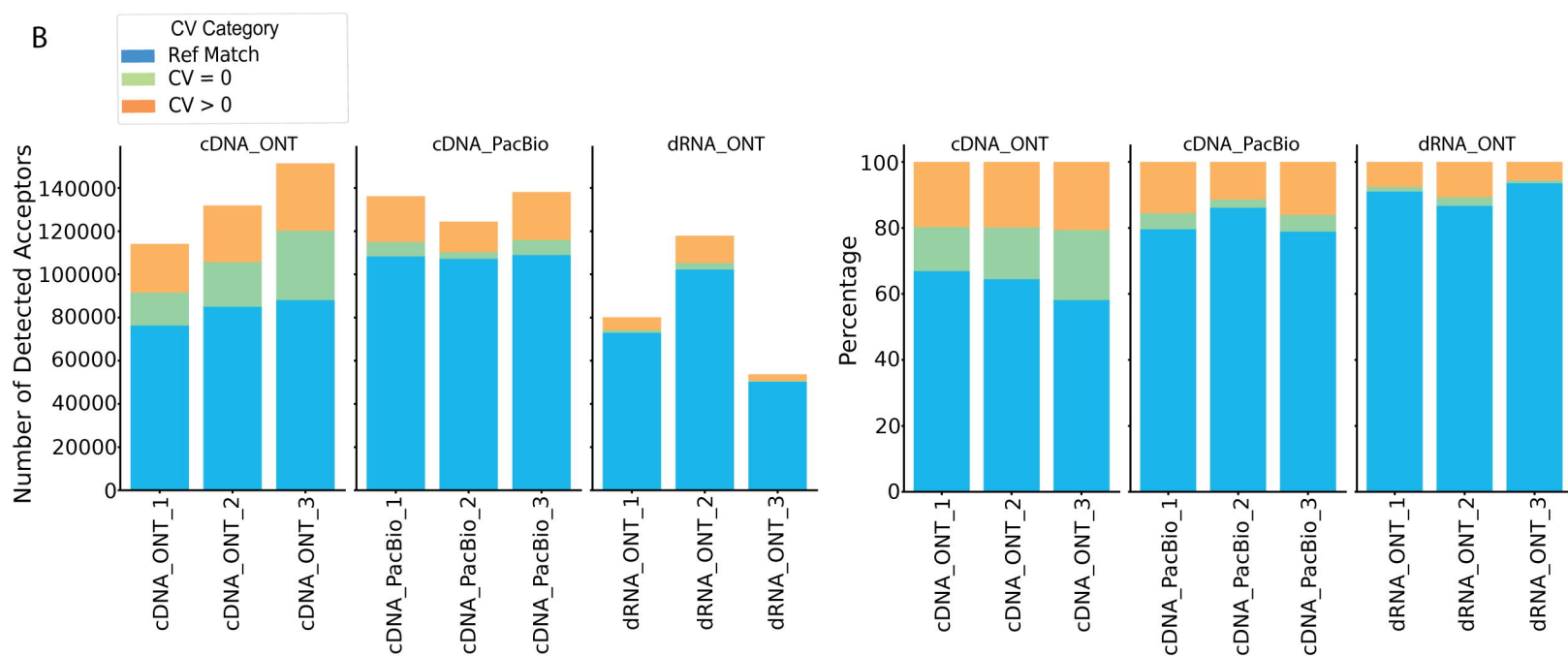




A



B



C

