

# Nanopore sequencing of RNA and cDNA molecules expands the transcriptomic toolbox in prokaryotes

Felix Grünberger<sup>1\*</sup>, Sébastien Ferreira-Cerca<sup>2,3</sup>, Dina Grohmann<sup>1,2\*</sup>

<sup>1</sup> Institute of Biochemistry, Genetics and Microbiology, Institute of Microbiology and Archaea Centre, Single-Molecule Biochemistry Lab & Biochemistry Centre Regensburg, University of Regensburg, Universitätsstraße 31, 93053 Regensburg, Germany

<sup>2</sup> Regensburg Center of Biochemistry (RCB), University of Regensburg, 93053 Regensburg, Germany

<sup>3</sup> Institute for Biochemistry, Genetics and Microbiology, Regensburg Center for Biochemistry, Biochemistry III, University of Regensburg, Universitätsstraße 31, 93053 Regensburg, Germany

\* To whom correspondence should be addressed.

Email: felix.gruenberger@ur.de

dina.grohmann@ur.de

## **Abstract**

High-throughput sequencing dramatically changed our view of transcriptome architectures and allowed for groundbreaking discoveries in RNA biology. Recently, sequencing of full-length transcripts based on the single-molecule sequencing platform from Oxford Nanopore Technologies (ONT) was introduced and is widely employed to sequence eukaryotic and viral RNAs. However, experimental approaches implementing this technique for prokaryotic transcriptomes remain scarce. Here, we present an experimental and bioinformatic workflow for ONT RNA-seq in the bacterial model organism *Escherichia coli*, which can be applied to any microorganism. Our study highlights critical steps of library preparation and computational analysis and compares the results to gold standards in the field. Furthermore, we comprehensively evaluate the applicability and advantages of different ONT-based RNA sequencing protocols, including direct RNA, direct cDNA, and PCR-cDNA. We find that cDNA-seq offers improved yield and accuracy without bias in quantification compared to direct RNA sequencing. Notably, cDNA-seq can be readily used for simultaneous transcript quantification, accurate detection of transcript 5' and 3' boundaries, analysis of transcriptional units and transcriptional heterogeneity. In summary, we establish Nanopore RNA-seq to be a ready-to-use tool allowing rapid, cost-effective, and accurate annotation of multiple transcriptomic features thereby advancing it to become a standard method for RNA analysis in prokaryotes.

**Keywords:** Nanopore, prokaryotes, RNA-seq, transcriptome

## **Introduction**

In the last decade, next-generation sequencing (NGS) technologies (1) revolutionised the field of microbiology (2), which is not only reflected in the exponential increase in the number of fully sequenced microbial genomes but also in the detection of microbial diversity in many hitherto inaccessible habitats based on metagenomics. Using transcriptomics, important advances were also possible in the field of RNA biology (3, 4) that shaped our understanding of the transcriptional landscape (5, 6) and RNA-mediated regulatory processes in prokaryotes (7). RNA sequencing (RNA-seq) technologies can be categorised according to their platform-dependent read lengths and the necessity of a reverse transcription and amplification step to generate cDNA (8).

Illumina sequencing yields highly accurate yet short sequencing reads (commonly 100-300 bp). Hence, sequence information is only available in a fragmented form, making full-length transcript- or isoform-detection a challenging task (9, 10). Sequencing platforms developed by Pacific Bioscience (PacBio) and Oxford Nanopore Technologies (ONT) solved this issue. Both sequencing methods are *bona fide* single-molecule sequencing techniques that allow the sequencing of long DNAs or RNAs (11, 12). However, the base detection differs significantly between the two methods. PacBio-sequencers rely on fluorescence-based single-molecule detection that identifies bases based on the unique fluorescent signal of each nucleotide during DNA synthesis by a dedicated polymerase (12). In contrast, in an ONT sequencer, the DNA or

RNA molecule is pushed through a membrane-bound biological pore with the aid of a motor protein attached to the pore protein called a nanopore. A change in current is caused by the translocation of the DNA or RNA strand through this nanopore, which serves as a readout signal for the sequencing process. Due to the length of the nanopore (version R9.4), a stretch of approximately five bases contributes to the current signal. Notably, only ONT-based sequencing offers the possibility to directly sequence native RNAs without the need for prior cDNA synthesis and PCR amplification (13). Direct RNA sequencing based on the PacBio platform has also been realised but requires a customised sequencing workflow using a reverse transcriptase in the sequencing hotspot instead of a standard DNA polymerase (14). Direct RNA-seq holds the capacity to sequence full-length transcripts and has been demonstrated as a promising method to discriminate and identify RNA base modifications (e.g. methylations (15–19)). ONT sequencing is a *bona fide* single-molecule technique and hence offers the possibility to detect molecular heterogeneity in a transcriptome (20). Recently, the technology was exploited to sequence viral RNA genomes (21–24) to gain insights into viral and eukaryotic transcriptomes (25–27) and to detect and quantify RNA isoforms in eukaryotes (19, 20, 28–30). Essentially, the requirements, but also the possibilities in eukaryotes and prokaryotes, are the same (31), with a poly(A) tail being an essential prerequisite, which is required to capture the RNAs. Using enzymatic polyadenylation of prokaryotic RNAs that in general lack poly(A) tails, the applicability of Nanopore RNA-seq has already been demonstrated by metatranscriptomic sequencing of bacterial food pathogens (32) and by accurate estimation of gene expression levels in *Klebsiella pneumoniae* (33). Despite these initial studies, a comprehensive analysis of the applicability of Nanopore RNA-seq for the analysis of prokaryotic transcriptomes is lacking.

In this study, we applied and benchmarked all currently available ONT library preparation methods to analyse RNAs in the prokaryotic model organism *Escherichia coli* K-12. These include direct sequencing of native RNAs, sequencing of direct cDNAs, and sequencing of PCR-amplified cDNAs. The goal was to create a robust workflow for the simultaneous determination of multiple transcriptional features. To this end, we analysed the reproducibility and comparability of transcript quantification, evaluated the accuracy of transcript boundary identification and the potential of long-read ONT RNA-seq to capture the complexity of bacterial transcriptional units. Noteworthy, due to the single-molecule resolution of ONT sequencing, in-depth analysis of transcription units becomes possible. In addition, we point out practical and technical considerations of the different methods such as the effects of rRNA depletion on the sequencing depth, the possibility to enrich for full-length transcripts in the cDNA protocols and the effects of read trimming.

## **Material & Methods**

**Cell growth and RNA extraction:** *Escherichia coli* K-12 MG1655 cells were grown in rich medium (10 g tryptone, 5 g yeast extract, 5 g NaCl per litre, pH 7.2) to an OD<sub>600nm</sub> of 0.5-0.6. To stabilize RNAs, two

volumes of RNeasy Lysis Buffer (Thermo Fisher Scientific) were immediately added to the cultures and stored at -20°C until cells were harvested by centrifugation at 4°C.

Total RNA of all samples except RNA001 was extracted using RNeasy Mini Kit (Qiagen) according to the manufacturer's instructions. RNA001 RNA was purified using the Monarch® Total RNA Miniprep Kit (New England Biolabs). The integrity of total RNA from *E. coli* was assessed via a Bioanalyzer (Agilent) run using the RNA 6000 Pico Kit (Agilent), and only RNAs with RNA integrity numbers (RIN) above 9.5 were used for subsequent treatments and sequencing.

**Poly(A) tailing, rRNA depletion and additional RNA treatment:** Next, RNAs were heat incubated at 70°C for 2 min and snap cooled on a pre-chilled freezer block before polyadenylating RNAs using the *E. coli* poly(A) polymerase (New England Biolabs). Briefly, 5 µg RNA, 20 units poly(A) polymerase, 5 µl reaction buffer and 1 mM ATP were incubated for 15 min at 37°C in a total reaction volume of 50 µl. To stop and clean up the reaction, poly(A)-tailed RNAs were purified following the RNeasy Micro clean-up protocol (Qiagen), which was used for all subsequent RNA clean-ups. The efficiency of poly(A)-tailing was evaluated via a Bioanalyzer run. Ribosomal RNA (rRNA) depletion was performed using the Pan-Prokaryote riboPOOL by siTOOLS, which effectively removes rRNAs from *E. coli*. For TEX-treated samples, partial digestion of RNAs that are not 5'-triphosphorylated (e.g. tRNAs, rRNAs) was achieved by incubation of the RNA with a Terminator 5'-Phosphate-Dependent Exonuclease (TEX, Lucigen). Therefore, 10 µg of RNA used in the RNA001 sample, were incubated with 1 unit TEX, 2 µl TEX reaction buffer and 0.5 µl RiboGuard RNase Inhibitor (Lucigen) in a total volume of 20 µl for 60 minutes at 30°C. Besides, 20 ng of rRNA-depleted samples subsequently used in the PCR-cDNA workflow (replicate 4 and 5), were only partially TEX-treated using the same enzyme and buffer concentrations but reducing the reaction time to 15 minutes. All reactions were terminated by adding EDTA and cleaned up following the RNeasy Micro clean-up protocol. Before library preparation, the extent of the remaining buffer and DNA contamination were tested by performing standard spectroscopic measurements (Nanodrop One) and using the Qubit 1X dsDNA HS assay kit (Thermo Fisher Scientific). Input RNAs were finally quantified using the Qubit RNA HS assay kit.

**Library preparation and sequencing:** Libraries for Nanopore sequencing were prepared from poly(A)-tailed RNAs according to protocols provided by Oxford Nanopore (Oxford Nanopore Technologies Ltd, Oxford, UK) for direct sequencing of native RNAs (SQK-RNA001, SQK-RNA002), direct cDNA native barcoding (SQK-DCS109 with EXP-NBD104) and PCR-cDNA barcoding (SQK-PCB109) with the following minor modifications: Agencourt AMPure XP magnetic beads (Beckman Coulter) in combination with 1 µl of RiboGuard RNase Inhibitor (Lucigen) were used instead of the recommended Agencourt RNAClean XP beads to clean up samples. For reverse transcription, Maxima H Minus Reverse Transcriptase (Thermo Fisher Scientific) was used for all cDNA samples and for the RNA002 samples (SuperScript III Reverse Transcriptase from Thermo Fisher Scientific used for RNA001 sample). The amount of input RNA, barcoding strategy, number of PCR cycles and extension times can be found in Supplementary Table 1 and are also summarized in part in Figure 1A.

Nanopore libraries were sequenced using either a MinION Mk1B connected to a laptop with the recommended specifications for Nanopore sequencing or a Mk1C. All samples were sequenced on R9.4 flow cells and the recommended scripts in MinKNOW to generate fast5 files with live-basecalling enabled. In case of an observed drop in translocation speed and subsequent reduced read quality, the flow

cells were refueled with flush buffer, as recommended by ONT. Flow cells were subsequently washed and re-used for further runs, provided a sufficient number of active pores left. To avoid cross-contamination of reads, a different set of barcodes was used for the next run. Also, the starting voltage of re-used flow cells was adjusted for the next run to account for the voltage drift during a sequencing run.

#### Data analysis:

*Basecalling, demultiplexing of raw reads and quality control of raw reads:* All fast5 reads were re-basecalled using guppy (ont-guppy-for-mk1c v4.3.4) in high-accuracy mode (rna\_r9.4.1\_70bps\_hac.cfg, dna\_r9.4.1\_450bps\_hac.cfg) without quality filtering. While standard parameters were used for basecalling fast5s from cDNA sequencing, fast5 files from RNA sequencing were basecalled with RNA-specific parameters (--calib\_detect, --reverse\_sequence and --u\_substitution). Next, basecalled fastq files from cDNA runs were demultiplexed in a separate step by the guppy suite command guppy\_barcode using default parameters and the respective barcoding kit. After that, relevant information from the guppy sequencing and barcode summary files were extracted to analyse the properties of raw reads (Supplementary Table 1). Please note that in Supplementary Table 2, all figures created from numerical data are referenced and linked to the corresponding code in the Github repository <https://github.com/felixgrunberger/microbepore>.

*Read alignment:* Files were mapped to the reference genome from *Escherichia coli* K-12 MG1655 (GenBank: U00096.3) (34), using minimap2 (Release 2.18-r1015, <https://github.com/lh3/minimap2>) (35). Output alignments in the SAM format were generated with -ax splice -k14 for Nanopore 2D cDNA-seq and -ax splice, -uf, -k14 for direct RNA-seq with i) -p set to 0.99, to return primary and secondary mappings and ii) with --MD turned on, to include the MD tag for calculating mapping identities. Alignment files were further converted to bam files, sorted and indexed using SAMtools (36). To evaluate the alignments, we first calculated the aligned read length by adding the number of M(atch) and I(nsertion) characters in the CIGAR string (13). Based on this, the mapping identity was defined as  $(1 - \text{NM}/\text{aligned\_reads}) * 100$ , where NM is the edit distance reported taken from minimap2. Read basecalling and mapping metrics can be found in Supplementary Table 1. To analyse single reads in more detail with respect to the RNA type (mRNA, rRNA, other ncRNA, unspecified) they map to, bam files were first converted back to FASTQ using bedtools v2.29.2 (37). Next FASTQ files were remapped to a transcriptome file using minimap2 with the previously mentioned parameters to assign single read names with feature IDs. To handle multi-mapping reads, only the mapping location with i) the highest overall identity or if identical ii) the position with most aligned bases was kept for every read id.

*Gene abundance estimation:* A publicly available short-read Illumina dataset (SRR1927169) obtained from RNA-seq data of *E. coli* K-12 grown under rich conditions was downloaded from Gene Expression Omnibus (GEO) GSE67218. Reads were first quality trimmed using trimmomatic v0.39 (38) (leading:20, trailing:20, slidingwindow:4:20, minlen:12) and mapped to the reference genome using bowtie2 (-N 0, -L 26) (39).

SMRT-Cap data obtained from sequencing data from rich-medium samples (SRR7533626, SRR7533627) were downloaded from GEO GSE117273 (40). PacBio reads were processed as described in the SMRT-Cap protocol using the pacbio\_trim.py script downloaded from <https://github.com/elitaone/SMRT-cappable-seq>. In short, reads were filtered and trimmed using the respective filter and poly functions. Next, reads were mapped to the *E. coli* K-12 genome using minimap2 with PacBio-specific (-ax map-pb) options (35). Bam

files from Illumina and SMRT-Cap sequencing were converted to FASTQ format and remapped to the gene file as described before.

To estimate gene abundances from ONT, short-read Illumina and SMRT-Cap libraries, Salmon (v.1.4.0) was applied in alignment-based mode (41). Transcripts per million (TPM) were re-calculated using the salmon-computed effective transcript length, after dropping reads mapping to rRNAs, that are variable between non-depleted and depleted RNA sets.

*Identification and trimming of full-length transcripts:* Full-length cDNA reads containing SSP and VNP primers in the correct orientation were identified using pychopper (v.2.5.0) with standard parameters using the default pHMM backend and autotuned cutoff parameters estimated from subsampled data (<https://github.com/nanoporetech/pychopper>). After a first round, a second round of pychopper was applied to the unclassified direct cDNA reads with DCS-specific read rescue enabled. Reads from rescued and full-length folders were merged and used for subsequent steps. To evaluate the influence of different trimming approaches on the accuracy of transcript boundary analysis, we applied additional 5' and 3' trimming steps using cutadapt v3.2 (42). To this end, polyA sequences were removed from the 3' ends (-a A{10}, -e 1, -j 0) and remaining SSP sequences were removed from the 5' ends (-g TTTCTGTTGGTGCTGATATTGCTGGG, -e 1, -j 0) of direct RNA and full-length cDNA reads. Finally, trimmed reads were mapped using minimap2 as described before. Reads with more than 10 clipped bases on either side were removed from the alignments using samclip (v.0.4.0, <https://github.com/tseemann/samclip>).

To assess the impact of trimmings on gene body coverage, a coverage meta-analysis was performed. First, a transcript file was created for all genes with an ONT-annotated primary 5' and 3' end (see next section). Based on this, strand-specific coverage files were created from the bam files and coverage analysis performed using a custom R script. The genomic coordinates and the counted reads per position were first scaled to values between 0 and 100 and the mean coverage distribution per normalised position was calculated. To evaluate the coverage profiles and the decay at the 5' or 3' ends, we calculated the quartile coefficient of variation (interquartile range/median) (compare 19) and additionally compared the mean coverage in the first and last 10% of the positions to the median values.

*Detection of transcript boundaries:* The determination of enriched 5' and 3' ends was carried out in the same way, but independently of each other, and is briefly explained in the following: First, strand-specific read ends in bedgraph format were created from bam files using bedtools genomecov (-5 or -3 option, -bga) (37). Next, the previously published Termseq\_peaks script (43) was used to call peaks for each sample individually without including replicates (<https://github.com/NICHD-BSPC/termseq-peaks>). This script is based on *scipy.signal.find\_peaks*, which is running in the background of Termseq\_peaks with lenient parameters (prominence=(None,None), width=(1,None), rel\_height=0.75). However, we deliberately used Termseq\_peaks since its ability to include replicates by applying an Irreproducible Discovery Rate method which can be applied to future studies. For end detection, only the leniently called peaks in the narrowPeak file were used after adding the number of counts for each position using bedtools intersect. Enriched positions were finally filtered and annotated based on the following criteria: i) For each peak the position with the highest number of reads was selected. ii) Positions within 20 bases were merged and only the position with the highest number of reads retained. iii) Positions with less than three reads were filtered out. iv) Positions were assigned based on their relative orientation to a gene and their respective peak height as primary (depending on 5' or 3'

detection: highest peak within 300 bases upstream or downstream of a gene, respectively), secondary (each additional peak 300 bases up/downstream of a gene) and internal (each peak in the coding range).

Reproducibility and comparability of primary 5' and 3' ends were evaluated based on Pearson coefficients calculated from pairwise complete observations. Additionally, 5' and 3' untranslated regions (UTR) were calculated based on the distance of the enriched primary site to the start or end of a coding region, respectively. The positions of primary sites called from direct RNA-seq data were corrected by 12 bases.

*Detection and quantification of transcriptional units:* Tables containing each read as a single row were created from the bam files using the R package Genomic alignments (44). Reads that mapped to the opposite strand of an annotated mRNA or ncRNA or that mapped to widely separated genomic positions were discarded. Next, all range overlaps sharing more than 100 bases were defined between the read table and the genomic feature table using the findOverlaps function from the GenomicRanges package. This way, multiple features can be assigned to each individual read. If their genomic positions are adjacent, the combination of features covered by a coverage-dependent number of reads (10 reads for PCR-cDNA replicate 4) are considered as a transcriptional unit. To enable a quantitative assessment of the transcriptional units and the respective context, the number of reads is first determined for each feature individually and then compared with the number of reads in each detected unit. We compared the transcriptional units with the operon tables from the RegulonDB database (45) and the SMRT-Cappable-seq study (40).

#### Public data:

In addition to the publicly available results from the SMRT-Cappable-seq study (40), the short-read Illumina data for gene expression comparison and the RegulonDB (45) mentioned above, we also compared ONT RNA-seq 5' ends with the results of a differential RNA-seq study (46) and 3' ends with Term-seq results (47).

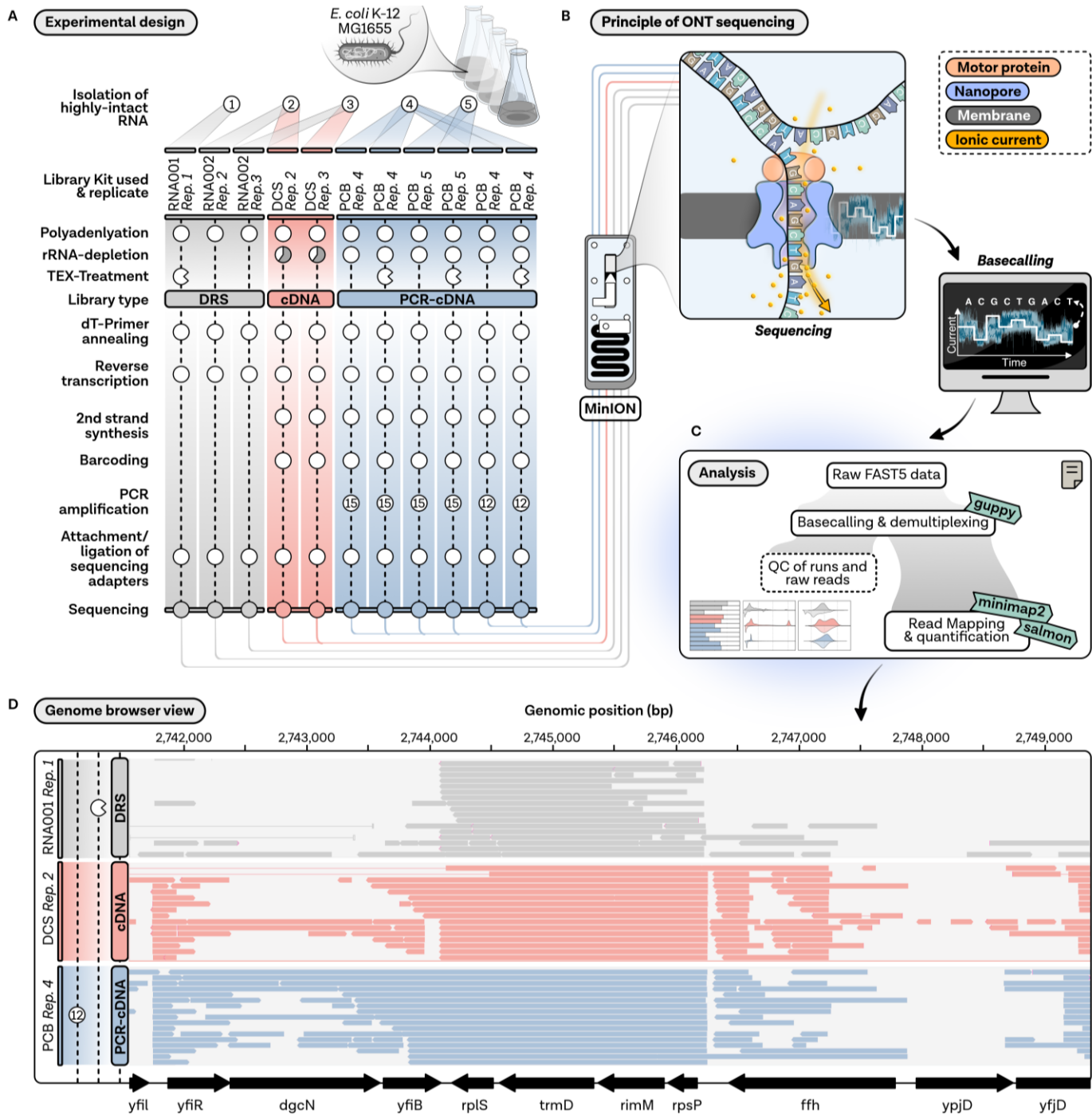
## **Results**

### *Experimental design for benchmarking Nanopore sequencing of RNA and cDNA molecules in prokaryotes*

Currently, three different protocols from ONT are available for the analysis of RNAs including i) direct sequencing of native RNAs (SQK-RNA002, referred to as DRS in this study), ii) direct sequencing of cDNAs (SQK-DCS109, referred to as cDNA in this study) and iii) sequencing of PCR-amplified cDNAs (SQK-PCB109, referred to as PCR-cDNA in this study) (Figure 1A, Supplementary Fig. 1). In short, all methods rely on polyadenylated RNAs as starting material since RNAs are either annealed to an oligo(dT) primer for the cDNA approaches or ligated to a double-stranded oligo(dT) splint adapter in the DRS approach. Although reverse transcription is optional for DRS, it is highly recommended by ONT and the community to resolve secondary structures in the RNA and to decrease the probability of pore blockage, which ultimately results in an increase in total throughput (20). However, only the RNA strand carries the motor protein and is subsequently sequenced. The cDNA protocols take advantage of the template-switching ability of the reverse transcriptase, which adds three non-templated Cs to the end of the cDNA (48). This allows the enrichment of full-length transcripts during the

analysis (Supplementary Fig. 1). After RNA digestion, the second strand is synthesized, followed by barcode ligation, PCR amplification in the PCR-cDNA protocol, attachment or ligation of sequencing adapters and sequencing.

We performed all three protocols using unfragmented total RNA prepared from the prokaryotic model organism *E. coli* K-12 strain MG1655 grown at 37°C in rich medium. The aim was to compare the results to other full-length sequencing protocols and platforms and discuss current limitations and best practices analysing prokaryotic transcriptomes using Nanopore sequencing of RNA and cDNA molecules. Two biological replicates for each library preparation method were sequenced on a MinION using R9.4 flow cells controlled by MinKNOW. The key steps of library preparation and sequencing are depicted in Figure 1A,B and are briefly summarized in the following: after purification of high-quality RNAs using silica-membrane columns with a cut-off size of about 200 nucleotides, RNAs were immediately polyadenylated to make them amenable for library preparation and to preserve the 3' ends from further degradation during the next steps of the library preparation. Since full-length sequencing of RNAs and cDNAs is dependent on the quality of the source material, we only used RNAs with integrity values (RIN) greater than 9.5. Also, Bioanalyzer analysis was used to confirm efficient polyadenylation based on a shift in ribosomal RNA peaks and to check fragment size of PCR-amplified cDNAs (Supplementary Fig. 2A,B). To increase the proportion of sequenced mRNAs, ribosomal RNAs, which usually make up the main part of the RNA pool, can be depleted. However, the input quantity requirements currently still make it challenging to use rRNA-depleted RNAs in a sensible and cost-efficient way, especially for DRS. The input amounts are currently listed to be 500 ng polyA+ (DRS), 100 ng polyA+ RNA (cDNA) and 1 ng (PCR-cDNA), respectively. Therefore, we used non-depleted RNA for DRS sequencing, a mix of depleted (40%) and non-depleted RNA (60%) for the cDNA protocol and fully depleted RNA for the PCR-cDNA approach. Additionally, we tested the compatibility with other RNA treatments using the commonly applied digestion of 5'-monophosphorylated non-primary RNAs with a 5'-Phosphate-dependent Terminator Exonuclease (TEX) as an example (Figure 1A). However, it should be noted that we deliberately chose reaction conditions not sufficient for complete digestion of all non-primary RNAs. The intention of this design was not to distinguish primary from processed transcripts but rather to minimize the rRNA content even further.



**Figure 1: Overview of generated datasets for benchmarking Nanopore sequencing of RNA and cDNA molecules in prokaryotes.** A, Five replicates of the prokaryotic model organism *Escherichia coli* strain K-12 MG1655 were sequenced using currently available RNA-seq protocols from Oxford Nanopore, including direct RNA sequencing (DRS), direct cDNA sequencing (cDNA) and sequencing of PCR-amplified cDNAs (PCR-cDNA). Different rRNA-depletion, additional treatment strategy (Terminator 5'-Phosphate-Dependent Exonuclease, TEX), kit names used (RNA001, RNA002, DCS109, PCB109) and key steps of the library preparation are outlined in the graphic workflow summary. B, Principle of Nanopore sequencing: An ionic current drives the cDNA or the RNA strand of a RNA/cDNA hybrid through the membrane-embedded Nanopore. The motor protein, attached during library preparation, unzips the double strands, and controls the translocation speed. Translocation of the strand alters the electric signal, which is used to determine the sequence. C, Basic workflow for analysing Nanopore reads including basecalling and demultiplexing using ONT-developed guppy, custom scripts to perform quality control of runs/reads, minimap2 (35) to align the reads to the reference genome and salmon (41) in alignment-based mode for gene quantification. D, Colourised integrative genomics viewer (49) snapshot of single reads from different ONT protocols for an exemplary region in *E. coli*.

### *Overall run and raw read characteristics and analysis of mapped reads*

Sequencing throughput on a single FLO-MIN106 flow cell is dependent on the kit chemistry and currently listed by ONT to typically range between 1 to 4 Gb for DRS, more than 8 Gb for cDNA and about 10 Gb for the PCR-cDNA kits. Considering that a higher yield could be expected for the cDNA kits and the (partial) depletion of ribosomal RNAs, cDNA runs were multiplexed and aborted as soon as a sufficient number of reads (> 0.5Gb) was reached. All sequencing parameters and run statistics are listed in Supplementary Table 1 and shown in Supplementary Fig. 3,4. The sequencing yield of unfiltered reads ranged between 0.09 and 2.21 million reads, or 0.08 Gb to 1.57 Gb, respectively (Supplementary Fig. 3). Read qualities, which are specified as mean qscore values, were similarly distributed within the three library types, showing median values of 8.8 (DRS), 9.7 (cDNA) and 10.5 (PCR-cDNA) (Supplementary Fig. 4A). As expected, the read length distributions of the individual samples were highly dependent on the effect of rRNA depletion (Supplementary Fig. 4B). Although we could confirm the reports of previous studies that very short direct RNA reads are associated with bad quality (13, 20), we did not see a pronounced effect in other library types or for very long RNAs in our datasets (Supplementary Fig. 4C).

We next aligned the unfiltered reads to the *E. coli* K-12 genome using minimap2 (Figure 1C). 71.4% (DRS), 64.7% (cDNA) and 48.9% (PCR-cDNA) of the reads mapped to the genome, which corresponds to 78.0% (DRS), 64.7% and 47.2% of the bases, respectively (Supplementary Fig. 5). The moderate numbers arise from short reads with low quality, which dominate the class of unmapped reads and are particularly common for the direct RNA datasets but also occur in the (PCR-)cDNA approaches (Supplementary Fig. 6A-D). The lower total number of mapped reads in the PCR-cDNA samples is due to the preference for over-amplification of short fragments in the PCR, which is more pronounced at higher cycle numbers. This suggests that successful sequencing can be estimated reasonably well already from the Bioanalyzer results of PCR-amplified cDNA (Supplementary Fig. 2B). Based on the length distribution, which is similar to the unamplified cDNA and the proportion of mapped reads (62%), we concluded that 12 PCR cycles are sufficient to obtain high quality sequencing data (Supplementary Fig. 5,6).

Most of the reads in the non-coding RNA class originate from the *ssrA* gene in our sample conditions, producing the transfer-messenger RNA (tmRNA), which explains the uniform length distribution. The tmRNA has tRNA-specific base modifications (50) that lead to an altered current profile, which presumably explains the lower read quality in the DRS approach (Supplementary Fig. 6B). After the RNA is transcribed into cDNA, the modifications are lost, and the quality of the sequenced reads increases significantly. As expected, the raw read length of rRNA-mapping reads was largely dependent on the pre-designed depletion efficiency and subsequent TEX

treatment (Supplementary Fig. 6A). Indeed, the number of reads mapping to ribosomal RNAs is significantly reduced in TEX-treated samples compared to the non-treated counterparts (Supplementary Fig. 7).

In the following, we will focus on mRNA-originating reads performing an in-depth analysis of transcriptomic features in *E. coli*. Reads mapping to mRNAs in fully rRNA-depleted libraries make up about 33% of all mapped reads (PCR-cDNA samples with 12 PCR cycles), which corresponds to 42% of all mapped bases (Supplementary Fig. 7). The aligned read length distribution of mRNA-mapping reads was similar between all library types with median values of 406 (DRS), 372 (cDNA) and 395 (PCR-cDNA) bases (Supplementary Fig. 8A). Therefore, all three library preparation methods allow long-read sequencing of RNAs capturing both ends of transcripts (Figure 1D). These advantages are particularly clear when looking at the mean aligned length of the 100 longest reads in each protocol, which are 4738 (DRS), 6567 (cDNA) and 6132 (PCR-cDNA) bases. At this point, it should be mentioned that the 100 shortest reads have mean lengths of 89 (DRS), 80 (cDNA) and 80 (PCR-cDNA) bases, which is caused by the mapping tool minimap2 used with standard parameters. As previously reported, the mapped read identity of direct RNA reads (88.1%) is substantially lower as compared to cDNA reads (96% cDNA, 94% PCR-cDNA) (Supplementary Fig. 8B) (13). However, we noticed that the read identity improved when using the RNA002 chemistry instead of the meanwhile outdated RNA001 kit. Although template-switching and second-strand synthesis enriches explicitly for full-length transcripts in all cDNA protocols, no clear difference was detected in the aligned length distribution. The difference in the number of PCR cycles leads to significant differences in mean read lengths (15 cycles: 310 bases, 12 cycles: 526 bases), although there is no effect on read quality and identity.

Comparing raw read length with aligned read lengths, we noticed that many reads in the cDNA protocol are twice as large as their mapped counterpart, which is caused by reverse transcription artefacts (51, 52) that generate 2D-like reads, containing both strand of a transcript (Supplementary Fig. 9A). Interestingly, the reverse complement part of the read has much lower quality scores than the reverse transcribed RNA. This is confirmed by a correlation analysis between the raw read qualities and the mapped read identity (Supplementary Fig. 9B). Direct RNA and PCR-cDNA reads with low quality led to a lower identity score. This is not observed in the direct cDNA dataset since the distribution is dominated by the low-quality peak of the reverse complement. In most of the cases, only the good-quality first part of the 2D-like read maps to the genome and the aligned read identity is high. The second part of the read, however, is discarded (Supplementary Fig. 9C,D). However, it should be noted that these 2D-like reads make up the majority of reads that have mapped to mRNAs in the cDNA libraries (Supplementary Fig. 9E). Although much less common, the artefact also occurs in PCR-cDNA reads and as

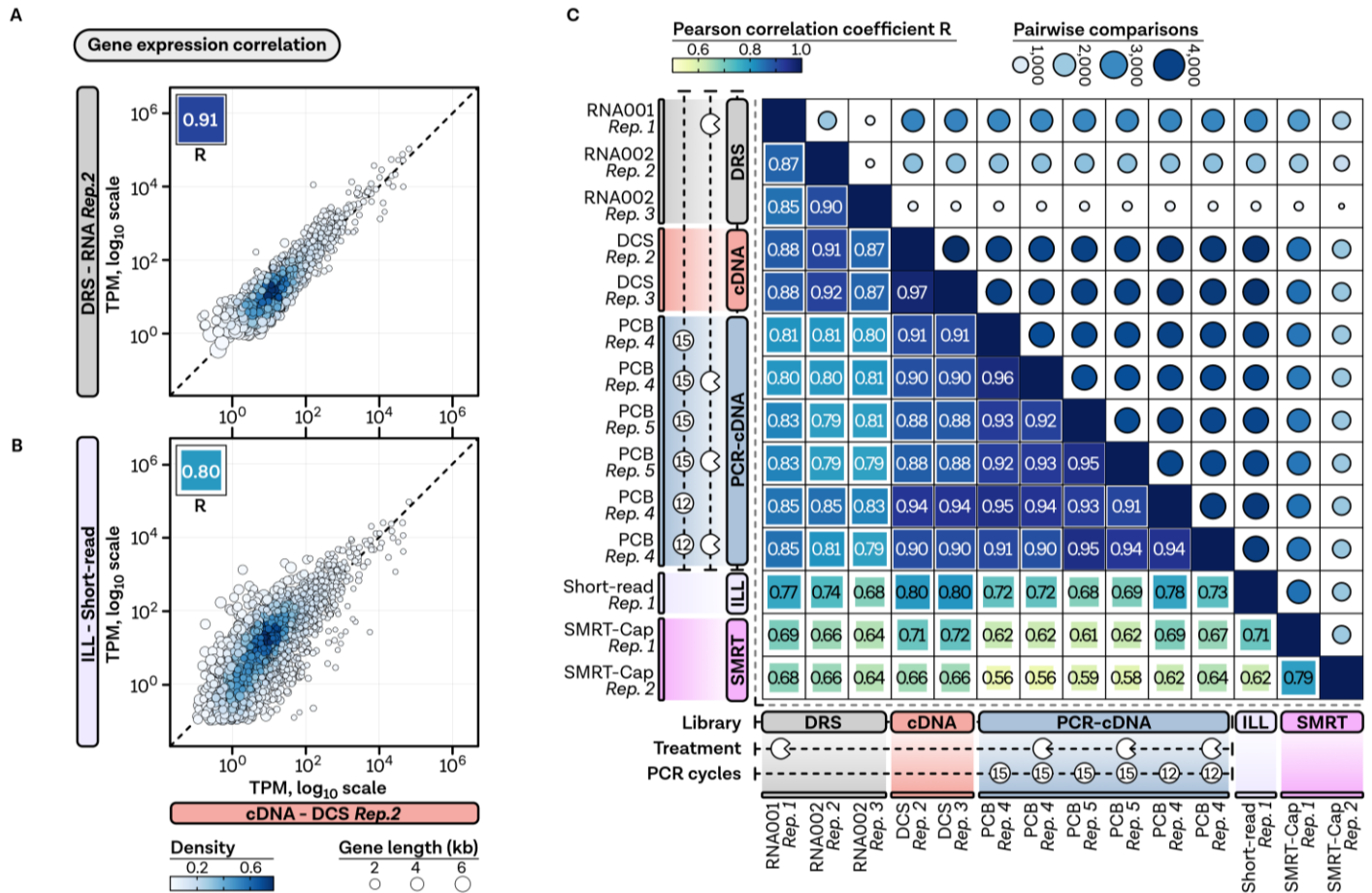
expected, is not found in direct RNA reads (Supplementary Fig. 9E).

#### *Reproducibility and comparability of gene quantification*

Since the first strand is always sequenced, 2D-like reads are not expected to distort the quantification of reads. To test this and to determine the overall comparability and robustness of the data in absolute quantitative terms, we compared the count data based on untrimmed, uncorrected Nanopore reads with published short-read (Illumina) and full-length long-read (SMRT-Cappable-Seq protocol, PacBio) cDNA sequencing data from *E. coli* sampled under very similar conditions (40). Since we only consider reads that map to mRNAs for this purpose, we first looked at the sequencing depth of each dataset to assess whether representative statements can be made. Sequencing depth was dependent on rRNA depletion, TEX treatment and the total number of reads sequenced. Therefore, sequencing depths between 0.2-fold (DRS, RNA002, replicate 2) and 52-fold (PCR-cDNA, 12 cycles, replicate 1) reflect the design of the particular experiment and are mostly comparable to the selected SMRT-Cappable (replicate 1: 51-fold, replicate 2: 7-fold) and short-read Illumina (70-fold) datasets (Supplementary Fig. 10A). Considering the sequencing strategies, (PCR-)cDNA Nanopore sequencing offers a more straightforward way to produce comprehensive data sets to analyse mRNA features. Almost 90% of known genes were covered by at least one read in all (PCR-)cDNA libraries. In contrast, direct RNA libraries only covered 70% (RNA001, replicate 3), 44% (RNA002, replicate 1) and 13% (RNA002, replicate 2) of the genes (Supplementary Fig. 10B). In order to evaluate how many reads are needed to cover at least 75% of all genes, we subsampled the reads of the representative rRNA-depleted PCR-cDNA sample (12 PCR cycles). We found that a

sequencing depth of about 10-fold is sufficient for this purpose, corresponding to 70,000 mRNA-mapping reads (Supplementary Fig. 10C).

Quantitatively, the data are highly reproducible between replicates and when using different library types (Figure 2A,C, Supplementary Table 3). Despite different sequencing platforms, protocols for sample preparation and sequencing depths, we observed a reasonably high correlation between expression data from published short-read Illumina RNA-seq data and ONT datasets (Figure 2B,C). Nevertheless, we found that a higher number of PCR cycles resulted in particularly GC-rich genes being underrepresented, leading to an overall more insufficient correlation in the PCR-cDNA datasets (Supplementary Fig. 11A,B). However, since we observed a similar effect with the non-amplified direct cDNA sample, which overall showed the best correlation to the Illumina data, other biases cannot be ruled out completely. For example, the SMRT-Cap protocol includes stringent size-selection filtering for fragments bigger than 1 kb. Consequently, from a purely quantitative perspective, the SMRT-Cap data are not fully comparable to the Nanopore data, but this may also be partly due to the sequencing depth. Taken together, the ONT data are highly consistent and allow a quantitative analysis of various transcriptomic features, which we will discuss in more detail below. However, the PCR bias is a critical point and researchers should carefully determine the number of PCR cycles required for their sample of choice.



**Figure 2: ONT sequencing of RNA and cDNA molecules is suitable for quantitative measurements.** **A**, Correlation between counts measured in transcripts per million (TPM) of cDNA replicate 2 and DRS replicate 2. Each point represents one gene, colour-coded by the density at the plot position. Gene lengths are indicated by circle size. Pearson correlation is given at the top left. **B**, Correlation between TPM counts of cDNA replicate 2 and the publicly available short-read Illumina dataset (ILL). **C**, Correlation matrix between all ONT, ILL and SMRT-Cappable-seq (SMRT) samples (40). Pearson correlation coefficients calculated from pairwise-complete observations are depicted and colour- and square-size coded. Additionally, the number of available pairwise comparisons is shown by circle size and colour in the upper right half of the plot.

### Identification and trimming of full-length transcripts

To accurately quantify and identify the number of full-length reads, we used Pychopper (<https://github.com/nanoporetech/pychopper>), a tool developed by ONT. This tool allows the detection and trimming of full-length cDNA reads based on SSP and VNP primers. In addition, it orients the sequenced reads (Supplementary Fig. 12A). As already evident from the length distribution, 2D-like reads make up a significant portion of the direct cDNA samples, which is confirmed by the low percentage of the Pychopper-detected full-length transcripts of about 34% in contrast to over 80% of full-length reads in all PCR-cDNA samples using standard settings (Supplementary Fig. 12B). However, a DCS-specific setting in Pychopper, which can handle 2D-like reads better and allows to rescue many reads, almost doubled the number of full-length reads detected (59%). When comparing

the aligned read lengths, we detected only a minimal difference between untrimmed and full-length-filtered reads, which is probably caused by random mapping of adapter sequences (Supplementary Fig. 12C). Despite Pychopper trimming, we observed that Adenine, caused by polyadenylation, and Guanine, caused by non-templated addition of nucleotides by the template-switching RT, are overrepresented at the 3' and 5' ends of cDNA reads, respectively (Supplementary Fig. 13A). To enable precise determination of transcript boundaries, we successfully trimmed off long poly(A) tails that are not expected to be found at the 3' ends of bacterial transcripts, and remaining SSP adapters from the 5' ends (Supplementary Fig. 13A,B).



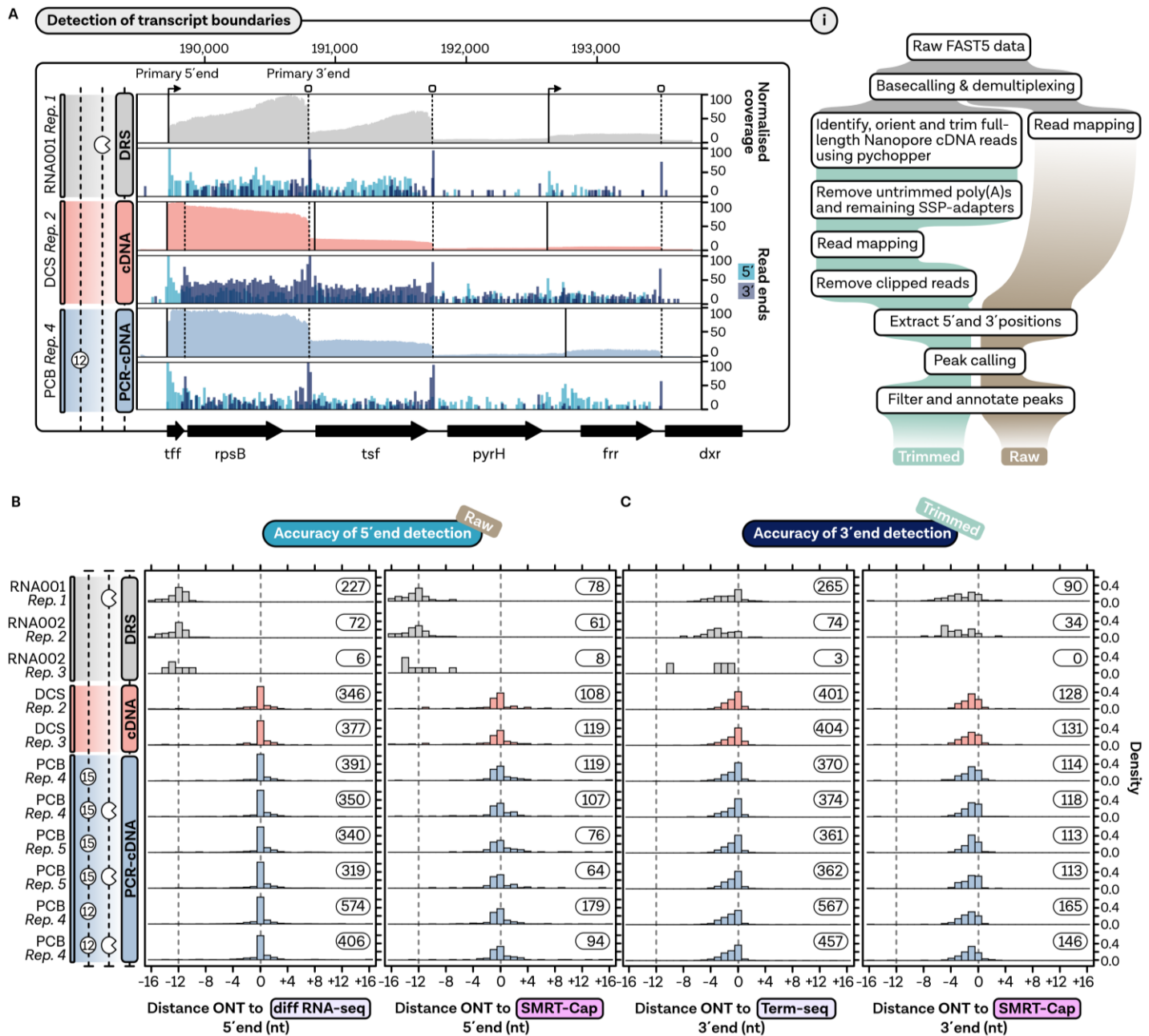
### *Identification of transcript boundaries*

Long-read ONT sequencing of RNA and cDNA molecules allows the simultaneous readout of 5' and 3' transcript boundaries (Figure 3A). Since full-length read starts and ends are expected to be enriched at functional relevant terminal positions and not randomly distributed, we applied a peak calling algorithm that detects all local maxima in reads first. In the next step, comparable to the evaluation of SMRT-Cap or short-read datasets, we defined the highest accumulation of 5' ends in a peak 300 bp upstream of an annotated gene as the primary 5'-transcript boundary (Figure 3A). Each additional peak in this region was designated as secondary and enriched intergenic peaks as internal. Because our samples do not only contain primary transcripts, we deliberately did not designate these ends as transcription start sites, although a considerable overlap is expected. We were able to define between 549 and 5,019 5' transcript ends in representative data sets, which varied depending on sequencing depth and trimming (Supplementary Fig. 14A, Supplementary Table 4). As described in other studies, the majority of enriched 5' ends are localized in internal regions. However, we could also identify up to 1,248 primary sites. Unexpectedly, untrimmed reads had a higher agreement in 5' ends at the single-nucleotide level to other comparable methods such as short-read differential RNA-seq and SMRT-Cap than trimmed reads (Figure 3B, Supplementary Fig. 14B) (40, 46). Ends determined by direct RNA sequencing are about 12 nt shorter, which is in line with previous observations (13, 19, 20) and can be rationalised by a lack of control of the RNA translocation speed after the motor protein falls off the 5' ends of the RNA (Figure 3B). PCR-cDNA and cDNA 5' ends are very clearly defined and predominantly end at the same base. In contrast, DRS leads to fuzzy 5' ends, presumably caused by a lower mapping accuracy. TEX treatment had neither a positive nor negative effect on 5' end detection or the number of reads starting at the enriched 5' ends. This may be due to the short treatment time and the digestion of the remaining ribosomal RNA leaving mRNA-mapping primary transcripts unaffected. Primary 5' transcript ends highly correlated between all different library types provided that enough reads support the enriched position (Supplementary Fig. 15). The moderate correlation (0.67) to SMRT-Cap 5' ends can mainly be attributed to the different library preparation approaches (Supplementary Fig. 16A): In contrast to our data, only primary transcripts are specifically captured and subsequently small transcripts naturally occurring in *E. coli* are intentionally lost due to size selection. The correlation drastically improved, when considering the positions of secondary 5' ends determined during ONT read analysis (Supplementary Fig. 16B). We found that for some genes, the SMRT-Cap primary site coincides with the ONT secondary, but not primary site. Although no specific enrichment for primary transcripts was performed for most of the samples, the 5' UTR distributions and the bacterial-typical nucleotide contents of upstream regions lead to the assumption that ONT sequencing is capable of accurately determining transcription start sites (Supplementary Fig. 17).

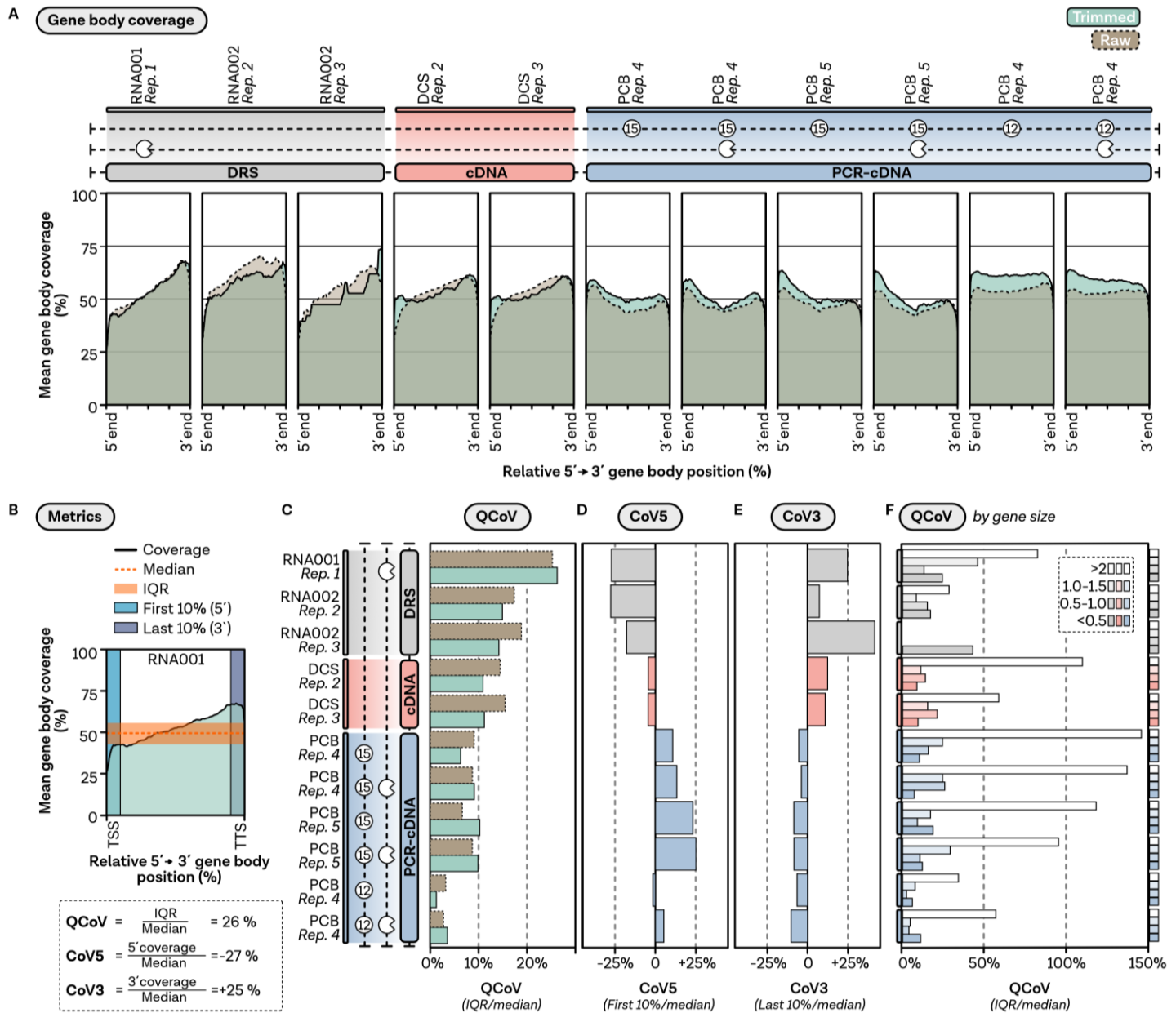
Peak enrichment analysis and 3' end annotations were performed as described for the 5' ends (Figure 3A). Overall, the number of enriched 3' ends found in the respective categories was slightly lower as compared to the 5' ends (Supplementary Fig. 18A, Supplementary Table 5). In contrast to the rather detrimental effect of trimming on the accuracy of 5' end detection, trimming increased the number of 3' ends that are identical to Term-seq (Supplementary Fig. 18B) (47). Although 3' end detection is highly reproducible and 3' ends overall highly correlate with SMRT-Cap detected ends, ONT 3' ends are fuzzier and tend to be up to 3 nucleotides shorter (Figure 3C, Supplementary Fig. 19). Since we cannot exclude that 3' to 5' exonucleases degrade RNAs after transcription, enriched sites may either represent genuine termination sites or enriched processed 3' ends. Nevertheless, the 3' UTR lengths of primary 3' ends and the poly(T) termination motif, which is typical for intrinsic terminators, suggest that most detected primary 3' ends are genuine transcription termination sites (Supplementary Fig. 20A,B).

### *Gene body coverage of long-read Nanopore reads*

In contrast to DRS, the cDNA protocols provide access to full-length transcripts due to the template-switching behaviour of the RT (Supplementary Fig. 12). Accordingly, it is expected that the 5' and 3' ends are covered to the same extent and that the coverage distribution over a gene is overall flat, which should improve an accurate transcriptional unit analysis. However, previous studies have shown that both DRS and direct cDNA reads are often truncated at the 5'-end. The reasons for this observation are still not completely clear but could be related to the fact that RNAs are directly sequenced starting from the 3' ends, to problems during template-switching, or sequencing-related issues like current spikes. To estimate the effect of the 3'-coverage bias, we looked at the gene body coverage profile between all samples and used previously introduced metrics, like the quartile coefficient of variation (QCoV) (19), to quantify coverage drops along the transcripts (Figure 4A,B). For the DRS and cDNA samples, we can confirm that 5' ends are less covered compared to the 3' ends (Figure 4A-E). Overamplification during PCR results in both ends being more enriched compared to the transcript centre. In contrast, at 12 cycles the reads are equally distributed across the gene body deviating on average less than 5% from the median coverage (Figure 4C). As expected, quality filtering and selection of full-length cDNA reads with both recognition adaptors results in an enrichment of the 5' ends for all cDNA samples (Figure 4D). However, we see that transcripts longer than 2 kb are less well uniformly covered (Figure 4E). It should be noted that these do not occur very often in our selected data that rely on the previous annotation of 5' and 3' ends, which could influence the distribution.



**Figure 3: Analysis of transcript boundaries detected using ONT RNA-seq methods.** **A**, Exemplary region in *E. coli* containing the non-coding RNA *tff* and 5 other genes. Coverage profiles of raw reads for the different library protocols have been normalised to 100. 5' read ends of raw reads (blue) and 3' read ends of trimmed reads (purple) are shown as histograms binned in a window of 10 bases, log-transformed and normalised to 100. Following the analysis pipeline depicted on the right we identified 5' and 3' enriched positions. Primary 5' (solid lines) and 3' ends (dashed lines) are shown for each dataset in the coverage plot. **B**, Accuracy of 5' end detection using raw reads assessed by comparison of distances between primary ONT 5' ends to differential RNA-seq primary transcription start sites (TSS) (46) and SMRT-Cap TSS (40). **C**, Accuracy of 3' end detection using trimmed reads assessed by comparison of distances between primary ONT 3' ends to short-read Term-seq primary transcription termination sites (TTS) (47) and SMRT-Cap TTS (40).



**Figure 4: Gene body coverage of raw and full-length enriched Nanopore reads.** **A**, Meta-analysis of gene body coverage profiles for genes that have an ONT-annotated 5' and 3' end. The gene bodies were scaled between the TSS and TTS to adjust for different transcript lengths. Coverage profiles show the mean values for each position after adjusting the coverage to values between 0 and 100 for each individual gene. Coverages are shown as area plots for calculated coverages based on raw (brown, dotted line) and trimmed (mint, solid line) reads. **B**, The decay from 5' or 3' ends and the overall coverage profiles were evaluated based on the quartile coefficient of variation (interquartile range/median, QCov) and by comparing the mean values of the first (CoV5) and last 10% (CoV3) of the coverage profiles to the median. Analysis of **C**, QCov calculated from trimmed (mint) and raw (brown) coverage profiles, **D**, CoV5 (trimmed), **E**, CoV3 (trimmed), and **F**, QCov (trimmed) values grouped by gene lengths that are indicated by transparency of the respective library type colour.

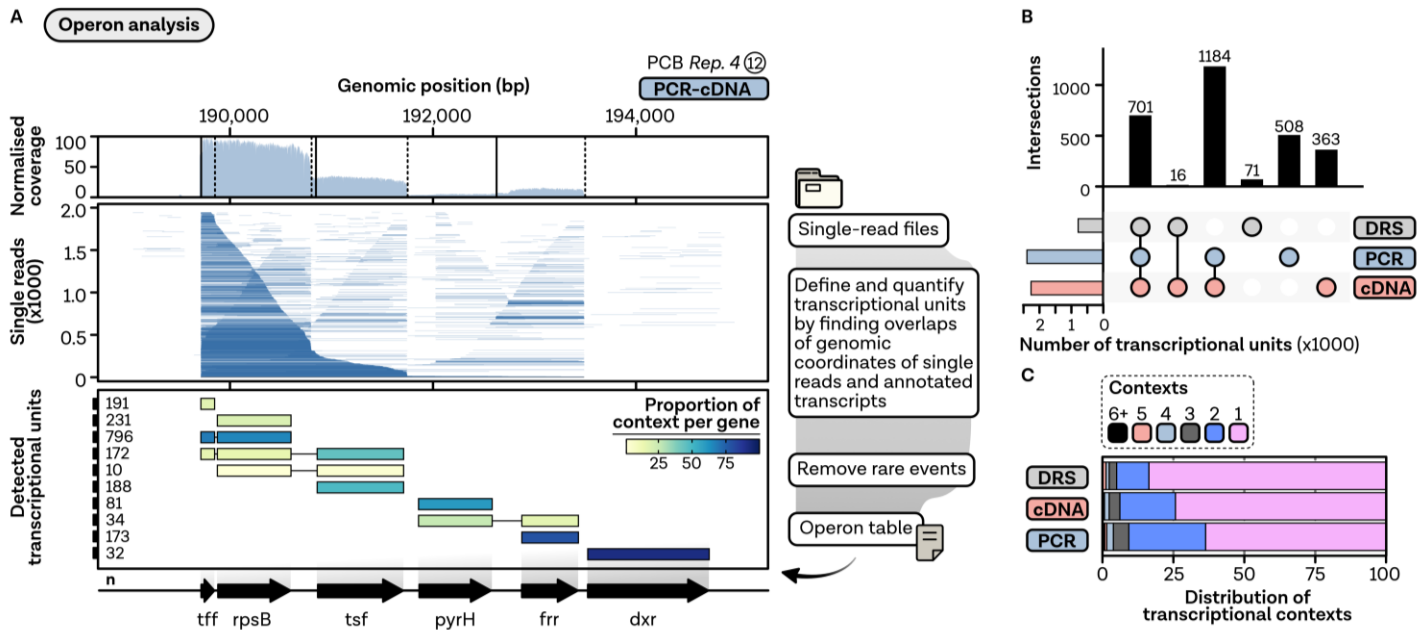
### *Nanopore sequencing captures the complexity of bacterial transcriptional units*

The distribution of reads over the gene body confirmed that ONT sequencing can cover both ends of a transcript. Since read lengths are theoretically only limited by the transcript size, ONT sequencing has the potential to accurately define complex transcriptional unit structures by finding overlaps between the mapping coordinates of individual reads and the transcript positions (Figure 5A). Following the annotation approach from the SMRT-Cap protocol, the unique combination of genes within a transcriptional unit was defined as the transcriptional context of a gene. Transcriptional unit prediction was performed exemplarily for one each of the DRS (RNA001 replicate 1), cDNA (DCS109 replicate 2) and PCR-cDNA (PCB109 replicate 4) libraries (Supplementary Table 6). Thereby, 788 (DRS), 2264 (cDNA) and 2433 (PCR-cDNA) unique transcriptional units were defined, respectively (Figure 5B). Mainly limited by the sequencing depth, the vast majority of defined transcriptional units (PCR-cDNA: 90%, cDNA: 83%, RNA: 90%) overlapped between the different protocols (Figure 5B). Hence, rare transcriptional unit variants stretching over multiple genes are not detected at low sequencing depth and stringent detection filters, which is also reflected in the mean number of genes encoded in a transcriptional unit: 1.14 for the DRS, 1.18 for the cDNA and 1.26 for the PCR-cDNA approach, respectively. This is in agreement with the observation that particularly long transcriptional units are underrepresented in our dataset. Therefore, the overall agreement with SMRT-Cap (43%) and the RegulonDB database (50%) is only moderate, which is presumably additionally heavily influenced by the respective detection algorithms, library preparation and sample conditions (Supplementary Fig. 21). Nevertheless, the distribution of transcriptional contexts in the PCR-cDNA dataset is in good agreement with the results from the SMRT-Cap analysis, showing that many genes are transcribed in more than one context (Figure 5C).

Note that without prior enrichment or treatment, quantification of the individual transcriptional contexts should consider that prokaryotic transcripts are subject to various degradation and

processing events: Therefore, it was not surprising that we captured a mix of 3' or 5' intact transcripts, which are often processed from the other end, as indicated by the ONT single-read tracks (Figure 5A, Supplementary Fig. 22,23). Effects that arise from RNA processing could be analysed in more detail when sequencing transcriptomes of exonuclease knock-out strains or with protocols that specifically enrich for primary transcripts (compare Send-seq and SMRT-Cap protocol) (40, 53). However, after the explicit enrichment of full-length transcripts and under the valid assumption that transcripts are not strongly degraded (compare RIN values) the extensive transcriptional heterogeneity is surprising. This can not only be seen in Figure 5A, but also in other examples, such as the RegulonDB-annotated operon *rpsP-rimM-trmD-rplS* (Supplementary Fig. 22) or a section of the genome containing many ribosomal proteins (Supplementary Fig. 23). The annotation of transcriptional units fits very well with the prediction of primary 5' and 3' end and illustrates that long-read ONT RNA-seq can more easily identify transcripts that arise from a shared promoter and have heterogeneous 3' ends. As already shown in the SMRT-Cap data, the *tff-rpsB-tsf* unit, which is identical to the operon annotated in the RegulonDB, is terminated in a stepwise manner. However, an additional termination site can be detected directly after the putative small RNA *tff*, which is otherwise lost through size selection (Figure 5A).

In summary, nanopore sequencing is capable of not only accurately detecting complex transcriptional unit structures but can also aid in quantification or in deciphering the unprecedented transcriptional heterogeneity, which may be improved by using specialized strains or conditions depending on the scientific question.



**Figure 5: Capability of ONT sequencing to capture complex bacterial transcriptional units.** **A**, Workflow, and visualisation of transcriptional unit analysis using long-read Nanopore data. After finding overlaps between single reads (compare single-read track sorted with increasing length of the read) and gene positions, transcriptional units are defined by the unique combination of genes that are covered by a read. The total number of reads assigned to a transcriptional unit is depicted on the left. The distribution of contexts per gene was calculated from the number of reads assigned to a gene in the respective context divided by the total number of reads per gene and is visualised using a colour code. **B**, UpSet-plot showing that the comparability and number of identically detected transcriptional units in the different library preparation methods is sequencing-depth dependent. **C**, Distribution of transcriptional contexts, which is defined as the number of transcriptional units a gene is part of.

## Discussion

In this study, we benchmarked all currently available kits from Oxford Nanopore for the analysis of RNAs, including direct sequencing of native RNA (RNA001, RNA002), direct cDNA (DCS109) and PCR-cDNA sequencing (PCB109, in the bacterial model organism *Escherichia coli* K-12. As a result, we demonstrate that multiple properties of the transcriptome can be examined simultaneously with high accuracy. This study therefore provides the first comprehensive analysis of ONT RNA-seq methods in prokaryotes. Furthermore, after extensive quality control of the sequenced libraries, we show that gene expression values are highly reproducible between library types and replicates and strongly correlate with the most commonly used short-read Illumina RNA-seq data. Additionally, we provide a bioinformatics workflow that allows accurate determination of transcript boundaries and quantitative analysis of transcriptional units applicable to all prokaryotes.

However, at present, some disadvantages of Nanopore RNA-seq should be considered: First, it must be ensured that the polyadenylation reaction in the organism of choice works equally effectively for all RNAs. Second, direct sequencing of RNAs requires a large amount of starting RNA material (> 10 µg) to yield enough mRNA (500 ng) left after effective rRNA depletion. Since the depletion kits are usually not designed for these quantities, the additional reactions are another cost factor.

Higher costs for DRS also originate from the slower sequencing speed, which negatively impacts throughput and the current lack of a barcoding options provided by ONT. Although there are already excellent options to build a custom set of DRS barcodes, this is not as straightforward to use as for cDNA libraries (54). Regarding 5' end detection, it has been shown multiple times that about 12 bases are missing from the DRS 5' ends. This observation can be explained by the motor protein falling off the end of a transcript resulting in a loss of control to guide the RNA through the nanopore, which is not the case for the cDNA data (13, 20). Another point of criticism that is repeatedly discussed is the comparatively low accuracy, especially for DRS, but also for cDNA datasets (13, 20, 55). Although this is not a significant problem for most questions, it affected the base-accurate trimming of adapter sequences and thus influenced the accuracy of the determination of the transcript ends. In particular, up to four more bases are trimmed off at the 3' ends since the homo-poly(A) sequence is usually low in quality and can only be trimmed inaccurately. Determining the 3' ends without trimming, which performs better at the 5' ends, performed even worse since long-read Nanopore mappers like minimap2 allow a higher number of errors (35). In general, the choice of the mapping tool should be well considered, as it greatly impacts the quality of the analysis. We applied the widely used and actively developed minimap2, which fails to align small RNAs (~80 bases cutoff) (35). While other mapping tools, like Magic-BLAST (56) or GraphMap2 (57) can align short transcripts, it is usually at the

expense of other aspects, and the method of choice dependent on the respective question. Despite or even because of these limitations, the Nanopore community is very active and interested in providing solutions for the problems discussed. Indeed, there are already promising applications that will also further improve ONT RNA-seq in prokaryotes in the future, like the error-correction of cDNA reads using isONcorrect (25) or the improvement of 5' end detection in DRS after 5' dependent adapter ligation (19).

Based on our results and considering the most cost-effective way to create and sequence libraries, we conclude that cDNA sequencing is the method of choice for most scientific questions, except for the analysis of RNA modifications (15). As only 1 ng of rRNA-depleted RNA is sufficient to generate PCR-cDNA libraries, cDNA-seq is highly preferable for organisms or conditions where the amount of RNA isolated is a crucial criterion. Our data clearly show that the number of cycles in the PCR should be controlled with special care. Otherwise, small AT-rich transcripts are preferentially amplified and sequenced, which distorts the quantification and further analyses. However, if this is handled correctly and the number of cycles is as low as possible, in our case 12, the PCR-cDNA data are highly comparable to the direct cDNA results. In any case, reverse transcription is a critical point for all cDNA. Nevertheless, the ONT-recommended Maxima H Minus Reverse Transcriptase (Thermo Fisher Scientific) performed quite well for our samples as documented in the gene body coverage data and the reproducibly good quantification. Another advantage of the enzyme used is that the reaction temperature can be increased to transcribe sequences with exceptionally high GC content or secondary structures.

In fact, there are already some sophisticated ways to profile full-length transcripts in *E. coli*, including the SMRT-Cappable-seq (40) and the SEnd-Seq (53) protocols. Comparison of SMRT-Cap and ONT data show that both datasets are highly congruent, although the repeatedly discussed size selection in the PacBio libraries plays a critical role and is a disadvantage. Unfortunately, despite the introduction of these methods, they have not yet been used in the prokaryotic community for further studies, although the reasons for this may well be diverse. However, we can imagine that the low initial costs of purchasing a MinION and the excellent performance could encourage some laboratories to use Nanopore RNA-seq in prokaryotes. The additional costs and IT infrastructure requirements are also limited, with basecalling of the data representing the highest computational effort for these analyses.

Taken together, a key advantage of ONT RNA-seq is that multiple features can be addressed simultaneously with high accuracy. This versatility distinguishes the technique from the various RNA-seq technologies designed to tackle only one specific question or biochemical assays. Strikingly, Nanopore sequencing is a *bona fide* single-molecule method. Hence,

molecular heterogeneity at the transcriptome level can be analysed, so that even minor RNA populations can be detected that are inevitably lost in ensemble sequencing approaches. However, we observed a complex transcription pattern with multiple possible RNA variants. Given that transcription and translation are coupled in *E. coli*, new questions about the translation efficiency and transcript stability of the transcript variants emerge (58–61). Furthermore, high-quality long-read RNA-seq data can be used to analyse degradation or processing patterns to gain new insights into mRNA decay in prokaryotes. With this study, we not only show the applicability of ONT RNA-seq in prokaryotes, but also provide representative long-read transcriptome data from *E. coli* and a robust bioinformatical workflow to the community that can be used to tackle various questions.

## Data availability:

To facilitate easier access basecalled and demultiplexed FASTQ, mapped BAM files from untrimmed reads and large read summary files are publicly available from <https://zenodo.org/record/4879174#.YLSkky221pQ>.

All scripts and code used in this work are available on GitHub (<https://github.com/felixgrunberger/microbepore>). Additionally, a more detailed documentation can be found at <https://felixgrunberger.github.io/microbepore>.

Sequencing files in original FAST5 format are publicly available in the Sequence Read Archive SRA (RNA001: PRJNA632538, all other datasets: PRJNA731531).

## Funding:

This work was supported by the Deutsche Forschungsgemeinschaft [SFB960 TPA7 to D.G., and SFB960 TPB13 to S.F.-C.]; Funding for open access charge: Deutsche Forschungsgemeinschaft.

## Conflict of Interest:

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

## Acknowledgements:

We thank all the members of the Ferreira-Cerca lab and of the Grohmann lab, especially Prof. Dr. Winfried Hausner and Martin Fenk for fruitful discussions.

## References

1. Levy, S.E. and Myers, R.M. (2016) Advancements in Next-Generation Sequencing. *Annu. Rev. Genomics Hum. Genet.*, **17**, 95–115.
2. Escobar-Zepeda, A., Vera-Ponce de León, A. and Sanchez-Flores, A. (2015) The Road to Metagenomics: From Microbiology to DNA Sequencing Technologies and Bioinformatics. *Front. Genet.*, **6**.
3. Wang, Z., Gerstein, M. and Snyder, M. (2009) RNA-Seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.*, **10**, 57–63.
4. Hör, J., Gorski, S.A. and Vogel, J. (2018) Bacterial RNA Biology on a Genome Scale. *Mol. Cell*, **70**, 785–799.
5. Croucher, N.J. and Thomson, N.R. (2010) Studying bacterial transcriptomes using RNA-seq. *Curr. Opin. Microbiol.*, **13**, 619–624.
6. Nowrousian, M. (2010) Next-Generation Sequencing Techniques for Eukaryotic Microorganisms: Sequencing-Based Solutions to Biological Problems. *Eukaryot. Cell*, **9**, 1300–1310.
7. Saliba, A.-E., Santos, S. and Vogel, J. (2017) New RNA-seq approaches for the study of bacterial pathogens. *Curr. Opin. Microbiol.*, **35**, 78–87.
8. Stark, R., Grzelak, M. and Hadfield, J. (2019) RNA sequencing: the teenage years. *Nat. Rev. Genet.*, **10.1038/s41576-019-0150-2**.
9. Byrne, A., Cole, C., Volden, R. and Vollmers, C. (2019) Realizing the potential of full-length transcriptome sequencing. *Philos. Trans. R. Soc. B Biol. Sci.*, **374**, 20190097.
10. Tilgner, H., Jahanbani, F., Blauwkamp, T., Moshrefi, A., Jaeger, E., Chen, F., Harel, I., Bustamante, C.D., Rasmussen, M. and Snyder, M.P. (2015) Comprehensive transcriptome analysis using synthetic long-read sequencing reveals molecular co-association of distant splicing events. *Nat. Biotechnol.*, **33**, 736–742.
11. Mikheyev, A.S. and Tin, M.M.Y. (2014) A first look at the Oxford Nanopore MinION sequencer. *Mol. Ecol. Resour.*, **14**, 1097–1102.
12. Eid, J., Fehr, A., Gray, J., Luong, K., Lyle, J., Otto, G., Peluso, P., Rank, D., Baybayan, P., Bettman, B., et al. (2009) Real-Time DNA Sequencing from Single Polymerase Molecules. *Science (80-. )*, **323**, 133–138.
13. Soneson, C., Yao, Y., Bratus-neuenschwander, A., Patrignani, A., Robinson, M.D. and Hussain, S. (2019) A comprehensive examination of Nanopore native RNA sequencing for characterization of complex transcriptomes. *Nat. Commun.*, **10**, 1–14.
14. Vilfan, I.D., Tsai, Y.-C., Clark, T.A., Wegener, J., Dai, Q., Yi, C., Pan, T., Turner, S.W. and Korlach, J. (2013) Analysis of RNA base modification and structural rearrangement by single-molecule real-time detection of reverse transcription. *J. Nanobiotechnology*, **11**, 8.
15. Begik, O., Lucas, M.C., Pyszcz, L.P., Ramirez, J.M., Medina, R., Milenkovic, I., Cruciani, S., Liu, H., Vieira, H.G.S., Sas-Chen, A., et al. (2021) Quantitative profiling of pseudouridylation dynamics in native RNAs with nanopore sequencing. *Nat. Biotechnol.*, **10.1038/s41587-021-00915-6**.
16. Liu, H., Begik, O., Lucas, M.C., Ramirez, J.M., Mason, C.E., Wiener, D., Schwartz, S., Mattick, J.S., Smith, M.A., Novoa, E.M., et al. (2019) Accurate detection of m6A RNA modifications in native RNA sequences. *Nat. Commun.*, **10.1038/s41467-019-11713-9**.
17. Smith, A.M., Jain, M., Mulrone, L., Garalde, D.R. and Akeson, M. (2019) Reading canonical and modified nucleobases in 16S ribosomal RNA using nanopore native RNA sequencing. *PLoS One*, **14**, e0216709.
18. Jenjaroenpun, P., Wongsurawat, T., Wadley, T.D., Wassenaar, T.M., Liu, J., Dai, Q., Wanchai, V., Akel, N.S., Jamshidi-Parsian, A., Franco, A.T., et al. (2021) Decoding the epitranscriptional landscape from native RNA sequences. *Nucleic Acids Res.*, **49**, e7.
19. Parker, M.T., Knop, K., Sherwood, A. V., Schurch, N.J., Mackinnon, K., Gould, P.D., Hall, A.J., Barton, G.J. and Simpson, G.G. (2020) Nanopore direct RNA sequencing maps the complexity of Arabidopsis mRNA processing and m6A modification. *Elife*, **9**.
20. Workman, R.E., Tang, A.D., Tang, P.S., Jain, M., Tyson, J.R., Razaghi, R., Zuzarte, P.C., Gilpatrick, T., Payne, A., Quick, J., et al. (2019) Nanopore native RNA sequencing of a human poly(A) transcriptome. *Nat. Methods*, **16**, 1297–1305.
21. Boldogkői, Z., Moldován, N., Balázs, Z., Snyder, M. and Tombácz, D. (2019) Long-Read Sequencing – A Powerful Tool in Viral Transcriptome Research. *Trends Microbiol.*, **27**, 578–592.
22. Keller, M.W., Rambo-Martin, B.L., Wilson, M.M., Ridenour, C.A., Shepard, S.S., Stark, T.J., Neuhaus, E.B., Dugan, V.G., Wentworth, D.E. and Barnes, J.R. (2018) Direct RNA Sequencing of the Coding Complete Influenza A Virus Genome. *Sci. Rep.*, **8**, 14408.
23. Viehweger, A., Krautwurst, S., Lamkiewicz, K., Madhugiri, R., Ziebuhr, J., Hölzer, M. and Marz, M. (2019) Direct RNA nanopore sequencing of full-length coronavirus genomes provides novel insights into structural variants and enables modification analysis. *Genome Res.*, **10.1101/483693**.
24. Wang, D., Jiang, A., Feng, J., Li, G., Guo, D., Sajid, M., Wu, K., Zhang, Q., Ponty, Y., Will, S., et al. (2021) The SARS-CoV-2 subgenome landscape and its novel regulatory features. *Mol. Cell*, **10.1016/j.molcel.2021.02.036**.
25. Sahlin, K., Sipos, B., James, P.L. and Medvedev, P. (2021) Error correction enables use of Oxford Nanopore technology for reference-free transcriptome analysis. *Nat. Commun.*, **12**, 1–13.
26. Zhao, L., Zhang, H., Kohnen, M. V., Prasad, K.V.S.K., Gu, L. and Reddy, A.S.N. (2019) Analysis of Transcriptome and Epitranscriptome in Plants Using PacBio Iso-Seq and Nanopore-Based Direct RNA Sequencing. *Front. Genet.*, **10**.
27. Tombácz, D., Moldován, N., Balázs, Z., Gulyás, G., Csabai, Z., Boldogkői, M., Snyder, M. and Boldogkői, Z. (2019) Multiple Long-Read Sequencing

- Survey of Herpes Simplex Virus Dynamic Transcriptome. *Front. Genet.*, **10**.
28. Byrne,A., Beaudin,A.E., Olsen,H.E., Jain,M., Cole,C., Palmer,T., DuBois,R.M., Forsberg,E.C., Akeson,M. and Vollmers,C. (2017) Nanopore long-read RNAseq reveals widespread transcriptional variation among the surface receptors of individual B cells. *Nat. Commun.*, **8**, 16027.
29. Dong,X., Tian,L., Gouil,Q., Kariyawasam,H., Su,S., De Paoli-Iseppi,R., Praver,Y.D.J., Clark,M.B., Breslin,K., Iminoff,M., *et al.* (2021) The long and the short of it: unlocking nanopore long-read RNA sequencing data with short-read differential expression analysis tools. *NAR Genomics Bioinforma.*, **3**, 1–11.
30. Seki,M., Oka,M., Xu,L., Suzuki,A. and Suzuki,Y. (2021) Transcript Identification Through Long-Read Sequencing. In.pp. 531–541.
31. Choi,S.C. (2016) On the study of microbial transcriptomes using second- and third-generation sequencing technologies. *J. Microbiol.*, **54**, 527–536.
32. Yang,M., Cousineau,A., Liu,X., Luo,Y., Sun,D., Li,S., Gu,T., Sun,L., Dillow,H., Lepine,J., *et al.* (2020) Direct Metatranscriptome RNA-seq and Multiplex RT-PCR Amplicon Sequencing on Nanopore MinION – Promising Strategies for Multiplex Identification of Viable Pathogens in Food. *Front. Microbiol.*, **11**, 1–14.
33. Pitt,M.E., Nguyen,S.H., Duarte,T.P.S., Teng,H., Blaskovich,M.A.T., Cooper,M.A. and Coin,L.J.M. (2020) Evaluating the genome and resistome of extensively drug-resistant *Klebsiella pneumoniae* using native DNA and RNA Nanopore sequencing. *Gigascience*, **9**, 1–14.
34. Riley,M., Abe,T., Arnaud,M.B., Berlyn,M.K.B., Blattner,F.R., Chaudhuri,R.R., Glasner,J.D., Horiuchi,T., Keseler,I.M., Kosuge,T., *et al.* (2006) *Escherichia coli* K-12: a cooperatively developed annotation snapshot--2005. *Nucleic Acids Res.*, **34**, 1–9.
35. Li,H. (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**, 3094–3100.
36. Li,H., Handsaker,B., Wysoker,A., Fennell,T., Ruan,J., Homer,N., Marth,G., Abecasis,G. and Durbin,R. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
37. Quinlan,A.R. and Hall,I.M. (2010) BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841–842.
38. Bolger,A.M., Lohse,M. and Usadel,B. (2014) Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics*, **10.1093/bioinformatics/btu170**.
39. Langmead,B. and Salzberg,S.L. (2012) Fast gapped-read alignment with Bowtie 2. *Nat. Methods*, **9**, 357–359.
40. Yan,B., Boitano,M., Clark,T.A. and Ettwiller,L. (2018) SMRT-Cappable-seq reveals complex operon variants in bacteria. *Nat. Commun.*, **9**, 3676.
41. Patro,R., Duggal,G., Love,M.I., Irizarry,R.A. and Kingsford,C. (2017) Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods*, **14**, 417–419.
42. Martin,M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet:journal*, **17**, 10.
43. Adams,P.P., Baniulyte,G., Esnault,C., Chegireddy,K., Singh,N., Monge,M., Dale,R.K., Storz,G. and Wade,J.T. (2021) Regulatory roles of *Escherichia coli* 5' utr and orf-internal rnas detected by 3' end mapping. *Elife*, **10**, 1–33.
44. Lawrence,M., Huber,W., Pagès,H., Aboyoun,P., Carlson,M., Gentleman,R., Morgan,M.T. and Carey,V.J. (2013) Software for Computing and Annotating Genomic Ranges. *PLoS Comput. Biol.*, **9**, 1–10.
45. Santos-Zavaleta,A., Salgado,H., Gama-Castro,S., Sánchez-Pérez,M., Gómez-Romero,L., Ledezma-Tejeida,D., García-Sotelo,J.S., Alquicira-Hernández,K., Muñoz-Rascado,L.J., Peña-Loredo,P., *et al.* (2019) RegulonDB v 10.5: Tackling challenges to unify classic and high throughput knowledge of gene regulation in *E. coli* K-12. *Nucleic Acids Res.*, **47**, D212–D220.
46. Thomason,M.K., Bischler,T., Eisenbart,S.K., Förstner,K.U., Zhang,A., Herbig,A., Nieselt,K., Sharma,C.M., Storz,G. and Storz,G. (2015) Global Transcriptional Start Site Mapping Using Differential RNA Sequencing Reveals Novel Antisense RNAs in *Escherichia coli*. *J. Bacteriol.*, **197**, 18–28.
47. Dar,D. and Sorek,R. (2018) High-resolution RNA 3-ends mapping of bacterial Rho-dependent transcripts. *Nucleic Acids Res.*, **46**, 6797–6805.
48. Matz,M., Shagin,D., Bogdanova,E., Britanova,O., Lukyanov,S., Diatchenko,L. and Chenchik,A. (1999) Amplification of cDNA ends based on template-switching effect and step-out PCR. *Nucleic Acids Res.*, **27**, 1558–1560.
49. Robinson,J.T., Thorvaldsdóttir,H., Winckler,W., Guttman,M., Lander,E.S., Getz,G. and Mesirov,J.P. (2011) Integrative genomics viewer. *Nat. Biotechnol.*, **29**, 24–26.
50. Himeno,H., Kurita,D. and Muto,A. (2014) TmRNA-mediated translation as the major ribosome rescue system in a bacterial cell. *Front. Genet.*, **5**, 1–13.
51. Tuiskunen,A., Leparç-Goffart,I., Boubis,L., Monteil,V., Klingström,J., Tolou,H.J., Lundkvist,A. and Plumet,S. (2010) Self-priming of reverse transcriptase impairs strand-specific detection of dengue virus RNA. *J. Gen. Virol.*, **91**, 1019–1027.
52. Perocchi,F., Xu,Z., Clauder-Münster,S. and Steinmetz,L.M. (2007) Antisense artifacts in transcriptome microarray experiments are resolved by actinomycin D. *Nucleic Acids Res.*, **35**.
53. Ju,X., Li,D. and Liu,S. (2019) Full-length RNA profiling reveals pervasive bidirectional transcription terminators in bacteria. *Nat. Microbiol.*, **4**, 1907–1918.
54. Smith,M.A., Ersavas,T., Ferguson,J.M., Liu,H., Lucas,M.C., Begik,O., Bojarski,L., Barton,K. and Novoa,E.M. (2020) Molecular barcoding of native RNAs using nanopore sequencing and deep learning. *Genome Res.*, **30**, 1345–1353.
55. Garalde,D.R., Snell,E.A., Jachimowicz,D., Sipo,B., Lloyd,J.H., Bruce,M., Pantic,N., Admassu,T., James,P., Warland,A., *et al.* (2018) Highly parallel direct RNA sequencing on an array of nanopores. *Nat. Methods*, **15**, 201–206.
56. Boratyn,G.M., Thierry-Mieg,J., Thierry-Mieg,D., Busby,B. and Madden,T.L. (2019) Magic-BLAST, an accurate RNA-seq aligner for long and short reads. *BMC Bioinformatics*, **20**, 405.
57. Sović,I., Šikić,M., Wilm,A., Fenlon,S.N., Chen,S. and Nagarajan,N. (2016) Fast and sensitive mapping of nanopore sequencing reads with GraphMap. *Nat. Commun.*, **7**.
58. Webster,M.W., Takacs,M., Zhu,C., Vidmar,V., Eduljee,A., Abdelkareem,M. and Weixlbaumer,A. (2020) Structural basis of transcription-translation coupling and collision in bacteria. *Science (80-. )*, **10.1126/science.abb5036**.
59. Wang,C., Molodtsov,V., Firlar,E., Kaelber,J.T., Blaha,G., Su,M. and Ebright,R.H. (2020) Structural basis of transcription-translation coupling. *Science (80-. )*, **10.1126/science.abb5317**.
60. Proshkin,S., Rachid Rahmouni,A., Mironov,A., Nudler,E., Rahmouni,A.R., Mironov,A. and Nudler,E. (2010) Cooperation Between Translating Ribosomes and RNA Polymerase in Transcription Elongation. *Science (80-. )*, **328**, 504–508.
61. Irastortza-Olaziregi,M. and Amster-Choder,O. (2021) Coupled Transcription-Translation in Prokaryotes: An Old Couple With New Surprises. *Front. Microbiol.*, **11**.