

1-19-2023

Semi-quantitative Detection of Pseudouridine Modifications and Type I/II I/li Hypermodifications in Human mRNAs Using Direct Long-Read Sequencing

Sepideh Tavakoli

Mohammad Nabizadeh


Amr Makhamreh

Howard Gamper

Caroline A McCormick

See next page for additional authors

Follow this and additional works at: <https://jdc.jefferson.edu/bmpfp>

 Part of the [Medical Biochemistry Commons](#), and the [Medical Molecular Biology Commons](#)

[Let us know how access to this document benefits you](#)

This Article is brought to you for free and open access by the Jefferson Digital Commons. The Jefferson Digital Commons is a service of Thomas Jefferson University's [Center for Teaching and Learning \(CTL\)](#). The Commons is a showcase for Jefferson books and journals, peer-reviewed scholarly publications, unique historical collections from the University archives, and teaching tools. The Jefferson Digital Commons allows researchers and interested readers anywhere in the world to learn about and keep up to date with Jefferson scholarship. This article has been accepted for inclusion in Department of Biochemistry and Molecular Biology Faculty Papers by an authorized administrator of the Jefferson Digital Commons. For more information, please contact: JeffersonDigitalCommons@jefferson.edu.

Authors


Sepideh Tavakoli, Mohammad Nabizadeh, Amr Makhamreh, Howard Gamper, Caroline A McCormick, Neda K Rezapour, Ya-Ming Hou, Meni Wanunu, and Sara H Rouhanifard

Semi-quantitative detection of pseudouridine modifications and type I/II hypermodifications in human mRNAs using direct long-read sequencing

Received: 27 June 2022

Accepted: 5 January 2023

Published online: 19 January 2023

 Check for updates

Sepideh Tavakoli^{1,5}, Mohammad Nabizadeh^{2,5}, Amr Makhamreh¹, Howard Gamper³, Caroline A. McCormick¹, Neda K. Rezapour⁴, Ya-Ming Hou³, Meni Wanunu^{1,4} & Sara H. Rouhanifard¹ ✉

Here, we develop and apply a semi-quantitative method for the high-confidence identification of pseudouridylated sites on mammalian mRNAs via direct long-read nanopore sequencing. A comparative analysis of a modification-free transcriptome reveals that the depth of coverage and specific k-mer sequences are critical parameters for accurate basecalling. By adjusting these parameters for high-confidence U-to-C basecalling errors, we identify many known sites of pseudouridylation and uncover previously unreported uridine-modified sites, many of which fall in k-mers that are known targets of pseudouridine synthases. Identified sites are validated using 1000-mer synthetic RNA controls bearing a single pseudouridine in the center position, demonstrating systematic under-calling using our approach. We identify mRNAs with up to 7 unique modification sites. Our workflow allows direct detection of low-, medium-, and high-occupancy pseudouridine modifications on native RNA molecules from nanopore sequencing data and multiple modifications on the same strand.

Enzyme-mediated RNA chemical modifications have been extensively studied on noncoding RNAs^{1,2}; however, messenger RNAs are also targets of RNA modification. Although modifications are less frequent in mRNAs than other RNAs³, mRNA modifications can potentially impact gene expression⁴, RNA tertiary structure formation⁵, or the recruitment of RNA-binding proteins⁶. Pseudouridine (ψ) is synthesized from uridine by one of more than a dozen pseudouridine synthases identified to date⁷. It was the first discovered RNA modification⁸ and represents 0.2–0.6% of total uridines in mammalian mRNAs³. Other uridine modifications on mRNAs include 5-methyl uridine⁹ and dihydrouridine^{10,11}, but these occur to a lesser extent on mRNAs. ψ -modified mRNAs are more resistant to RNase-mediated degradation¹²,

and have the potential to modulate splicing¹³, immunogenicity¹⁴, and translation^{15,16} in vivo. Further, ψ modifications of RNAs are responsive to cellular stress, leading to increased RNA half-life^{17,18}. A critical barrier to understanding the precise biological functions of pseudouridylation is the absence of high-confidence methods to map ψ -sites in mRNAs. ψ modifications do not affect Watson-Crick base pairing¹⁹, thereby making them indistinguishable from uridine in hybridization-based methods. In addition, this modification bears the same molecular weight as the canonical uridine, making digested nucleotides challenging to detect directly by mass spectrometry^{20,21}, and even more difficult to detect within a specific mRNA sequence that is isolated from cells due to the high input and purity requirements.

¹Department of Bioengineering, Northeastern University, Boston, MA, USA. ²Department of Mechanical Engineering, Northeastern University, Boston, MA, USA. ³Department of Biochemistry and Molecular Biology, Thomas Jefferson University, Philadelphia, PA, USA. ⁴Department of Physics, Northeastern University, Boston, MA, USA. ⁵These authors contributed equally: Sepideh Tavakoli, Mohammad Nabizadeh. ✉e-mail: s.rouhanifard@northeastern.edu

Ψ is conventionally labeled using N-cyclohexyl-N'-b-(4-methylmorpholinium) ethylcarbodiimide (CMC), a reagent that modifies the N1 and N3 positions of ψ , N1 of guanine, and the N3 of uridine²². Treatment with a strong base removes the CMC from all the sites except for the N3 position of ψ . Recently, the use of an RNA bisulfite reaction was demonstrated for ψ -specific labeling^{23,24}. Chemical labeling of ψ combined with next-generation sequencing^{3,18,23} has yielded over 2000 putative ψ sites within mammalian mRNAs, but different methods have identified different sites with somewhat limited overlap²⁵, pointing to a need for alternative detection methods that do not require CMC. Additionally, since these methods are combined with short read sequencing, combinatorial analysis of multiple modifications on one transcript is impossible. Here we aim to develop a direct and orthogonal method for ψ detection on mRNAs without relying on intermediate chemical reactions, but with the ability to detect previously annotated sites and to uncover previously unreported sites.

Recently, several studies have reported using nanopore-based direct RNA sequencing to directly read RNA modifications^{26–31}. In these reports, ion current differences for different k-mer sequences (k = 5 nucleotides) as an RNA strand is moved through the pore suggest the presence of a modified RNA base. Detection of ψ using nanopores was also confirmed for ribosomal RNAs (rRNAs)²⁸, on the *Saccharomyces cerevisiae* transcriptome²⁹, and for viral RNAs³¹, as indicated by a U-to-C base-calling error at various sequence sites. Algorithms for ψ quantification have been produced^{29,30} using combinatorial sequences that contain many ψ sites within close proximity, and control RNAs that contain many natural RNA modifications also in close proximity (e.g., rRNA). While the control molecules in these studies allow analysis of many k-mers, the accuracy of quantifying ψ occupancy at a given site can be highly dependent on the nucleotide sequence surrounding the modification. Moreover, sequence context is particularly important for quantification of RNA molecules wherein the secondary structure can influence the kinetics of translocation through the nanopore³².

Here, we describe a nanopore-based method to identify known ψ modifications and map uridine modifications in a HeLa transcriptome. We compare the sequence alignment to identical negative controls without RNA modifications and show that the number of reads and specific k-mer sequences are critical parameters for defining ψ sites and for assigning significance values based on these parameters. Our approach recapitulates 198 previously annotated ψ sites, 34 of which are detected by 3 independent methods, thus providing a ground truth list of ψ modifications in HeLa cells. Our approach also reveals 1691 putative sites of uridine-modification that have not been reported previously. We show that these previously unreported sites tend to occur within k-mer sequences that are often recognized by pseudouridine synthases, such as PUS7 and TRUB1.

Applying our algorithm for detecting ψ modifications using rRNAs which have been comprehensively annotated by mass spectrometry, we assign 38/43 ψ modifications. We demonstrate, however, that rRNA is not suitable for benchmarking RNA modifications by nanopore sequencing due to the frequent clustering of RNA modifications nearby to the ψ -site, which interferes with the accuracy of basecalling, thus leading to false positives. Additionally, we synthesize and analyze five 1000 mer RNA controls containing either uridine or ψ within the sequence context of a known pseudouridylated position in the human transcriptome. This analysis reveals that U-to-C mismatch errors are systematically under-called for the detection of ψ , enabling us to identify 40 high-occupancy ψ sites, which we denote as hypermodified type I.

Further, we identify 38 mRNAs with up to 7 uridine modification sites, which are confirmed by single-read analysis. Combined, we report a workflow that enables direct identification and semi-quantification of the ψ modification on native mRNA molecules. The long nanopore reads allow for the detection of multiple

modifications on one transcript, which can shed light on cooperative effects of mRNA modifications as a mechanism to modulate gene expression.

Results

Nanopore analysis of an unmodified HeLa transcriptome generated by in vitro transcription

To identify putative sites of mRNA ψ modifications, we extracted RNA from HeLa cells and prepared two libraries for direct RNA sequencing (Fig. 1a). The first (direct) library consists of native mRNAs, and the second is an in vitro transcribed (IVT) mRNA control library in which polyadenylated RNA samples were reverse transcribed to cDNAs, then IVT back into RNA using canonical nucleotides to delete the RNA modifications. Each library was prepared for sequencing using Oxford Nanopore's Direct RNA Sequencing Kit, and then sequenced on a separate MinION flowcell and basecalled using *Guppy 3.2.10*. Three replicates of the native RNA library produced an average of ~1.2 million poly(A) RNA strand reads, of which ~800,000 have a read quality of 7, with an average of N50 read length (defined as the shortest read length needed to cover 50% of the sequenced nucleotides) of 850 bases and a median length of ~670 bases (Supplementary Fig. 1). Further, we compared the coverage for individual genes in the direct RNA libraries and found that transcripts per million (TPMs) were very similar ($R^2 = 1.06$ for replicates 1 and 2 and 0.97 for replicates 2 and 3; Supplementary Fig. 1) Similarly, two replicates of the IVT library produced an average of 1.6 million reads that passed the quality filter, with N50 of 890 and a median length of 710 bases. Alignment was performed using *minimap2*^{17,33} and the reads for each library were subsequently aligned to the GRCh38 human genome reference.

Basecalling accuracy is used to identify uridine modifications in RNA

To define differences between the IVT and direct libraries for uridine modification detection, any source of error other than the uridine modification itself must be minimized, including misalignments to the GRCh38 human genome reference. We minimized incorrect alignment by only considering the primary alignment of each read (i.e., the alignment with the highest mapping quality). Also, the reads with mapping quality score of 20 or higher which corresponds to the probability of correct alignment of 99% or higher were used for the downstream analysis. The second potential source of error is the presence of single-nucleotide polymorphisms (SNPs), whereby the base is different from the reference genome. We identified likely SNP sites based on an equivalent U-to-C mismatch percentage in both the IVT and the direct RNA sequencing samples (Supplementary Fig. 2). In the case of a modified RNA nucleotide, the U-to-C mismatch percentage would be significantly higher in the direct RNA sequencing sample relative to the IVT control at the site of modification (Supplementary Fig. 2). The third source of error is a systematic basecalling error, whereby the basecalling algorithm fails to identify the correct base. To assess the basecalling accuracy using the *Guppy 3.2.10* algorithm, we calculated the error in the IVT control by comparing the basecalling to the reference genome (Fig. 1b). Since the IVT control contains only the canonical RNA nucleotides, these errors were independent of RNA modifications. We confirmed that the basecaller could reliably identify unmodified and aligned nucleotides with an average U-to-C error of 2.64%.

To confirm the quality of the IVT unmodified transcriptome, we compared the coverage for individual genes in the IVT and direct RNA libraries and found that TPMs were very similar ($R^2 = 0.96$; Fig. 1c). We also compared the distribution of read lengths for the IVT and direct RNA libraries and found that the samples were overlapping (Fig. 1d). Likewise, the coverage for individual transcripts was similar for IVT and direct RNA libraries (Fig. 1e), thus validating the IVT library as an unmodified transcriptome control.

Direct RNA nanopore sequencing identifies pseudouridines in mRNA via systematic U-to-C base-calling errors

We then examined specific locations on human mRNAs that have been previously identified as ψ sites by chemical-based methods (Fig. 1f). We selected 5 genes as examples: *ID11* (chr10:1044099)^{3,17,23}, *PARP4* (chr13:24426505)^{3,23}, *PSMB2* (chr1:35603333)^{3,17,23}, *MCM5* (chr22:35424407)^{3,17}, and *PABPC4* (chr1:39565149)^{3,17}, representing a range of different k-mers with a putative ψ in the center nucleotide (GUUCA, GUUCA, GUUCG, UGUAG, and GUUCC respectively). We chose a range of k-mers because specific sequences can influence the accuracy of base-calling (Supplementary Fig. 3). We detected a systematic U-to-C mismatch error at the reported ψ site in duplicates of each gene by direct RNA sequencing (*ID11* (chr10:1044099): $96.06 \pm 1.16\%$, *PARP4* (chr13:24426505): $91.71 \pm 7.56\%$, *PSMB2* (chr1:35603333): $81.07 \pm 1.68\%$,

MCM5 (chr22:35424407): $54.82 \pm 4.96\%$, *PABPC4* (chr1:39565149): $55.08 \pm 3.97\%$). We confirmed that the IVT samples maintained the standard base-calling error at each site (3.75%, 4.54%, 1.67%, 5.26% and 8.34% respectively; Fig. 1c).

Calculating the significance of U-to-C mismatch as a proxy for modification is dependent on mismatch percentage at a given site, the number of reads, and the surrounding nucleotides

To further improve the use of the U-to-C mismatch error as a proxy for U modifications we needed to minimize the error that occurs from other factors. We observed that the base quality on sites that have 3 or fewer reads is low, relative to the rest of the population, which would create bias in the downstream analysis (Fig. 2a). One reason for the lower quality of these sites is their proximity to the start/end of the

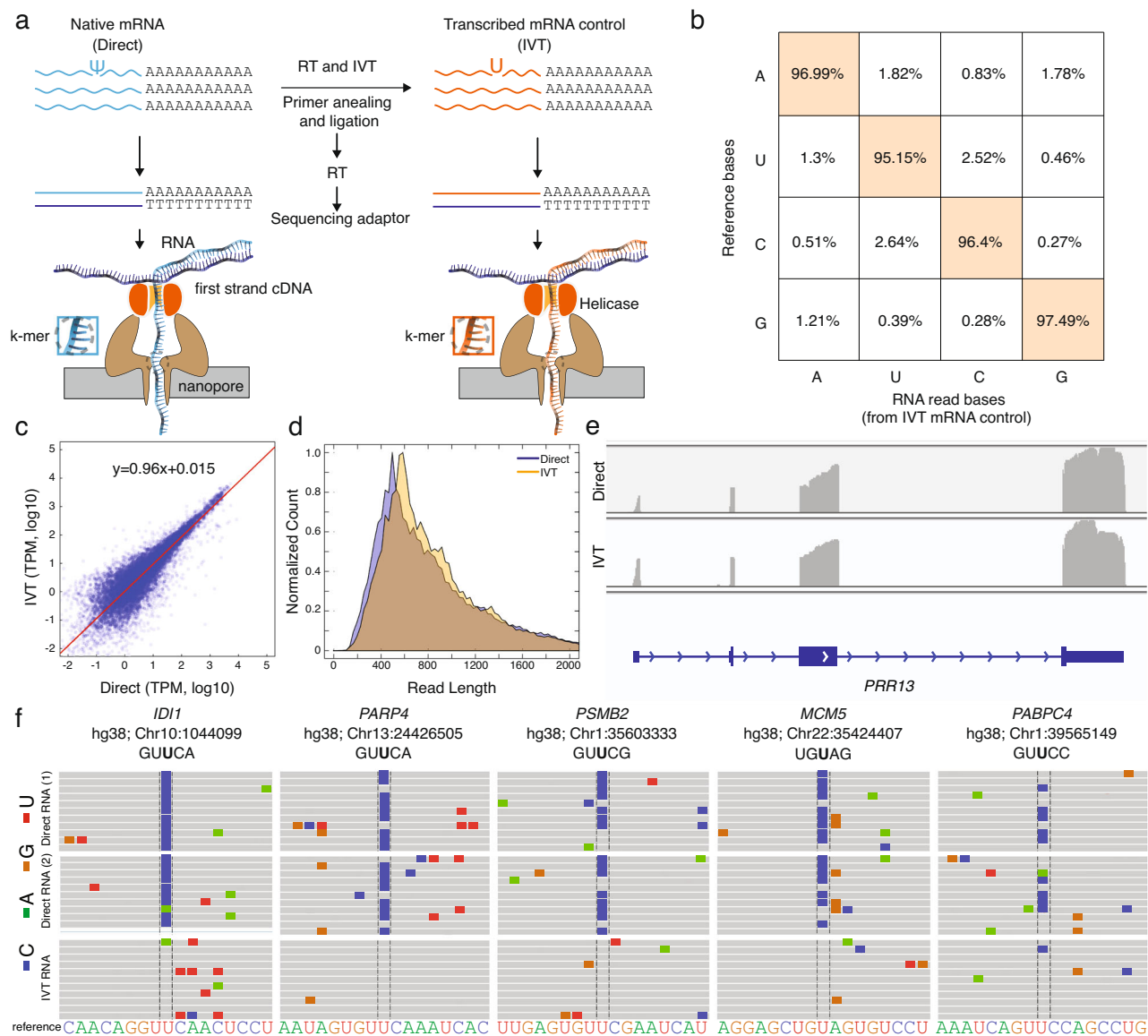


Fig. 1 | Nanopore native poly(A) RNA sequencing pipeline to identify ψ -modified sites. **a** Library preparation for Nanopore sequencing of native poly(A)-containing mRNAs (direct) and sequencing of in vitro transcribed (IVT) control. **b** The accuracy of called bases of in vitro transcribed (IVT) control samples. The x-axis shows called bases from nanopore reads and the y-axis is the base identity from the reference sequence at the same position that the nanopore reads are aligned to. **c** \log_{10} (TPM) of direct versus the \log_{10} (TPM) of IVT. **d** Normalized count

of different read lengths for direct reads (blue) versus IVT reads (orange). **e** IGV snapshot of *PRR13* in direct (top) and IVT (bottom). **f** Representative snapshot of aligned nanopore reads to the hg38 genome (GRCh38.p10) at previously annotated ψ sites. Correctly aligned bases are shown in gray, miscalled bases are shown in colors (cytosine, blue; adenine, green; guanine, orange; uridine, red). Genomic reference sequence is converted to sense strand and shown as RNA for clarity.

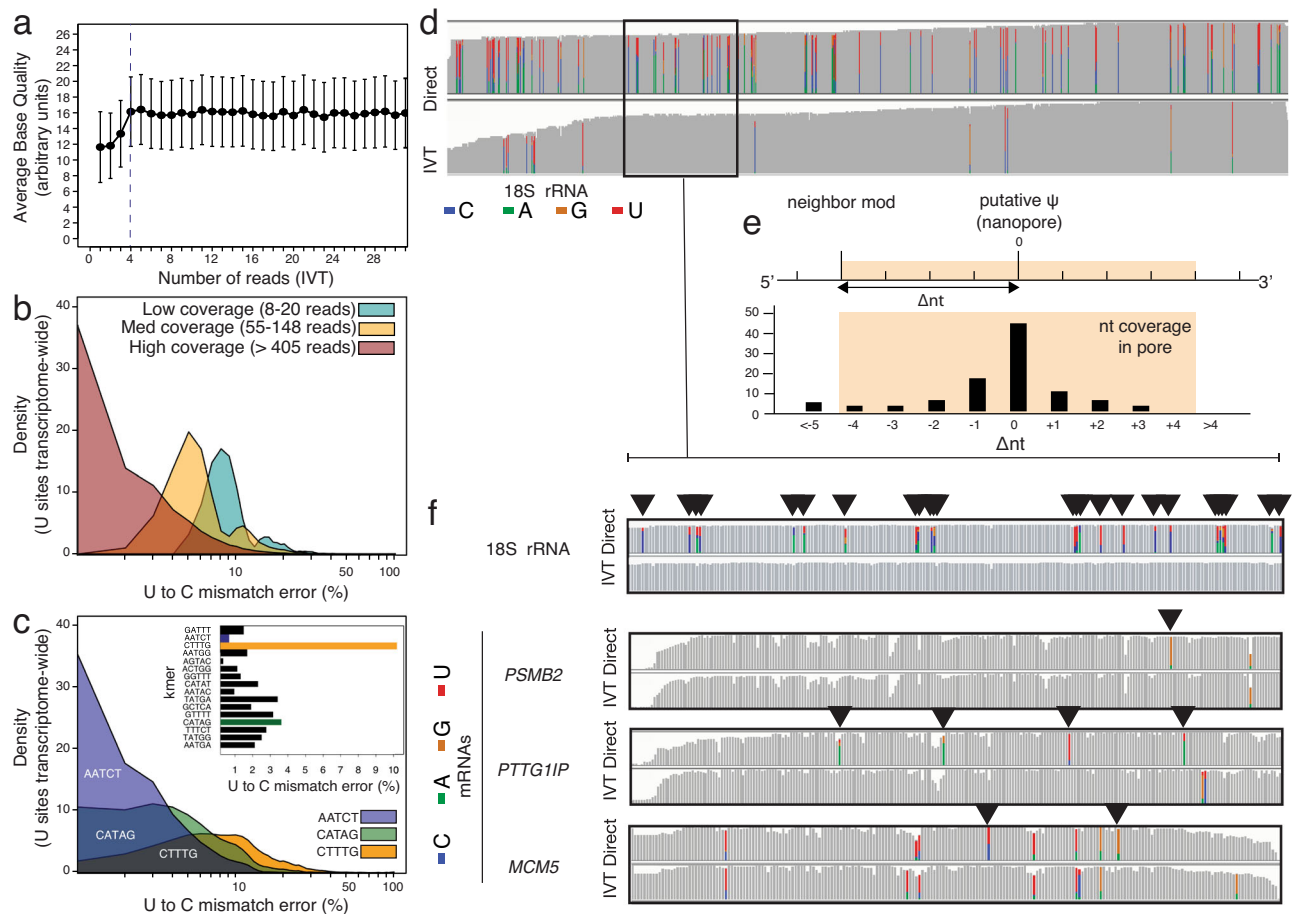


Fig. 2 | Basecalling errors can be used to detect RNA modifications if specific k-mer and coverage are considered, and the density of satellite modifications on rRNAs is different from mRNAs. **a** Average base quality for different numbers of reads using IVT reads. Data are represented as mean \pm standard deviation. **b** Distribution of U-to-C mismatch percentage for three populations based on read coverage (low coverage, teal; medium coverage, yellow; high coverage, red). **c** Distribution of U-to-C mismatch percentage for three populations based on 5-mer sequences (AATCT, blue; CATAG, green; CTTTG, yellow). **d** IGV snapshot of 18S rRNA for Direct (Upper) and IVT (lower). Correctly aligned bases are shown in gray,

miscalled bases are shown in colors (cytidine, blue; adenine, green; guanine, orange; uridine, red). **e** Schematic figure in which delta nucleotide is the distance to the putative modification position. **f** (Top) the IGV callout of a representative 200mer section of 18S rRNA and IGV snapshots of 200mer regions within 3 mRNAs with putative modifications. Presence of a basecalling error in the Direct and not in IVT denotes a putative site of modification indicated by a black triangle. Correctly aligned bases are shown in gray, miscalled bases are shown in colors (cytidine, blue; adenine, green; guanine, orange; uridine, red).

aligned section of their corresponding reads. It is common for the aligner to clip a few mismatched bases from the start/end of reads (known as soft-clipping). Therefore, to ensure sufficient coverage in both the direct RNA and IVT samples, we set a minimum threshold of 7 reads represented from each biological replicate before evaluation for site modification. We show that up to 3 bases adjacent to the soft clipped site usually yield lower base quality, and thus are not reliable regions to obtain information from (Supplementary Fig. 4).

To further investigate the U-to-C mismatch errors near the start/end of the alignment, we gathered the data for all the canonical uridine sites from our IVT control sample (>3 million uridine sites transcriptome-wide). For each of these positions, we calculated the U-to-C mismatch percentage, the number of aligned reads, and analyzed the surrounding bases of each site. We tabulated their 5-mers for which the target uridine site falls in the center. As expected, higher error rates were observed among low coverage sites (Fig. 2b). In addition, the surrounding bases of a site influenced the mismatch error (Fig. 2c). For example, uridine sites within the CUUUG k-mer, on average, showed a 10% mismatch error in the IVT reads, while uridine sites within the AAUCU k-mer had <0.4% average mismatch error. The average U-to-C mismatch of the specific k-mer in IVT is an important factor to be considered because it is essential to prevent a misinterpretation of the

inherent error of a k-mer as a site of modification. Therefore, the significance of the U-to-C mismatch percentage of a site must be interpreted based on a combination of the mismatch percentage and the number of reads in the direct RNA sample, as well as the average U-to-C mismatch error of the equivalent k-mer in the IVT sample.

It is important to ensure that the targets are not selected based on errors from other sources such as SNPs, basecalling, or alignment. In the cases that the IVT error at a specific position is higher than the average error of that k-mer, the mismatch error from the direct RNA reads is compared with error at the specific site rather than with the average error of that k-mer in IVT. To account for standard basecalling errors, we compare the direct reads to the IVT replicate, using the highest U-to-C error at that site in the direct reads.

Detection of annotated ψ sites ($p < 0.001$) on human rRNA

Human rRNA has been extensively annotated using mass spectrometry and is commonly used to benchmark the specificity of RNA modification. We generated and analyzed direct rRNA and IVT rRNA libraries from HeLa cells. A total of 43 previously annotated rRNA positions from the 18S and 5.8S subunits had sufficient coverage for analysis, while the 28S in the large subunit was omitted from analysis due to insufficient coverage. Of these sites, 38 (88.4%) were detected

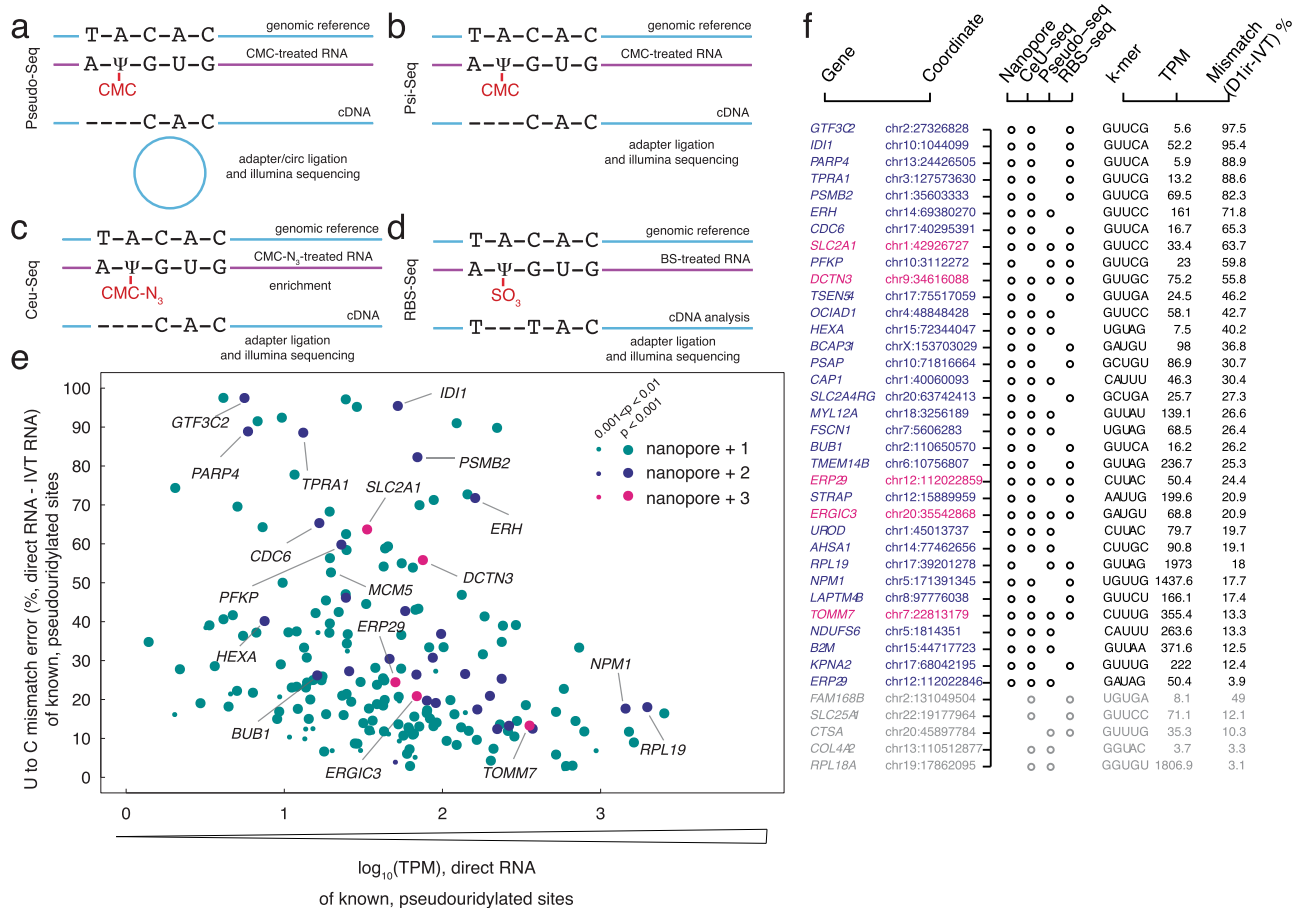


Fig. 3 | Previously annotated ψ modifications in the human transcriptome are validated by nanopore sequencing. **a** The schematic workflow of the CMC-based methods that have detected ψ modification in the human transcriptome. **a** Pseudo-Seq, **(b)** Ψ -Seq, **(c)** CeU-Seq, and **(d)** modified bisulfite sequencing (RBS-Seq). **e** U-to-C mismatch error (%) of the merged replicates of direct RNA of known ψ sites versus the \log_{10} (TPM) of merged direct RNA sequencing replicates. All targets shown are identified from the direct RNA sequencing data as likely to be pseudouridylated based on a P value calculation and are previously annotated by at least one previous

method. teal: annotated by one previous method, blue: annotated by two previous methods, magenta: annotated by three previous methods. P value was calculated by a one-sided F test (specific formula listed in Methods). $0.001 < p < 0.01$ is defined as significant and shown with a small dot; $p < 0.001$ is defined as highly significant and shown with a large dot. **f** The annotation of the genes containing a reported ψ modification by two or more previous methods (blue: annotated by two previous methods, magenta: annotated by three previous methods). The sites that are not validated by our nanopore method are shown in gray.

as ψ ($p < 0.001$; Fig. 2d, e, Supplementary Data 1). In addition to the known ψ sites that were determined by our algorithm, we detected 72 targets that are not on the list of previously detected ψ positions. Further inspection reveals that 10/72 of those positions exist in annotated rRNA positions not as ψ , but as 5-methyl uridine. A majority of the remaining positions (53/62) are within 4 bases of another previously annotated modification in each rRNA (Fig. 2d, e, Supplementary Fig. 5, Supplementary Data 1).

Comparison of mismatch error as a proxy for modification density on rRNA and mRNAs

Due to the extensive annotation and quantification of ψ sites on rRNA using mass spectrometry, they are considered the gold standard in the field for RNA modification detection and quantification. However, the distance between known modifications falls within the 12-nucleotide window of the biological nanopore, thus likely disrupting expected signal patterns, as shown in Fig. 2e. A comparison of basecalling error at sites in the direct reads and no error in the IVT control demonstrates the high density of modifications across the 18 S rRNA sequence with 25 errors within a ~200 nt window (Fig. 2d, f). In contrast, the 200 nt regions flanking putative modification sites on mRNAs (*PSMB2*, *PTTG1IP* and *MCMS5*) have 1–4 mismatch errors in the direct RNA samples as compared to the IVT controls.

Identification of previously annotated ψ sites on HeLa mRNA ($p < 0.001$)

Previous studies have identified putative ψ sites on human mRNA using biochemical methods including CMC^{3,17,18} and RNA bisulfite²³ (Fig. 3a–d). To evaluate the accuracy of our nanopore-based method in identifying previously annotated ψ sites, we generated a list of 334 putative ψ positions focusing on those that produced a minimum of 7 reads. We assigned p values (Fig. 2d) for each of the previously annotated ψ positions and found 232 positions with $p < 0.01$ (Fig. 3e). Among these positions, 198 sites were previously annotated by one other method and 34 were previously annotated by two or more methods²³ which we define as ground truth due to the identification of the same site by all three methods, i.e., 2+ methods and nanopore (Fig. 3f, Supplementary Data 2). For sites with sufficient coverage, our algorithm for determining ψ positions from direct RNA nanopore libraries had the highest overlap with Pseudo-seq (87.8%), followed by RBS-seq (77.9%), and had lowest overlap with CeU seq (67.6%).

Additional analysis of the positions that were previously annotated by 2 or more independent biochemical methods revealed only 5 positions that have sufficient coverage for analysis but were not identified as having a ψ by our algorithm (Fig. 3f). These positions include *COL4A2* (chr13:110512877), *RPL18A* (chr19:17862095), *CTSA* (chr20:45897784), and *SLC25A1* (chr22:19177964), each of which had a

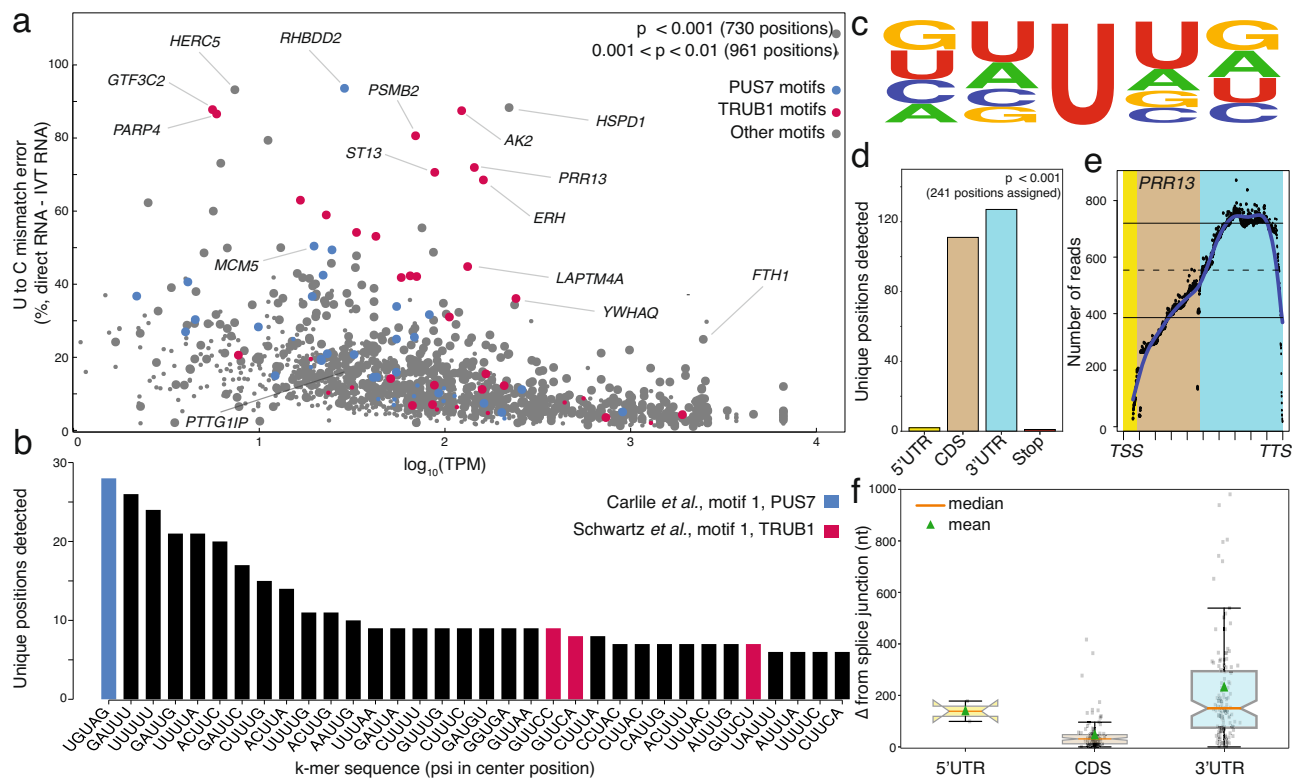


Fig. 4 | Nanopore sequencing detects uridine modifications transcriptome-wide. **a** The U-to-C mismatches detected by nanopore sequencing versus the $-\log_{10}(\text{TPM})$ of the merged direct RNA datasets. large dot: the detected targets identified by the significance factor of two out of three biological replicates. P value was calculated by a one-sided F test (specific formula listed in Methods). $0.001 < p < 0.01$ is defined as significant and shown with a small dot; $p < 0.001$ is defined as highly significant and shown with a large dot. blue: Targets with PUS7 motif, red: Targets with TRUB1 motif, and gray: Targets with the motifs other than PUS7 or TRUB1. **b** The k-mer frequency of the most frequently detected targets with higher confidence. **c** The sequence motif across the detected ψ modification for all

detected k-mers generated with kplogo⁴⁰. **d** The distribution of detected ψ sites in the 5' untranslated region (5' UTR; yellow), 3' untranslated region (3' UTR; blue), and coding sequence (CDS; brown). Sites were chosen based on a p value < 0.001 from (a). **e** The read depth of the reads aligned to *PRR13* versus the relative distance to the transcription start site (TSS) and transcription termination site (TTS). **f** The boxplots represent the distance of each detected site from the nearest splice junction for sites in the 5' UTR, 3' UTR, or CDS after reads were assigned to a dominant isoform using FLAIR³⁵. Median is shown with an orange line and mean is shown with a green triangle. Whiskers terminate at maxima/minima or a distance of 1.5 times the IQR away from the upper/lower quartile.

low mismatch error in direct RNA sequencing. Additionally, *FAM168B* (chr2:131049504) had high error in the corresponding IVT control, indicative of k-mer specific noise within those sequence contexts.

Transcriptome-wide detection of modified uridine sites in human cells

Next, we sought to apply our method for de novo detection of transcriptome-wide uridine modifications. We broadly encompassed ψ , DHU (dihydrouridine), and possibly Um (2'-O-methyl uridine) as uridine modifications that can occur in human transcriptome^{9,10}. To minimize the inclusion of sites of random error, we calculated significance based on the higher error of two replicates of analysis. We also required that two out of three direct replicates have the $p \leq 0.01$ to be defined as a site of uridine modification. Using this algorithm, we detected 1691 uridine-modification sites ($p < 0.01$), including 730 positions with a p value cutoff of 0.001 (Fig. 4a, Supplementary Data 3). Gene ontology analyses (GO Molecular Function 2021) were performed on genes with $p < 0.001$ using the enrichr website, showing that the RNA binding group has the highest normalized percentage of these genes. The enrichment of RNA binding group is also found in the GO of all identified transcripts (Supplementary Fig. 6, Supplementary Data 6). We then determined if Um was among the identified uridine-modified sites in our nanopore method. From a meta-analysis of known Um sites previously identified in HeLa cells⁹, we found no overlap with the sites called from our algorithm, indicating that the U-to-C base-calling does not report on the Um modification.

Uridine modifications are enriched within k-mer motifs of pseudouridine synthases

We assessed the k-mer frequencies for sites of uridine modification with $p < 0.001$ (Fig. 4b) and found that the k-mer UGUAG is the most highly represented, and that the k-mer GUUCN is among the most frequently detected. Note that UGUAG is the motif for PUS7 binding³⁴, while GUUCN is the motif for TRUB1²⁵, strongly indicating that they are sites of ψ modifications. To evaluate the sequence conservation of nucleotides within k-mers bearing a modification in the center position, we plotted the sequencing logo and found that the surrounding positions do not show any nucleotide preference (Fig. 4c).

Distribution of putative ψ -modified sites on mRNA sequences

We characterized the distribution pattern of sites that are most likely ψ modifications on mature mRNA transcripts and observed that ~60% of them are located on the 3' untranslated region (UTR) and 35% on the coding sequence (CDS), with very few targets detected in the 5' UTR (Fig. 4d). The limited detection of ψ sites in the 5' UTR is likely due to the low observed coverage in the 5' end of the RNA (i.e., near the transcription start site and covering a majority of the 5' UTR in many cases; Fig. 4e). Low coverage in the 5' ends of RNA is expected since the enzyme motor releases the last -12 nucleotides, causing them to translocate at speeds much faster than the limit of detection²⁶. Compared to the rest of the transcripts, there is also a sharp drop in coverage at the tail end of the 3' UTR, near the transcription termination site (Fig. 4e). Interestingly, we found one

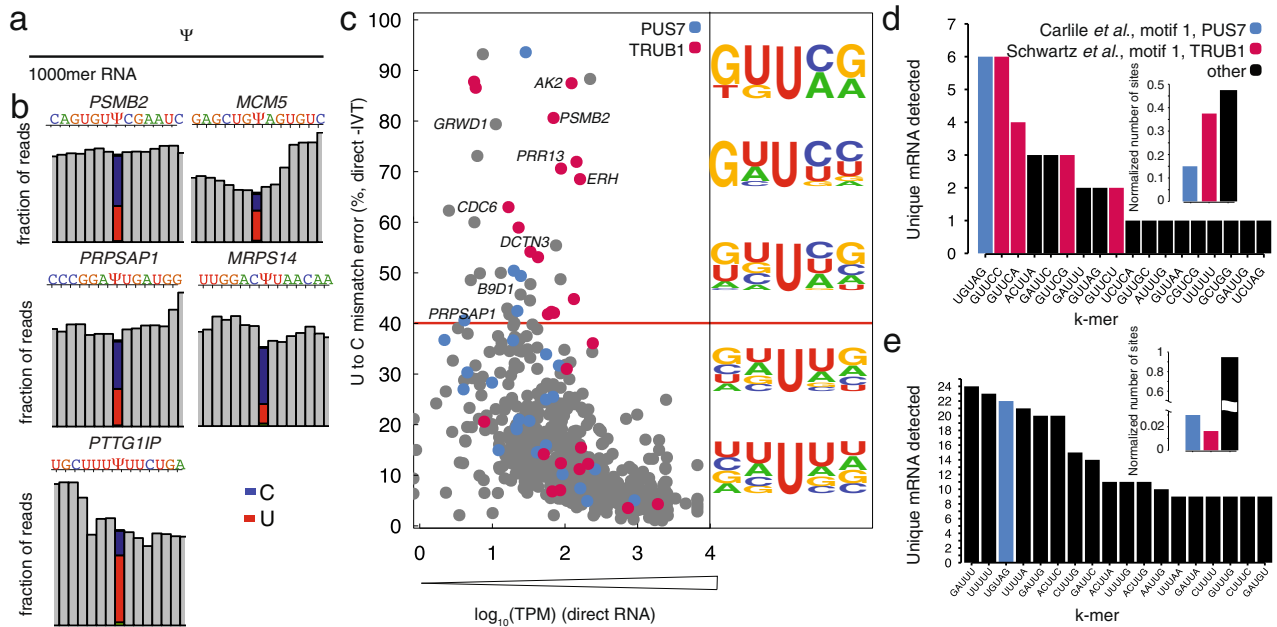


Fig. 5 | Solid-phase synthesis of 1000 mer RNA standards that maintain the sequence context of putative ψ sites demonstrates a systematic undercalling of the modification. **a** A pair of 1000-mer synthetic RNA oligos were designed, one containing 100% uridine and the other containing 100% ψ in the sequence context of a natural transcript. **b** The frequency histograms of 13 nucleotides surrounding the detected ψ position in the middle of a k-mer in four different mRNAs: *PSMB2*, *MCM5*, *PRPSAP1*, and *MRPS14*, and *PTTG1P*. **c** The U-to-C mismatches of the

detected ψ position for merged replicates of direct RNA seq versus $-\log_{10}(\text{significance})$. The targets with U-to-C mismatch of higher than 40% are defined as hypermodified type I. The sequence motifs for different mismatch ranges are shown. **d** K-mer frequency is shown for hypermodified type I and non-hypermodified ψ sites with the highest occurrence. **e** Distribution of U-to-C mismatches higher than 40% in mRNA regions. Insets for d and e show relative fractions of k-mers that are substrates of PUS7, TRUB1 and other motifs.

example of a putative ψ site within a transcription stop site: *GAGE2A* (chrX:49596731).

We calculated the distance of each putative ψ site from the closest splice site for high confidence ψ sites ($p < 0.001$). Prior to extracting the distance of the nearest splice junction for each target, we used the RNA isoform analysis tool, FLAIR³⁵ to bin the reads comprising high confidence modified targets into their respective dominant isoform. Overall, targets in the 3' UTRs are separated from a splice site by a longer distance relative to targets in CDS (Fig. 4f). Considering the vast differences in sequence length between CDSs and between the 3' UTRs, we observed a higher correlation between the splice distance of CDS-positioned targets and CDS length as compared to 3' UTR-positioned targets (Supplementary Fig. 6).

Systematic under calling of ψ percentage based on site-specific synthetic controls

To explore the quantitative potential of the U-to-C basecalling error as a proxy for pseudouridine, we constructed five synthetic mRNAs, each 1000 mer, bearing a pseudouridine at the nanopore detected site (Fig. 5a). These controls were designed to recapitulate the 1000 mer sequence flanking a natural ψ site in the human transcriptome, considering that the long-range sequence context can influence the current distributions of nanopore and basecalling for a given k-mer. Two of the chosen targets (*PSMB2*; chr1:35603333^{3,17,23} and *MCM5*; chr22:35424407^{3,17}) were annotated as ψ by two or more previous methods and the other three targets (*MRPS14*; chr1:175014468, *PRPSAP1*; chr17: 76311411, and *PTTG1P*; chr21:44849705) were detected de novo using the U-to-C mismatch error and our p value cutoff. For each site, we constructed a pair of RNA transcripts, in which the center position of the k-mer is either a uridine or a pseudouridine. We ran these synthetic controls through the nanopore directly and measured the U-to-C mismatch error for each. If the mismatch error were a perfect proxy for ψ , we expected to see 100% U-to-C mismatch in these synthetic controls. In contrast, we observed 38.17% U-to-C mismatch

error for *PSMB2*, 32.16% for *MCM5*, 69.64% for *PRPSAP1*, 69.35% for *MRPS14*, and 30.08% for *PTTG1P* (Fig. 5b). These results indicate a systematic under calling of ψ using our algorithm.

Sites with >40% U-to-C mismatch error are classified as type I hypermodification

We define hypermodification type I as a specific site within a transcript in which at least every other copy has a uridine modification. We therefore reasoned that a 40% mismatch error was an appropriate cutoff because the basecaller is systematically under-calling the modification. We further reasoned that, while somewhat arbitrary, at 40% mismatch error, we encompass all of the sites with 50% occupancy and above. From our de novo uridine-modification detection analysis, we identified 40 unique sites of hypermodification type I including *AK2* (chr1: 33014553), *ID11*(chr10:1044099), *GTF3C2*(chr2:196789267), *RHBDD2* (chr7:75888787), *HSPD1* (chr2:197486726) that show close to 100% mismatch error (Supplementary Data 4).

To assess the sequence conservation of nucleotides within k-mers bearing a putative ψ modification in the center position, we selected all unique sites with a U-to-C mismatch error above 40% (Supplementary Data 4). We found that the -1 position has a strong preference for uridine and the -2 position has a strong preference for guanosine. This preference becomes more significant as the mismatch percentage increases (Fig. 5c). The +1 position has a strong preference for cytosine especially at sites with > 80% U-to-C mismatch error.

We then assessed the k-mer frequencies for detected positions with a U-to-C mismatch error at >40% (Fig. 5d) and those with an error less than 40% (Fig. 5e). Among k-mers with a frequency >40% mismatch error, the GUUCN k-mer, representing the TRUB1 motif, is the most prominent (30/105 sites around 29%). The k-mer UGUAG, representing the PUS7 motif, is also prominent (5/105 sites around 4.8%). In contrast, the k-mers UGUAG (13/712, 1.8%), GUUCN, and all others occurred at a similar frequency as the most abundant non-hypermodified targets (15/712, 2.1%).

We assessed the location of putative ψ on the transcript and found that type I hypermodified sites are biased toward 3' UTRs, which is the same as sites that are not hypermodified (Supplementary Fig. 7). No significant difference was observed in the splice distance of type I hypermodified sites between sites in the 3' UTR and those in CDS regions of mRNA when compared to not-hypermodified sites (Supplementary Fig. 7).

Messenger RNAs with >1 uridine modifications are classified as type II hypermodification

We define hypermodification type II as the mRNAs that can be modified at two or more positions. Using only the sites with a high probability of modification ($p < 0.001$), we identified 104 mRNAs with modifications at 2 unique positions, 27 with 3 positions, 4 with 4 positions, 5 with 5 positions, 1 with 6 positions and 1 mRNA with 7 positions (Fig. 6a, Supplementary Data 5). For the mRNAs that are modified at 2 positions, we found no correlation between the mismatches in position 1 and position 2 ($R^2 = 0.039$), as expected because this error is highly dependent on the k-mer sequence. To demonstrate future applications of long-read sequencing for assessing modifications on the same transcript, we plotted every read for two mRNAs (*CHTOP* and *PABPC4*) and labeled each site using the called base (Fig. 6b). We observed that these mismatches could happen on the same read. For example, 17% and 7% of the reads had a U-to-C mismatch in both positions for *CHTOP* and *PABPC4* respectively. We plotted the distribution of type II hypermodifications across the body of the transcript and found a slight clustering of sites in the 3' UTR (Fig. 6c).

Discussion

We show here that systematic U-to-C basecalling errors detected from direct nanopore sequencing of transcriptomes can serve as indicators for annotated sites of ψ modification. In this work, we provide a foundation for identifying modified U sites with high confidence based on two approaches. In the first, U-to-C mismatch errors in native transcriptomes are compared against a corresponding unmodified transcriptome as a negative control to eliminate standard basecalling errors that occur in canonical bases. We weigh the transcriptome wide average of U-to-C errors in k-mers to minimize false positives due to low coverage that is an issue with direct mRNA sequencing. In the second, we use a set of long synthetic RNA controls with precisely positioned ψ modifications to aid in our discovery of systematic undercalling of ψ modifications, pointing to limitations of basecalling-guided RNA modification detection algorithms. Our approach is distinct from the ELIGOS algorithm³⁶, primarily because we consider the average U-to-C error of unmodified k-mers, enabling analysis of low coverage sites that may show as significant error due to random nanopore basecalling error. In addition, we determine individual modification sites by the exclusive presence of U-to-C mismatches rather than including all other substitutions, deletions, and insertions at a given site; thus, reducing false positive detection.

We demonstrate that this method can identify a majority of ψ sites that were detected by CMC and bisulfite-based next-generation sequencing platforms. Importantly, we produce a ground truth list of ψ positions in mRNAs that are detected by nanopore and by at least one other method that has previously annotated ψ in HeLa cells. This ground-truth list consists of 198 mRNA positions detected by nanopore and one other method, and 34 positions detected by nanopore and two other methods (Fig. 3), constituting a conservative list of targets that we suggest biologists focus on for performing functional analysis. This is a significant advance in the field because the current methods to identify ψ primarily rely on CMC modification and there is a danger in orthogonal methods relying on the same chemical. Although alternative chemical methods have been introduced, the overlap among them is relatively low. Our demonstration that

nanopore sequencing can detect modifications at the same site as chemical methods strengthens the identification at each site and demonstrates the rigor of our approach. These results suggest that the reported ground truth list can serve as a first step for developing biological assays to determine the specific function of individual ψ modifications in mRNAs.

Among the methods that we used to validate our data, Pseudo-seq shows the highest overlap between the detection targets. However, some targets that the other methods detected were not detected by our method. We demonstrate that these sites may be due to problematic k-mers that show a high error in the IVT control, or very low coverage, thereby not meeting the criteria for being identified as a modified site. Alternatively, artifacts from CMC labeling may account for this, including incomplete CMC adduct removal from unmodified uridines, reverse-transcriptase read-through of CMC-modified ψ sites, or uneven amplification of low-occupancy ψ -sites. Another potential reason for the differences could be batch differences between cell lines leading to differential occupancy at a given moment of mRNA extraction. We also observed several targets that were detected by our nanopore method that were not detected by other methods. While we are confident that these sites are modified due to differences between the native RNA versus the IVT control, we cannot rule out the possibility of other uridine modifications (as shown in the identification of 5'-methyl-uridines from the rRNA data). Nonetheless, many of these targets are likely ψ due to the enrichment of k-mer sequences that are motifs for ψ synthases PUS7 and TRUB1. Further studies are necessary to resolve this issue.

We note that our approach cannot rule out other types of uridine modifications, such as DHU^{10,11} or Um⁹ that can occur in mRNA. The estimated frequency and stoichiometry of DHU in eukaryotic mRNAs is low ($< 0.01\%$ ¹⁰) relative to ψ (0.2–0.6%³), but this modification may also fall within uridine-rich tracks¹⁰. Therefore, we recommend that these modified uridine sites be further explored as putative DHU modifications using appropriate DHU-seq¹⁰. While the frequency of Um is comparable to ψ (0.15%⁹), we performed a meta-analysis of known Um sites from HeLa cells⁹ and found no overlap with the sites called from our algorithm and de novo analysis (data not shown), indicating that the U-to-C basecalling in nanopore sequencing likely does not report on Um. A better discrimination among different types of uridine modifications occurring as a U-to-C mismatch from direct RNA nanopore sequencing will require additional synthetic controls and improved machine-learning algorithms. Parallel analysis of a transcriptome of interest with the transcriptome of a cell line lacking or reducing the expression of one of the ψ synthases will help to assign a specific modification de novo for putative ψ sites^{17,18,29}.

Previous studies have demonstrated the importance of long-range interactions³¹ for accurate calling of ψ modifications by direct RNA sequencing. To account for the contributions of long-range interactions, we have validated our method by analysis of five synthetic 1000 mers, each containing a site-specific ψ found within a natural target sequence in the human transcriptome. We find that the U-to-C basecalling error is systematically under-called at the modified site. Based on this finding, we defined hypermodification type I as sites that have >40% U-to-C mismatch error with the rationale that higher modification frequencies are more likely to have biological significance. We also define hypermodification type II as mRNAs bearing multiple U modification sites in a specific transcript. Finally, we show that U modifications can occur up to seven times on a single transcript.

A fully quantitative measure of ψ occupancy at a given site would require high-coverage sequencing runs of a comprehensive set of every possible ψ -containing k-mer within its natural sequence context (an estimated 13 nucleotides surrounding the modified site). While high-throughput controls have previously been generated^{29,31}, all Us were modified in those studies and consequently, these are not the ideal controls for detection of single ψ modifications within the natural

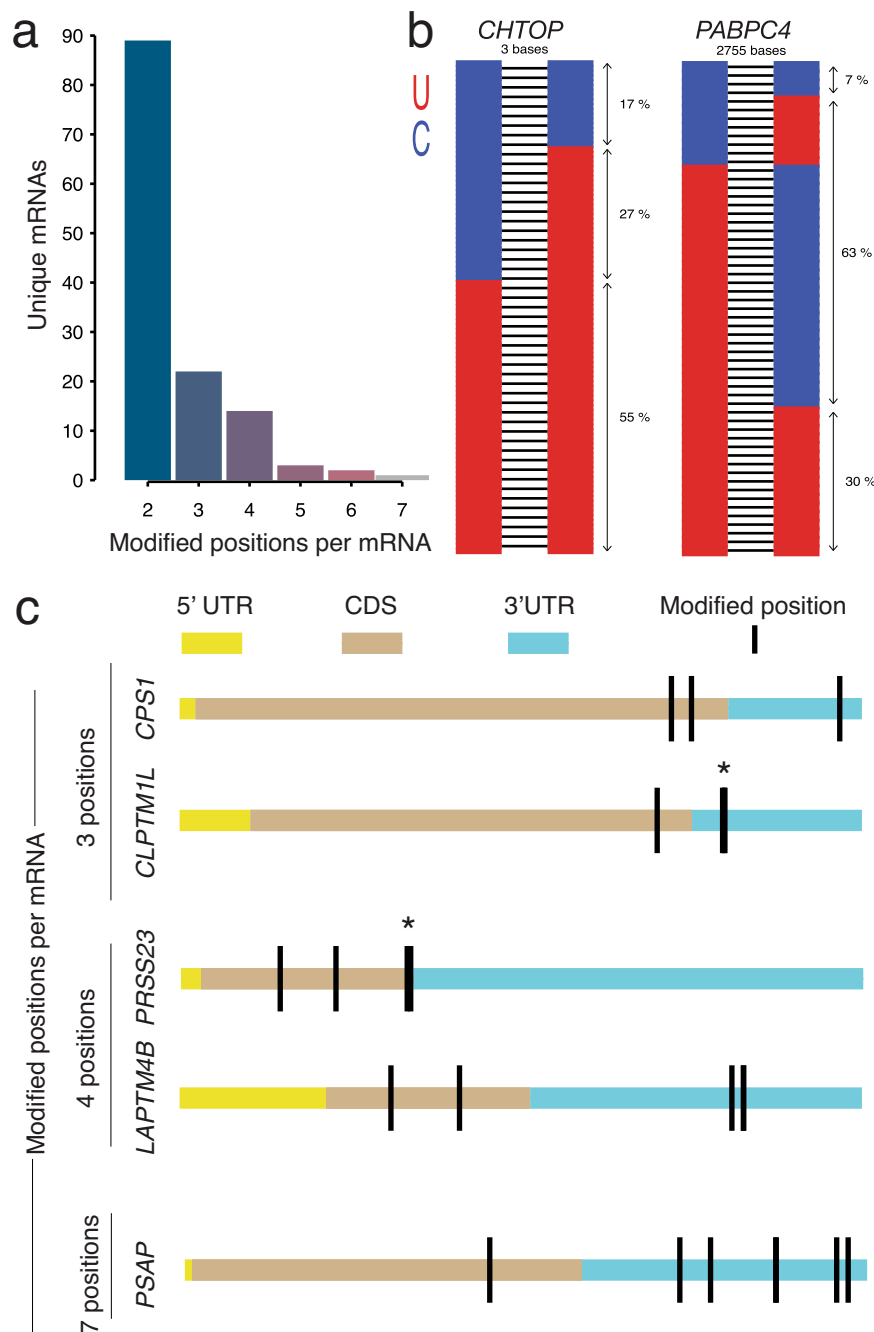


Fig. 6 | Type II hypermodification is defined as the mRNA targets that contain two or more uridine-modified positions. **a** Unique mRNAs that are classified as hypermodification type II positions and the number of modified positions possible on each. **b** Two examples of hypermodified type II transcripts (*CHTOP* –

chr1:153,645,392-153,654,395 & *PABPC4* – chr1:39,562,394-39,565,149) with two modified positions indicating U-to-C mismatch on a single read for long reads that cover both positions. **c** Examples of type II hypermodification with three or more modified positions distributed across each gene.

sequence contexts as demonstrated in the comparison of rRNA modifications to mRNA modifications (Fig. 2). In addition, rRNA-based controls, which contain high frequencies of conserved modification sites, are not ideal for assessing mRNA modifications, because the density of modifications in rRNA is much higher than that in mRNA (see Fig. 2). Although synthetic controls are most critical to resolve ambiguities of ψ detection by nanopore sequencing, the preparation of a large number of such molecules by individual synthesis of each is not feasible for a single laboratory. Future work will consider innovative methods to make libraries of all putative sites within their sequence contexts in order to quantitatively evaluate ψ profiles in transcriptomes.

The synthetic controls that we have generated demonstrate that the basecalling error is reliable in the calling of ψ at a given site. By setting a cutoff in U-to-C mismatch for a given site, we conservatively draw a list of high-confidence sites that are pseudouridylated with high occupancy, and thus, have a higher likelihood of leading to a measurable phenotype in the cell and conferring a functional impact on the cellular physiology. Our work serves as a foundation for detection and analysis of ψ and other uridine modifications on mRNAs with sequence specificity and single-molecule resolution. Future work should include an expansion of synthetic controls and training of an updated basecaller to improve our ability to distinguish, validate and quantify U modifications in transcriptomes.

Methods

Cell culture

HeLa cells were cultured in Dulbecco's modified Eagle's medium (Gibco, 10566024), supplemented with 10% Fetal Bovine Serum (FB12999102, FisherScientific) and 1% Penicillin-Streptomycin (Lonza, 17602E). To extract sufficient poly-A RNA, three confluent, 10 cm dishes were used for each experiment.

Total RNA extraction and Poly(A) RNA isolation

The total RNA extraction protocol was performed using a method that is the combination of total RNA extraction using TRIzol (Invitrogen, 15596026) and PureLink RNA Mini Kit (Invitrogen, 12183025). Cells were washed with 3 ml ice-cold PBS and added with 2 ml of TRIzol in each 10 cm dish and incubated at room temperature for 5 min. Every 1 ml of lysed cells in TRIzol was transferred to a LoBind Eppendorf tube and vortexed for 30 s and added with 200 μ l chloroform (Acros Organics, 423555000) in each tube and mixed by shaking for 15 s and incubated at room temperature for 3 min. Then the samples were centrifuged at 12000 \times G for 15 min at 4 $^{\circ}$ C and 0.4 ml of aqueous supernatant was transferred to a fresh LoBind Eppendorf tube and an equal volume of 70% ethanol was added to the solution followed by vortexing. In the following steps, PureLink RNA Mini Kit (Invitrogen, 12183025) and the protocol were performed according to the manufacturer's recommended protocol. Briefly, the solution was transferred to a pure link silica spin column and flow-through was discarded (every two microtubes were loaded on one column). The columns were washed with 0.7 ml of wash buffer I once and then with 0.5 ml wash buffer II twice. The total RNA was eluted using 50 μ l nuclease-free water. The RNA concentration was measured using a NanoDrop 2000/2000c Spectrophotometer.

NEBNext Poly(A) mRNA Magnetic Isolation Module (E7490L) was used to select poly(A) mRNA. The protocol was followed according to the manufacturer's protocol. The only modification was pooling 5 samples and performing the experiment in microtubes instead of PCR tubes. A total of 15 samples (3 microtubes) was used in each experiment to get enough Poly-A RNA product. The products were eluted from the NEBNext polyA magnetic isolation (NEB, E7490S) in Tris buffer. The three samples were pooled and ethanol precipitated to get to the concentration that is required for the sequencing step.

In vitro transcription, capping, and polyadenylation

Primers used are listed in Supplementary Data 7. cDNA-PCR Sequencing Kit (SQK-PCS109) kit was used for reverse transcription and strand-switching. Briefly, VN primer and Strand-Switching Primer were added to 50 ng poly-A RNA. Maxima H Minus Reverse Transcriptase (Thermo scientific, EP0751) was used to produce cDNA. IVT_T7 Forward and reverse primers were added to the product and PCR amplified using LongAmp Taq 2X Master Mix (NEB, M0287S) with the following cycling conditions: initial denaturation 30 s @ 95 $^{\circ}$ C (1 cycle), denaturation 15 s @ 95 $^{\circ}$ C (11 cycles), annealing 15 s @ 62 $^{\circ}$ C (11 cycles), extension 15 min @ 65 $^{\circ}$ C (11 cycles), final extension 15 mins @ 65 $^{\circ}$ C (1 cycle), and hold @ 4 $^{\circ}$ C. Each PCR product was added with 1 μ l of Exonuclease 1 (NEB, M0293S) and incubated at 37 $^{\circ}$ C for 15 min to digest any single-stranded product, followed by 15 min at 80 $^{\circ}$ C to inactivate the enzyme. Sera-Mag beads (9928106) were used according to the Manufacturer's protocol to purify the product. The purified product was then IVT using HiScribe T7 High yield RNA Synthesis Kit (NEB, E2040S) and purified using Monarch RNA Cleanup Kit (NEB, T2040S). The product was eluted in nuclease-free water and poly-A tailed using E. coli Poly(A) Polymerase (NEB, M0276). The product was purified once again using an RNA Cleanup Kit and adjusted to 500 ng polyA RNA in 9 μ l NF water to be used in the Direct RNA library preparation.

For rRNA IVT, total RNA was poly-A tailed using E. coli Poly(A) Polymerase (NEB, M0276) and purified using RNA Cleanup kit (NEB,

T2040S) then poly-A selected using NEBNext polyA magnetic isolation (NEB, E7490S). The poly-A tailed total RNA (50 ng) was then IVT according to the above protocol.

Synthetic sequence design

We constructed five synthetic 1000 mer RNA oligos, each with a site-specifically placed k-mer. Two versions of each RNA were prepared, one with 100% uridine and the other with 100% ψ at the central position of the k-mer. The uridine-containing RNAs were prepared by T7 transcription from G-block DNAs (synthesized by Integrated DNA Technologies), whereas the ψ -containing RNAs were prepared by ligation of left and right RNA arms (each 500 nts in length) to a 15 mer RNA bearing a ψ in the central position (synthesized by GeneLink). A T7 promoter sequence with an extra three guanines was added to all the DNA products to facilitate in vitro transcription. In addition, a 10 nt region within 30 nt distance of ψ was replaced by a barcode sequence to allow parallel sequencing of the uridine- and ψ -containing samples. Finally, each left arm was transcribed with a 3' HDV ribozyme that self-cleaved to generate a homogeneous 3'-end. Full-length RNA ligation products were purified using biotinylated affinity primers that were complementary to both the left and right arms.

Direct RNA library preparation and sequencing

The RNA library for Direct RNA sequencing (SQK-RNA002) was prepared following the ONT direct RNA sequencing protocol version DRCE_9080_v2_revH_14Aug2019. Briefly, 500 ng poly-A RNA or poly-A tailed IVT RNA was ligated to the ONT RT adaptor (RTA) using T4 DNA Ligase (NEB, M0202M). Then the product was reverse transcribed using SuperScriptTM III Reverse transcriptase (Invitrogen, 18080044). The product was purified using 1.8 \times Agencourt RNAClean XP beads, washed with 70% ethanol and eluted in nuclease-free water. Then the RNA: DNA hybrid was ligated to RNA adapter (RMX) and purified with 1 \times Agencourt RNAClean XP beads and washed twice with wash buffer (WSB) and finally eluted in elution buffer (ELB). The FLO-MIN106D was primed according to the manufacturer's protocol. The eluate was mixed with an RNA running buffer and loaded to the flow cell. MinKnow (19.12.5) was used to perform sequencing. Three replicates from different passages and different flow cells were used for each biological replicate. For Direct rRNA library preparation, total RNA was poly-A tailed using E. coli Poly(A) Polymerase (NEB, M0276) and purified using RNA Cleanup kit (NEB, T2040S) as in the above protocol.

Base-calling, alignment, and signal intensity extraction

Multi-fast5 files were basecalled in real-time by guppy (3.2.10) using the high accuracy model.

Then, the reads were aligned to the genome version hg38 using minimap 2 (2.17) with the option "-ax splice -uf -kl4". The sam files were converted to bam using samtools (2.8.13). Bam files were sorted by samtools sort and indexed using samtools index and visualized using IGV (2.8.13). The bam files were sliced using samtools view and the signal intensities were extracted using nanopolish eventalign (0.13.2).

Gene ontology and sequencing logo analysis

Gene ontology (GO) analysis of Molecular Function 2021 was performed using enrichr website³⁷⁻³⁹. The sequence motifs are generated by kpLogo website⁴⁰.

Modification detection and analysis

A summary of the base calls of aligned reads to the reference sequence was obtained using the *Rsamtools* package. Mismatch frequency was then calculated for a list of verified ψ sites. We observed that U-to-C mismatch frequency shows a better separation between the modified (IVT) and (potentially) modified (Direct) samples.

We know from our control sample that U-to-C mismatch frequency depends on both the molecular sequence and coverage (Fig. 2.

a–c). Therefore, the significance of an observed mismatch percentage at each site was calculated accordingly via Eq. (1):

$$p(N, N_{mm, dseq}, p_0) = \sum_{N_{mm} = N_{mm, dseq}}^N = \binom{N}{N_{mm}} \times p_0^{N_{mm}} \times (1 - p_0)^{N - N_{mm}}$$

where the significance of the mismatch frequency at each U site was calculated using the sequence-dependent expected error and the read coverage at that site.

Statistical analysis

All experiments were performed in multiple, independent experiments, as indicated in figure legends. All statistics and tests are described fully in the text or figure legend.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The data that support this study are available from the corresponding author upon reasonable request. Sequences were aligned to genome version hg38. All FASTQ files and Fast5 raw data generated in this work have been made publicly available in NIH NCBI SRA under the BioProject accession [PRJNA777450](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA777450).

Code availability

All code used in this work is publicly available at github.com/RouhanifardLab/PsiNanopore.

References

- Roundtree, I. A., Evans, M. E., Pan, T. & He, C. Dynamic RNA Modifications in Gene Expression Regulation. *Cell* **169**, 1187–1200 (2017).
- Taoka, M. et al. Landscape of the complete RNA chemical modifications in the human 80S ribosome. *Nucleic Acids Res.* **46**, 9289–9298 (2018).
- Li, X. et al. Chemical pulldown reveals dynamic pseudouridylation of the mammalian transcriptome. *Nat. Chem. Biol.* **11**, 592–597 (2015).
- Mellis, I. A., Gupte, R., Raj, A. & Rouhanifard, S. H. Visualizing adenosine-to-inosine RNA editing in single mammalian cells. *Nat. Methods* **14**, 801–804 (2017).
- Spitale, R. C. et al. Structural imprints in vivo decode RNA regulatory mechanisms. *Nature* **519**, 486–490 (2015).
- Wang, X. et al. N6-methyladenosine-dependent regulation of messenger RNA stability. *Nature* **505**, 117–120 (2014).
- Hunter, S. et al. InterPro in 2011: new developments in the family and domain prediction database. *Nucleic Acids Res.* **40**, D306–D312 (2012).
- Waldo, E. C. & Elliot, V. Nucleoside-5'-Phosphates from Ribonucleic Acid. *Nature* **167**, 483–484 (1951).
- Dai, Q. et al. Nm-seq maps 2'-O-methylation sites in human mRNA with base precision. *Nat. Methods* **14**, 695–698 (2017).
- Finet, O. et al. Transcription-wide mapping of dihydrouridine reveals that mRNA dihydrouridylation is required for meiotic chromosome segregation. *Mol. Cell* **82**, 404–419.e9 (2022).
- Draycott, A. S. et al. Transcriptome-wide mapping reveals a diverse dihydrouridine landscape including mRNA. *PLoS Biol.* **20**, e3001622 (2022).
- Anderson, B. R. et al. Nucleoside modifications in RNA limit activation of 2'-5'-oligoadenylate synthetase and increase resistance to cleavage by RNase L. *Nucleic Acids Res.* **39**, 9329–9338 (2011).
- Price, A. M. et al. Direct RNA sequencing reveals m6A modifications on adenovirus RNA are necessary for efficient splicing. *bioRxiv* 865485 <https://doi.org/10.1101/865485> (2019).
- Karikó, K., Buckstein, M., Ni, H. & Weissman, D. Suppression of RNA recognition by Toll-like receptors: the impact of nucleoside modification and the evolutionary origin of RNA. *Immunity* **23**, 165–175 (2005).
- Anderson, B. R. et al. Incorporation of pseudouridine into mRNA enhances translation by diminishing PKR activation. *Nucleic Acids Res.* **38**, 5884–5892 (2010).
- Eyler, D. E. et al. Pseudouridylation of mRNA coding sequences alters translation. *Proc. Natl Acad. Sci. U. S. A.* **116**, 23068–23074 (2019).
- Schwartz, S. et al. Transcriptome-wide mapping reveals widespread dynamic-regulated pseudouridylation of ncRNA and mRNA. *Cell* **159**, 148–162 (2014).
- Carlile, T. M. et al. Pseudouridine profiling reveals regulated mRNA pseudouridylation in yeast and human cells. *Nature* **515**, 143–146 (2014).
- Charette, M. & Gray, M. W. Pseudouridine in RNA: what, where, how, and why. *IUBMB Life* **49**, 341–351 (2000).
- Mengel-Jørgensen, J. & Kirpekar, F. Detection of pseudouridine and other modifications in tRNA by cyanoethylation and MALDI mass spectrometry. *Nucleic Acids Res.* **30**, e135 (2002).
- Addepalli, B. & Limbach, P. A. Mass spectrometry-based quantification of pseudouridine in RNA. *J. Am. Soc. Mass Spectrom.* **22**, 1363–1372 (2011).
- Ho, N. W. & Gilham, P. T. Reaction of pseudouridine and inosine with N-cyclohexyl-N'-beta-(4-methylmorpholinium)ethylcarbodiimide. *Biochemistry* **10**, 3651–3657 (1971).
- Khoddami, V. et al. Transcriptome-wide profiling of multiple RNA modifications simultaneously at single-base resolution. *Proc. Natl Acad. Sci. U. S. A.* **116**, 6784–6789 (2019).
- Dai, Q. et al. Quantitative sequencing using BID-seq uncovers abundant pseudouridines in mammalian mRNA at base resolution. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-022-01505-w> (2022).
- Safra, M., Nir, R., Farouq, D., Vainberg Slutskin, I. & Schwartz, S. TRUB1 is the predominant pseudouridine synthase acting on mammalian mRNA via a predictable and conserved code. *Genome Res.* **27**, 393–406 (2017).
- Workman, R. E. et al. Nanopore native RNA sequencing of a human poly(A) transcriptome. *Nat. Methods* **16**, 1297–1305 (2019).
- Liu, H. et al. Accurate detection of m6A RNA modifications in native RNA sequences. *Nat. Commun.* **10**, 4079 (2019).
- Smith, A. M., Jain, M., Mulroney, L., Garalde, D. R. & Akeson, M. Reading canonical and modified nucleobases in 16S ribosomal RNA using nanopore native RNA sequencing. *PLoS ONE* **14**, e0216709 (2019).
- Begik, O. et al. Quantitative profiling of pseudouridylation dynamics in native RNAs with nanopore sequencing. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-021-00915-6> (2021).
- Huang, S. et al. Interferon inducible pseudouridine modification in human mRNA by quantitative nanopore profiling. *Genome Biol.* **22**, 330 (2021).
- Fleming, A. M., Mathewson, N. J., Howpay Manage, S. A. & Burrows, C. J. Nanopore Dwell Time Analysis Permits Sequencing and Conformational Assignment of Pseudouridine in SARS-CoV-2. *ACS Cent. Sci.* <https://doi.org/10.1021/acscentsci.1c00788> (2021).
- Pyle, A. M. Translocation and unwinding mechanisms of RNA and DNA helicases. *Annu. Rev. Biophys.* **37**, 317–336 (2008).
- Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
- Carlile, T. M. et al. mRNA structure determines modification by pseudouridine synthase 1. *Nat. Chem. Biol.* **15**, 966–974 (2019).
- Tang, A. D. et al. Full-length transcript characterization of SF3B1 mutation in chronic lymphocytic leukemia reveals downregulation of retained introns. *Nat. Commun.* **11**, 1438 (2020).

36. Jenjaroenpun, P. et al. Decoding the epitranscriptional landscape from native RNA sequences. *Nucleic Acids Res.* **49**, e7 (2021).
37. Chen, E. Y. et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinforma.* **14**, 128 (2013).
38. Kuleshov, M. V. et al. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Res.* **44**, W90–W97 (2016).
39. Xie, Z. et al. Gene Set Knowledge Discovery with Enrichr. *Curr. Protoc.* **1**, e90 (2021).
40. Wu, X. & Bartel, D. P. kpLogo: positional k-mer analysis reveals hidden specificity in biological sequences. *Nucleic Acids Res.* **45**, W534–W538 (2017).

Acknowledgements

S.H.R. acknowledges support from a Seed Networks Award from the Chan Zuckerberg Initiative CZF2019-002424 and NIH 5R01HG011087-02. M.W. acknowledges support from NIH R01HG10087 and Oxford Nanopore Technologies. Y.M.H. acknowledges support from NIH HG011120.

Author contributions

S.T., M.W., and S.H.R. conceived of the research. S.T. designed and performed the experiments. S.T., M.N., A.M., C.A.M., and N.R. analyzed the data with guidance from M.W. and S.H.R. H.G. designed and synthesized synthetic RNA controls with guidance from Y.M.H. S.T. wrote the paper with guidance from Y.M.H. and S.H.R.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-023-35858-w>.

Correspondence and requests for materials should be addressed to Sara H. Rouhanifard.

Peer review information *Nature Communications* thanks Kunqi Chen and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at <http://www.nature.com/reprints>

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

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

© The Author(s) 2023