

Detecting Changes in the Epitranscriptome by Re-Interpreting M6aNET

Theodore M. Nelson

Introduction

Recent advances in nanopore technologies have permitted the first ever sequencing of RNA molecules, via a commercially available platform.¹ These sequencers allow for antibody-independent detection of RNA modifications on a transcriptome-wide basis, based on differential current signal within a nanopore, an engineered membrane-like protein.² Machine learning classification software, such as *m6anet*, have shown admirable applicability and sensitivity in resolving these modifications, firstly due to their dispensing of a negative control lacking m6a modifications, and secondly due to their performance in comparative studies.³

In training *m6anet*, the original source data for direct-RNA peaks were identified with another program, xPORE, by comparing sequenced *METTL3*-KO and WT HEK293T cell lines.⁴ Candidates DRACH sites were identified and labeled via m6ACE-Seq data, to generate a training dataset for *m6anet*, which outputs probability estimates for each covered DRACH motif covered within a human sample transcriptome. Within this article, I argue that these probability estimates are unreliable, given that they generalize across different RNA contexts and diverse cell types, which were not originally part of the training data.

VIRMA is a recently characterized component of the m6a methyltransferase complex, along with *METTL3*, regulating m6a deposition in the 3' untranslated region near protein stop codons.⁵ *VIRMA*-knockdown has been associated with a decrease in m6a deposition.⁶ I analyzed *VIRMA*-knockdown samples from direct-RNA sequencing runs from Oxford Nanopore Technologies 004 platform, comparing them with WT controls and in-vitro transcribed RNA, to propose a novel technique to identify differential methylation in a context-dependent manner.

Instead of assigning methylation status based on *m6anet*, I identified differentially methylated sites by examining the distribution of probabilities within an individual site. By examining these predictions, I can determine whether the model recognizes two distinct probability distributions, irrespective of their assigned mean probabilities.

By determining methylation sets in this way, I aim to propose this methodology as a more sensitive way of identifying *VIRMA*-knockdown affected m6a sites, generalizable to a wider range of biological questions, conditions, and direct-RNA sequencing runs.

Methods

Data. I analyzed six samples from the B lymphoblast cell line MM1S (Table 1) and one in-vitro-transcribed HCV sample, sequenced with the Oxford Nanopore Technologies Direct RNA Sequencing Kit (SQK-RNA004), according to manufacturer's instructions. Sequences were basecalled with Dorado software version 7.2.13+fba8e8925, client-server API version 16.0.0, with base calling model `rna_rp4_130bps_hac_prom.cfg`. Minimap2 version 2.26 was used to

generate alignments.⁷ These were reconnected to the original squiggle alignments, kept in the .blow5 format, with f5c version 1.4.⁸ M6anet version 2.1.0 was then applied with the HEK293T_RNA004 model to generate DRACH motif position predictions for human transcripts within the hg38 refMrna.fa.gz file (<http://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/>).²

Table 1: Sample Metadata for Samples

Sample Name	Condition
S007	MM1S-R1 SCRAMBLE
S008	MM1S-R1 gVIRMA 1
S009	MM1S-R1 gVIRMA 3
S010	MM1S-R2 SCRAMBLE
S011	MM1S-R2 gVIRMA 1
S012	MM1S-R2 gVIRMA 3
HCV_IVT_004	In-Vitro Transcribed HCV RNA

Caption: R1 and R2 refer to two biological replicates; SCRAMBLE, gVIRMA 1, and gVIRMA 3 refer to different CRISPR guides with either non-*VIRMA* targeting (SCRAMBLE) or *VIRMA*-specific targeting (gVIRMA 1 and gVIRMA 3). HCV_IVT_004 refers to in-vitro transcribed HCV RNA.

Identifying Epitranscriptome Modifications. The `data.indiv_proba.csv` file output by *m6anet* was transferred to an *R* version 4.3.1 environment, wrapped with *RStudio*. Custom scripts were written to generate the remaining analysis and graphs. Of note, the *diptest* package was utilized to calculate test statistics for the degree of unimodality within the distributions for sites with Hardigan's Dip Test.⁹

Results

Single-Molecule Distribution of M6aNET Predictions Resolves Limited Biological Variability

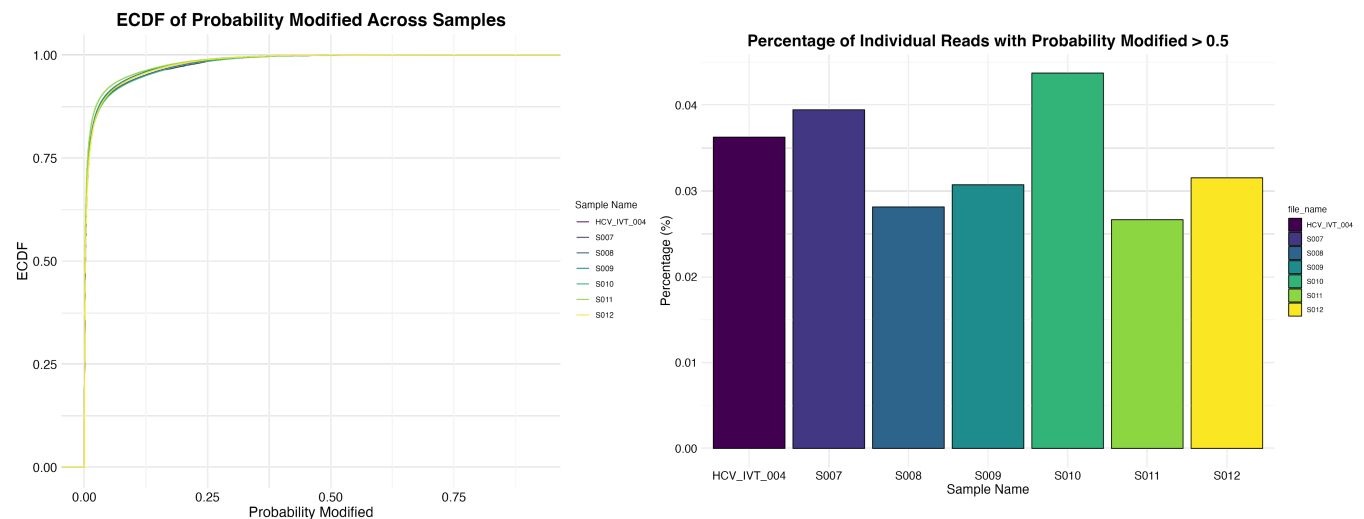
When examining the distribution of probabilities that a particular m6a site on an individual modified, I notice little difference between the *VIRMA*-knockdown (*VIRMA*-KD) conditions [S008, S009, S011, S012] and the wildtype conditions [S007, S010] (Figure 1). Notably, these distributions also did not exhibit any noticeable difference from the distribution of individual site probabilities from an HCV in-vitro transcribed sample, a system which lacks the capacity to form any modifications.

To resolve the capacity of *m6anet* to call high-confidence modifications, I selected an arbitrary cutoff of .5 to excise the modification calls of highest confidence, encompassing the top .05 percent of individual molecules. While the wildtype conditions did exhibit a slight .01 percentage increase when compared to the *VIRMA*-KD, the *VIRMA*-KD also demonstrated a

lower percentage of predicted reads than the HCV in-vitro transcribed sample. Although *VIRMA*-KD is expected to have lower levels of m6a, it should still demonstrate greater methylation probabilities than in-vitro transcribed RNA.

Together, these suggest that the native m6a probabilities output by *m6anet* are not providing meaningful resolution into the epitranscriptomics of these particular sample systems.

Figure 1: Distribution of Site Probabilities at Single-Molecule Resolution, predicted by *m6anet*



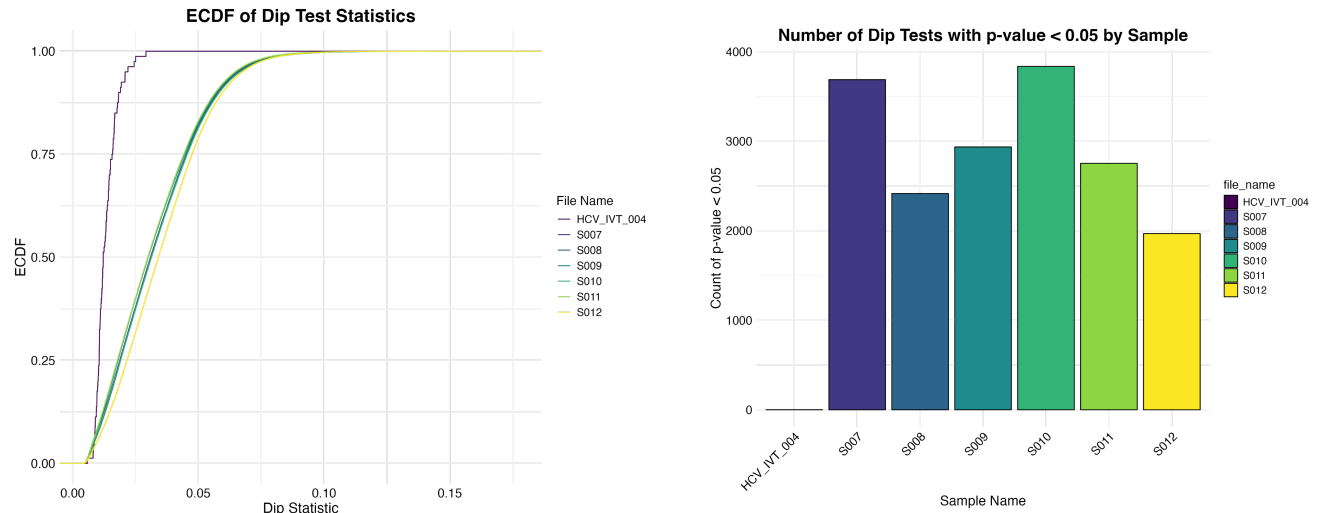
Caption: ECDF of the probability of modification for 1/20 of the individual molecules, randomly extracted from each sample (left), and a bar plot demonstrating the percentage of individual single-molecules with a high probability of modification (right)

Hardigan's Dip Test Resolves Potential Biological Variability in *VIRMA*-KO Dataset

On the assumption that *m6anet* would assign different methylation probabilities to individual molecules in different methylation states, I applied Hardigan's Dip test to test whether the individual distributions would be multi-modal.⁹ Hardigan's Dip test calculates a statistic, between zero [uni-modal] and one [multimodal], which quantifies the degree to which the cumulative distribution function of a multimodal function needs to be altered in order to become a unimodal distribution. I assume that sites with appreciable methylation will be bimodal, with separate peaks representing unmethylated and methylated states.

After calculating the dip test statistic for all sites within each sample, I noted a difference in the distribution between the biological samples and the in-vitro-transcribed RNA, suggesting the presence of epi-transcriptomic modifications was driving multimodal distributions (Figure 2). I additionally observed many sites with statistically significant multimodal distributions, with the wildtype samples having a greater number of multimodal distributions than *VIRMA*-KD conditions. These trends are affected by but are independent from the amount of coverage within each sample (data not shown).

Figure 2: Distribution of Site Modality Test Statistics, calculated with Hardigan's Dip test



Caption: ECDF of the dip statistic for individual sites (left) and a bar plot demonstrating the number of sites with a statistically significant dip statistic (right)

To understand whether these bi-modal distributions were related to each other in potentially biological relevant conditions, I examined the intersections between the statistically significant dip tests for S007, S008, S009, S010, S011, and S012.

Each of the samples had a large number of sites unique to that specific sample, suggesting potential individualized m6a dynamics. The largest overlap set was between the two wild-type sets, potentially implicating sites of methylation affected by *VIRMA*-KD (Figure 3). The second largest set of overlap consisted of all six samples, implying shared sites of m6a methylation.

The figure displays the distribution of intersection sizes for sets S007 through S012. The top part is a bar chart showing the intersection size for each set, and the bottom part is a dot plot showing the distribution of set sizes for each intersection size.

Intersection Size Data:

Set	Intersection Size
S007	1508
S008	1300
S009	952
S010	875
S011	851
S012	637
S007	526
S008	301
S009	167
S010	147
S011	142
S012	127
S007	121
S008	117
S009	115
S010	100
S011	92
S012	91
S007	88
S008	87
S009	82
S010	79
S011	71
S012	67
S007	64
S008	57
S009	56
S010	53
S011	52
S012	49
S007	48
S008	47
S009	47
S010	38
S011	37
S012	35
S007	33
S008	33
S009	32
S010	32

Set Size Distribution:

Intersection Size	S007	S008	S009	S010	S011	S012
1508	1	0	0	0	0	0
1300	0	1	0	0	0	0
952	0	0	1	0	0	0
875	0	0	0	1	0	0
851	0	0	0	0	1	0
637	0	0	0	0	0	1
526	1	0	0	0	0	0
301	0	1	0	0	0	0
167	0	0	1	0	0	0
147	0	0	0	1	0	0
142	0	0	0	0	1	0
127	0	0	0	0	0	1
121	1	0	0	0	0	0
117	0	1	0	0	0	0
115	0	0	1	0	0	0
100	0	0	0	1	0	0
92	0	0	0	0	1	0
91	0	0	0	0	0	1
88	1	0	0	0	0	0
87	0	1	0	0	0	0
82	0	0	1	0	0	0
79	0	0	0	1	0	0
71	0	0	0	0	1	0
67	0	0	0	0	0	1
64	0	0	0	0	0	0
57	0	0	1	0	0	0
56	0	0	0	1	0	0
53	0	0	0	0	1	0
52	0	0	0	0	0	1
49	0	0	0	0	0	0
48	0	0	0	0	0	0
47	0	0	0	0	0	0
38	0	0	0	0	0	0
37	0	0	0	0	0	0
35	0	0	0	0	0	0
33	0	0	0	0	0	0
32	0	0	0	0	0	0

Future work will quantify the multimodal peaks in order to allow for quantitative comparisons across samples, allowing for a determination of potential differences within the shared data set between *VIRMA*-KD and wildtype samples.

References

1. Garalde, D. R. *et al.* Highly parallel direct RNA sequencing on an array of nanopores. *Nat. Methods* **15**, 201–206 (2018).
2. Hendra, C. *et al.* Detection of m6A from direct RNA sequencing using a multiple instance learning framework. *Nat. Methods* **19**, 1590–1598 (2022).
3. Zhong, Z.-D. *et al.* Systematic comparison of tools used for m6A mapping from nanopore direct RNA sequencing. *Nat. Commun.* **14**, 1906 (2023).
4. Pratanwanich, P. N. *et al.* Identification of differential RNA modifications from nanopore direct RNA sequencing with xPore. *Nat. Biotechnol.* **39**, 1394–1402 (2021).
5. Yue, Y. *et al.* VIRMA mediates preferential m6A mRNA methylation in 3'UTR and near stop codon and associates with alternative polyadenylation. *Cell Discov.* **4**, 10 (2018).
6. Barros-Silva, D. *et al.* VIRMA-Dependent N6-Methyladenosine Modifications Regulate the Expression of Long Non-Coding RNAs CCAT1 and CCAT2 in Prostate Cancer. *Cancers* **12**, 771 (2020).
7. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
8. Gamaarachchi, H. *et al.* GPU accelerated adaptive banded event alignment for rapid comparative nanopore signal analysis. *BMC Bioinformatics* **21**, 343 (2020).
9. Hartigan, J. A. & Hartigan, P. M. The Dip Test of Unimodality. *Ann. Stat.* **13**, (1985).