

**Training Data Diversity Enhances the Basecalling of Novel RNA
Modification-Induced Nanopore Sequencing Readouts**

Supplementary Information

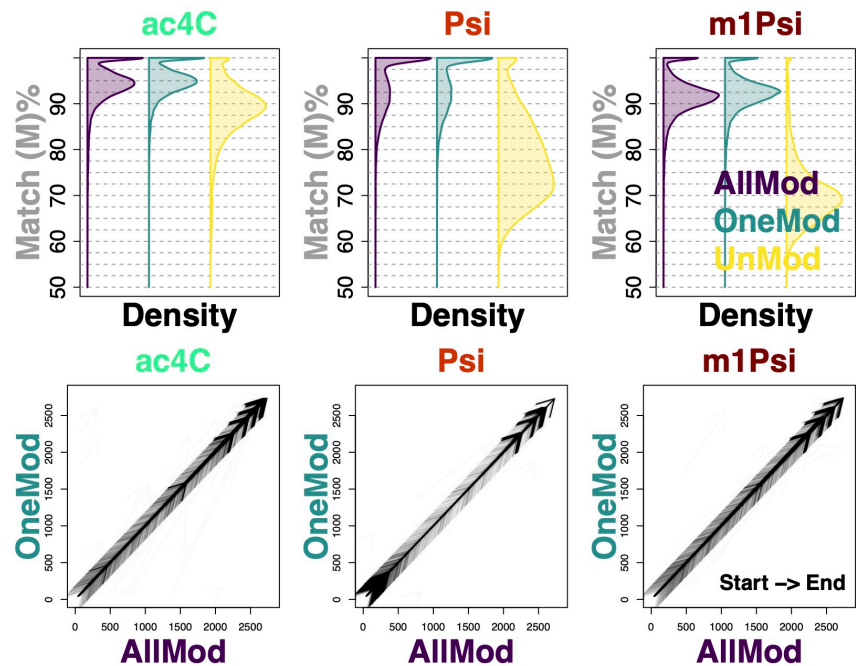
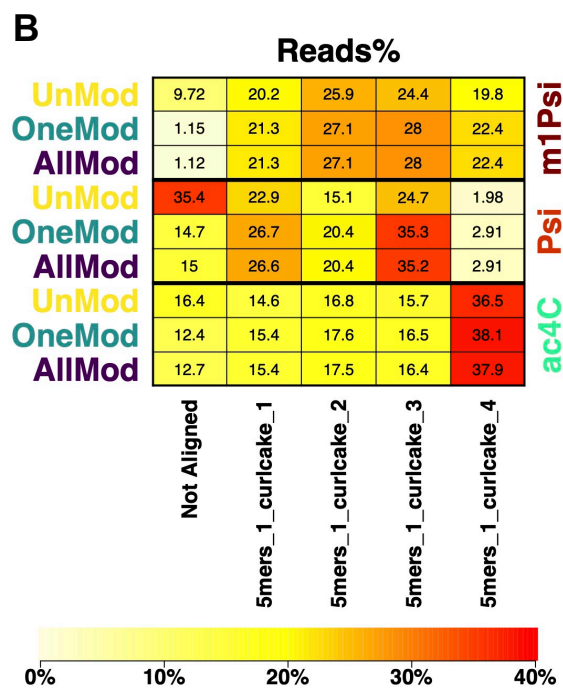
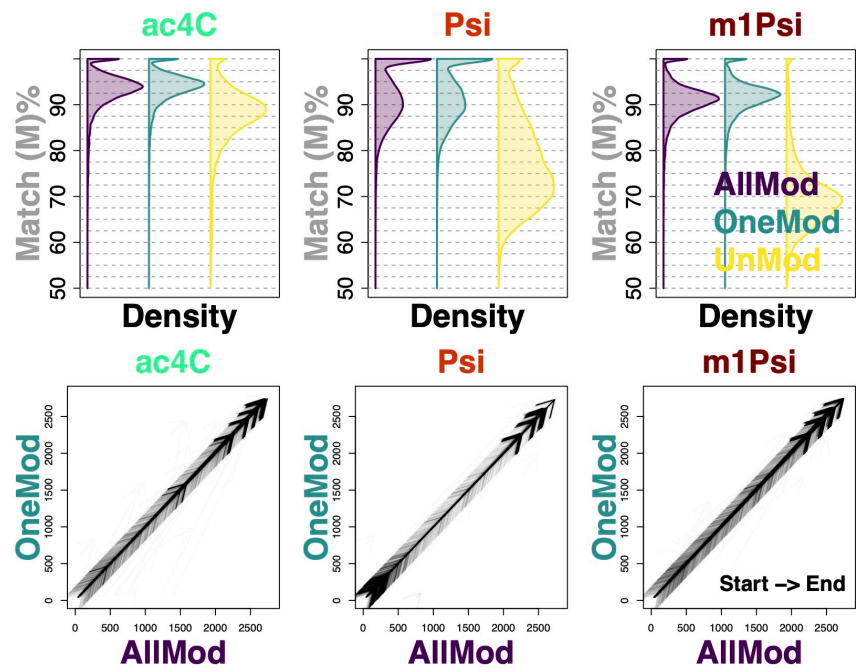
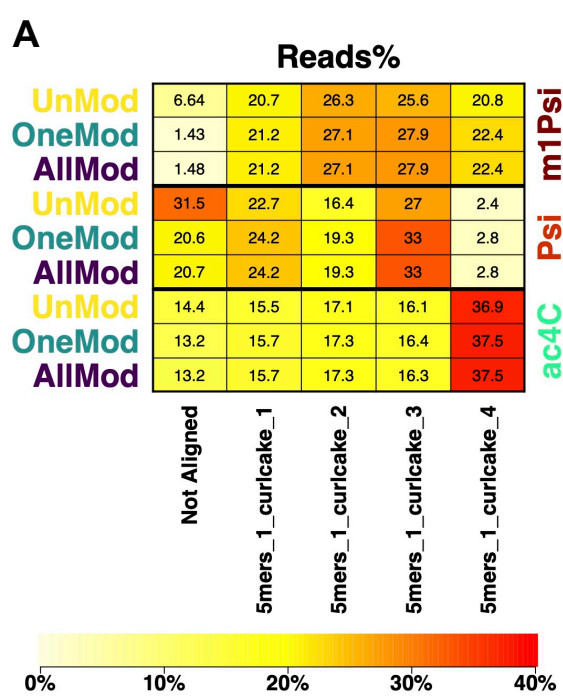
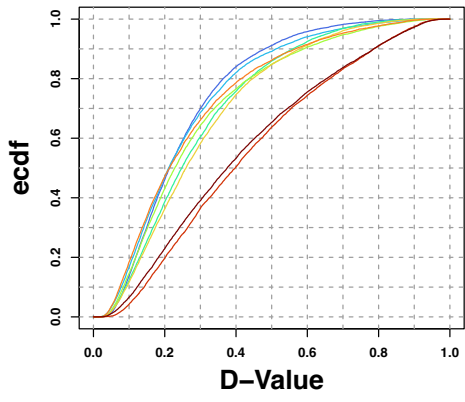
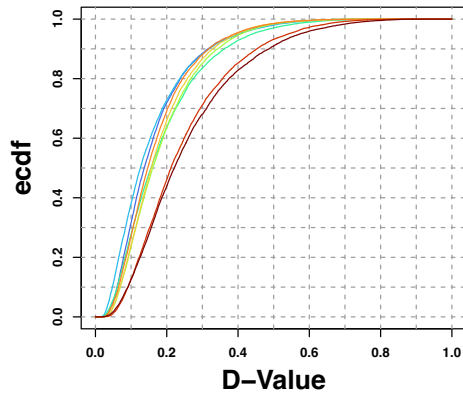


Figure S1. Bonito and Dorado trained with diverse known modifications gain the capability to basecall novel modifications. The basecalling performance of Bonito and Dorado were presented in (A) and (B), respectively. Same as in Figure 2B-D, the mappability, per-read CIGAR match fraction and alignment consistency were used as quantification metrics. AllMod, the basecaller trained by all the modifications except for the one to be basecalled; OneMod, the basecaller trained with only the modification to be basecalled; UnMod, the basecaller trained by only unmodified reads.

Mean



SD



Dwell

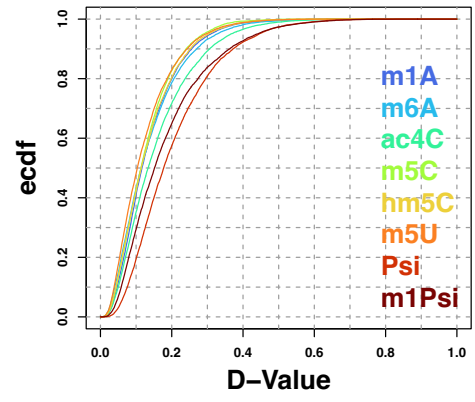


Figure S2. Modifications deviate nanopore sequencing signal features. Mean, SD and Dwell denote the mean, standard deviation and dwell time of Remora signal events. Acronyms denote different modification oligo categories. The KS-test D-values measure Mean, SD and Dwell distribution differences between modified and unmodified oligos. Ecdf denotes the empirical cumulative distribution function.

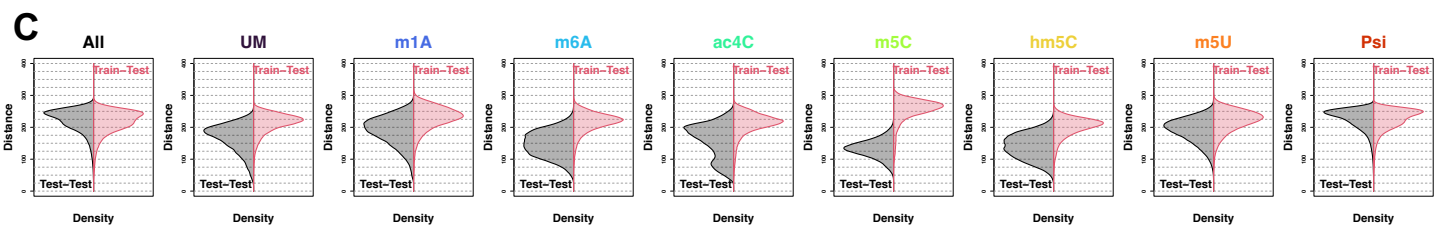
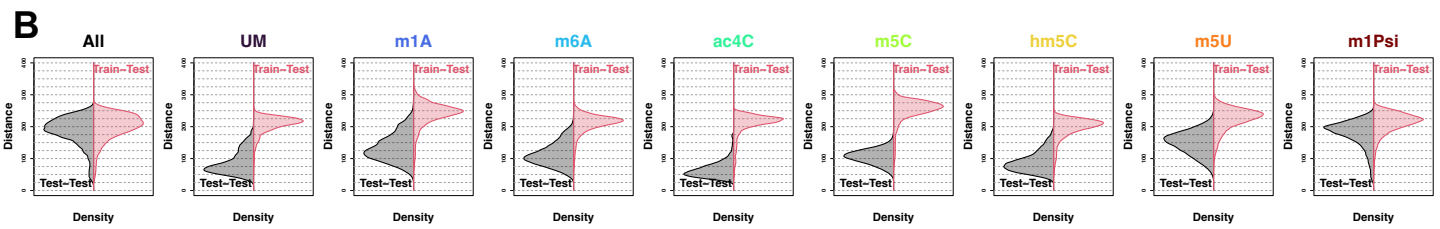
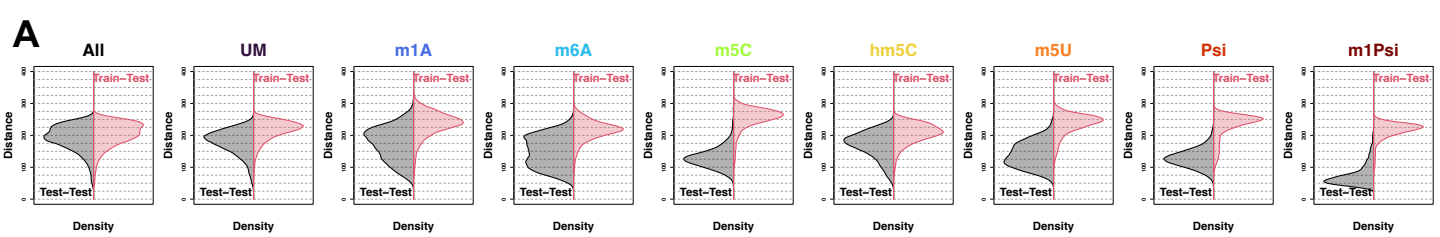


Figure S3. Quantifying the read-level similarity in the representation space. Reads were projected into the representation space and then subjected to PCA. Within the PC space, Euclidean distance was calculated to quantify the read-level similarity. Analyses for ac4C, Psi and m1Psi test reads were presented in (A), (B) and (C), respectively. All, the jointly-trained basecaller by all the oligo groups except for the test; other acronyms denote individually-trained basecallers. Train-Test and Test-Test represent the distance between training and test reads and among test reads, respectively. We took Test-Test groups as baselines to quantify the Train-Test similarity.

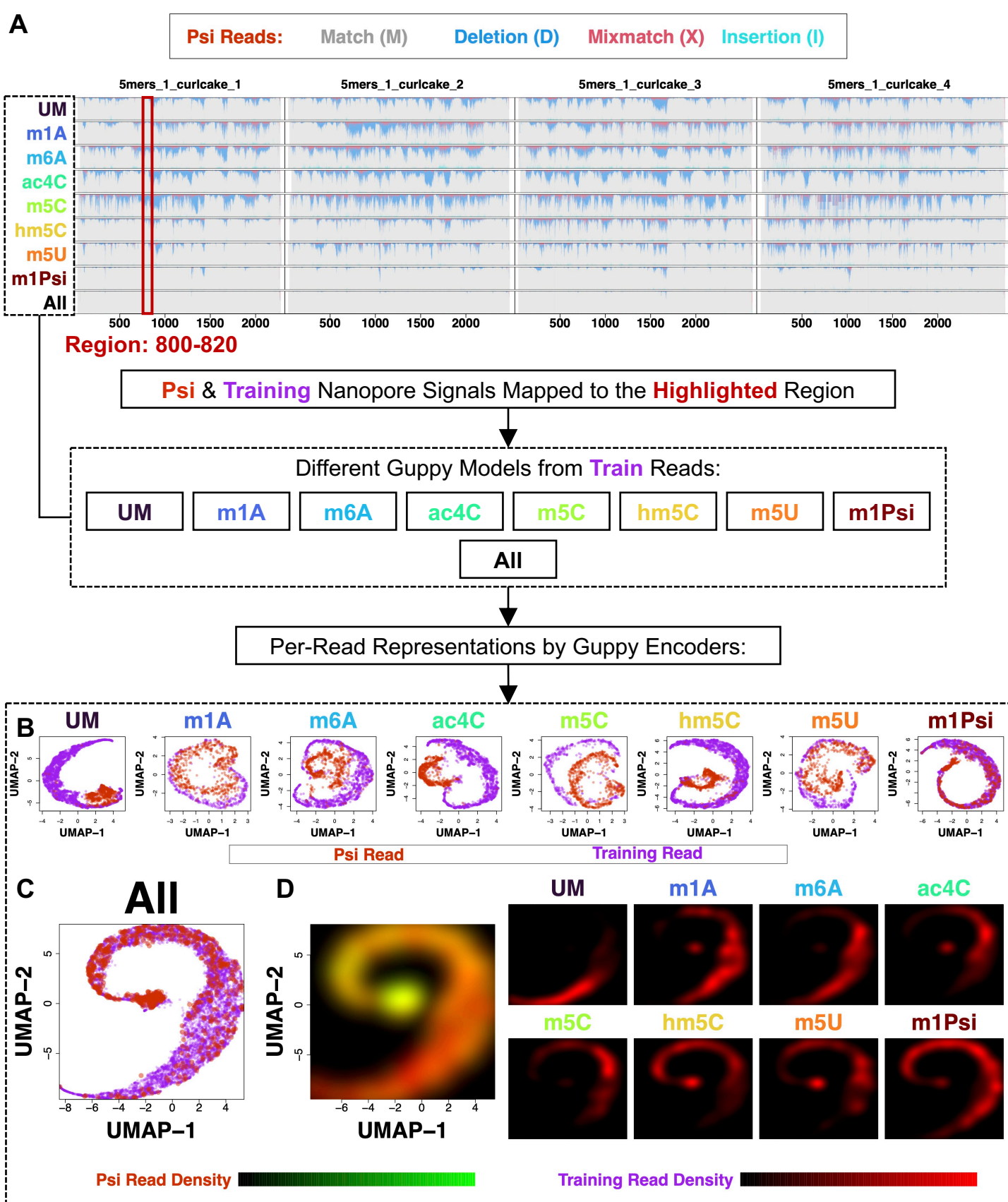


Figure S4. A generalizable representation space produced from diverse training modifications facilitates the basecalling of the out-of-sample Psi modification. (A) Performance of individually and jointly-trained basecallers on Psi reads was visualized with the genome viewer graph, which shows per-nucleotide CIGAR fractions. All, the jointly-trained basecaller by all the oligo types except for Psi; other acronyms denote individually-trained basecallers. For individually (B) and jointly-trained (C) basecallers, read fragments mapped to the boxed region were first converted as representation vectors with different basecaller encoders, then visualized by a UMAP plot. Train denotes reads used for training the corresponding basecaller. (D) Spatial distributions of different oligo types in the UMAP space as shown in (C). Black-to-green and red palette denotes Psi and training reads, respectively.

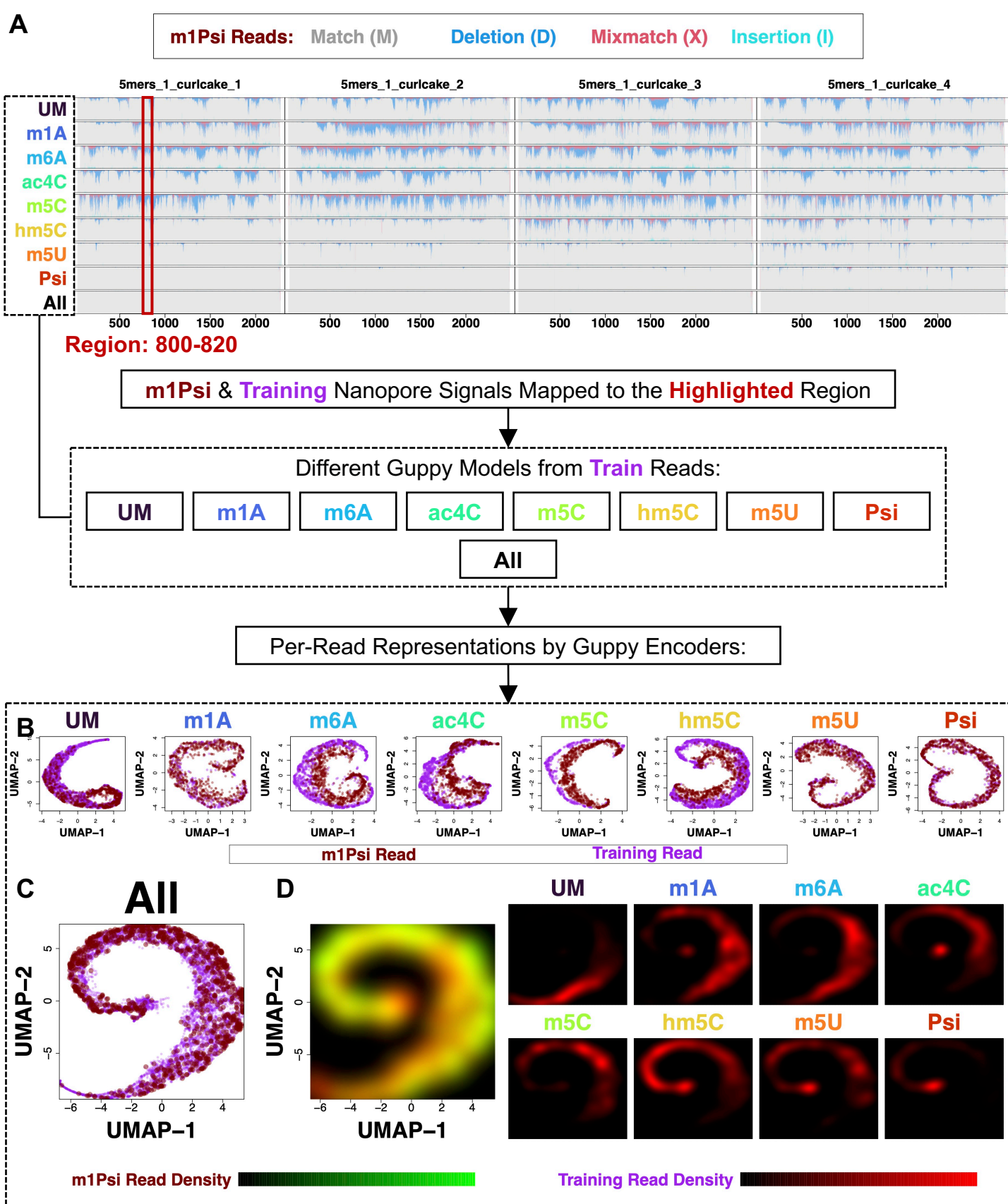
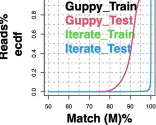


Figure S5. A generalizable representation space produced from diverse training modifications facilitates the basecalling of the out-of-sample m1Psi modification.

(A) Performance of individually and jointly-trained basecallers on m1Psi reads was visualized with the genome viewer graph, which shows per-nucleotide CIGAR fractions. All, the jointly-trained basecaller by all the oligo types except for m1Psi; other acronyms denote individually-trained basecallers. For individually (B) and jointly-trained (C) basecallers, read fragments mapped to the boxed region were first converted as representation vectors with different basecaller encoders, then visualized by a UMAP plot. Train denotes reads used for training the corresponding basecaller. (D) Spatial distributions of different oligo types in the UMAP space as shown in (C). Black-to-green and red palette denote m1Psi and training reads, respectively.

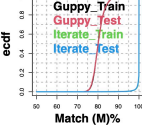
UM

Smers_1_curlcable.4	21.4	21.3	23.3	23.5
Smers_1_curlcable.3	19.4	19.2	21.2	21.2
Smers_1_curlcable.2	19.8	20.2	21.2	21.9
Smers_1_curlcable.1	23.8	23.1	25.2	26.1
Not Aligned	15.5	15.8	8.42	8.94



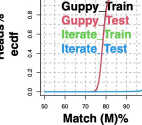
m1A

Smers_1_curlcable.4	3	2.68	13.7	14
Smers_1_curlcable.3	4.3	4.34	20.2	19.9
Smers_1_curlcable.2	2.6	2.6	14.3	15
Smers_1_curlcable.1	2.7	3.72	27.8	28.8
Not Aligned	0.4	0.74	23.3	23.3



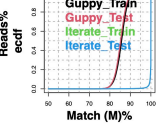
ac4C

Smers_1_curlcable.4	8.35	6.41	38.3	37.6
Smers_1_curlcable.3	7.27	5.39	17.3	16.4
Smers_1_curlcable.2	8.59	8.27	17.2	17.2
Smers_1_curlcable.1	6.13	4.27	15.2	15.2
Not Aligned	60.7	27.4	11	13.1



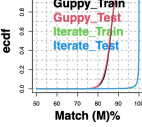
m5C

Smers_1_curlcable.4	13.3	13	27.3	27.6
Smers_1_curlcable.3	9.3	8.4	18.9	18.1
Smers_1_curlcable.2	9.09	8.58	17.2	17.9
Smers_1_curlcable.1	11	12	21	21.7
Not Aligned	95.3	47	15.7	14.8



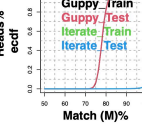
hm5C

Smers_1_curlcable.4	7.32	7.54	12.7	12.7
Smers_1_curlcable.3	3.45	3.83	5.88	6.24
Smers_1_curlcable.2	19.3	18.2	20.2	20.2
Smers_1_curlcable.1	28.4	24.4	38.4	38.1
Not Aligned	45.6	45.8	18.2	21.2



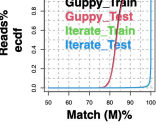
m6A

Smers_1_curlcable.4	6.28	6.3	20.4	20.4
Smers_1_curlcable.3	9.49	9.6	22.8	22.7
Smers_1_curlcable.2	5.99	6.01	18.2	18.2
Smers_1_curlcable.1	5.92	5.75	24.4	24.1
Not Aligned	75.3	72.4	13.8	13.8



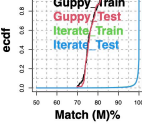
m5U

Smers_1_curlcable.4	7.33	10.4	16.8	20.9
Smers_1_curlcable.3	9.96	14.2	19.4	23.1
Smers_1_curlcable.2	13.1	21	24.4	30
Smers_1_curlcable.1	16.2	3.24	38.8	7.61
Not Aligned	65.2	62.2	11.3	10.8



Psi

Smers_1_curlcable.4	0	0.00564	2.68	2.79
Smers_1_curlcable.3	0.04	0.0272	31.7	33.4
Smers_1_curlcable.2	0.205	0.208	16.8	18.7
Smers_1_curlcable.1	0.355	0.365	22.1	24.3
Not Aligned	0.4	0.42	7.7	20.8



m1Psi

Smers_1_curlcable.4	0.445	0.465	22.8	22.4
Smers_1_curlcable.3	1.9	2.23	28.2	27.9
Smers_1_curlcable.2	2.17	2.24	27.2	27.1
Smers_1_curlcable.1	7.86	7.98	31.1	31.1
Not Aligned	0.4	0.4	1.53	1.61

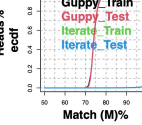


Figure S6. The iterative basecalling of RNA oligos. Guppy and iterate denote Guppy and 4th iteration (final iteration) basecallers, respectively. Train and test denote training and the independent test datasets, respectively. UM denotes unmodified reads. Match denotes the CIGAR category M. Ecdf, empirical cumulative distribution function. Start and End denote the alignment direction.