

UC Santa Cruz
Graduate Research Symposium 2016

Title

Cytosine Methylation Variant Calling with MinION Nanopore Sequencing

Permalink

<https://escholarship.org/uc/item/4669g90b>

Author

Rand, Arthur

Publication Date

2016-05-17

Peer reviewed

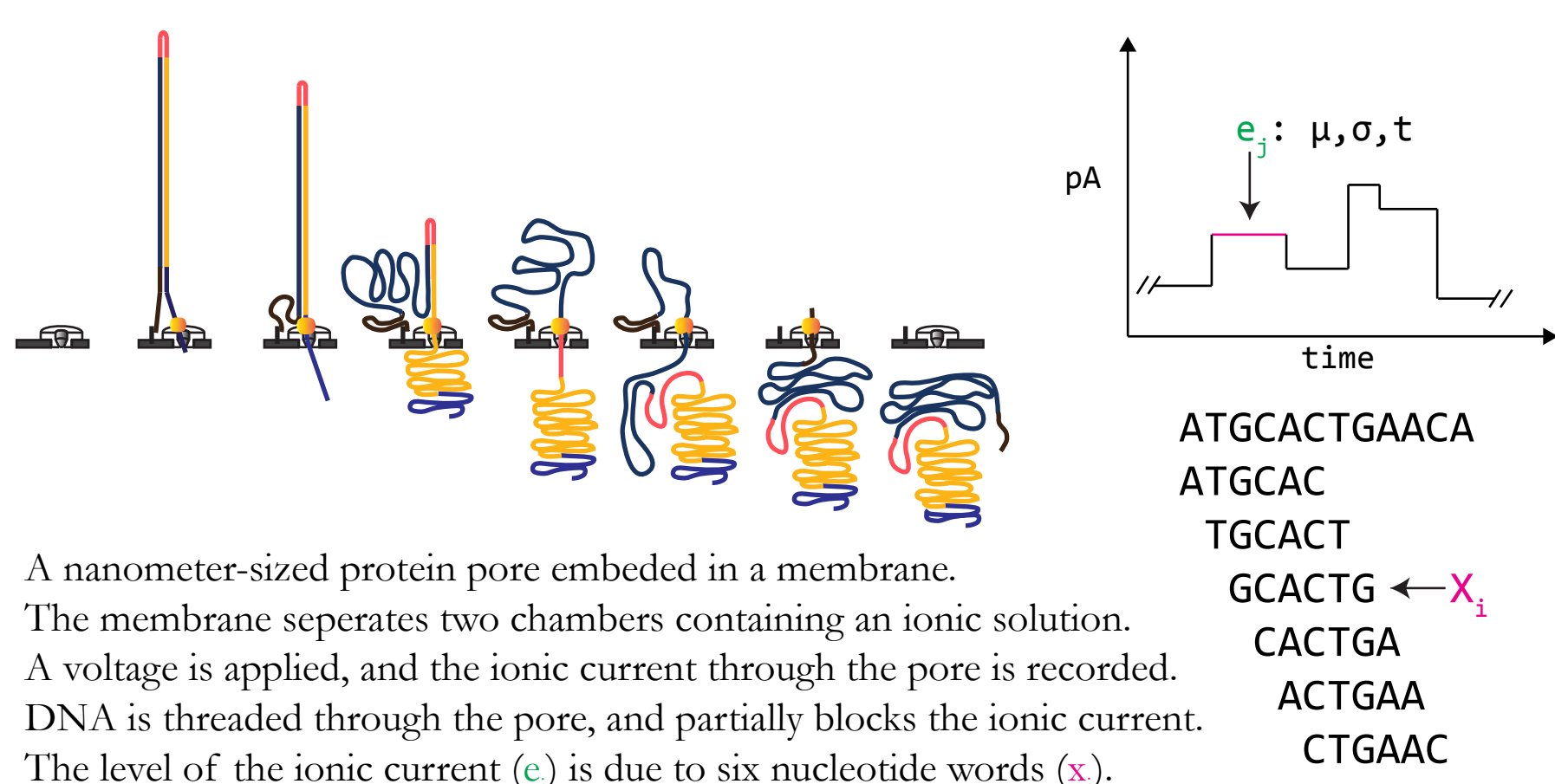
Cytosine Methylation Variant Calling with MinION Nanopore Sequencing

Arthur C. Rand, Miten Jain, Jordan Eizenga, Audrey Musselman-Brown, Hugh E. Olsen, Mark Akeson and Benedict Paten
Department of Biomolecular Engineering, University of California, Santa Cruz

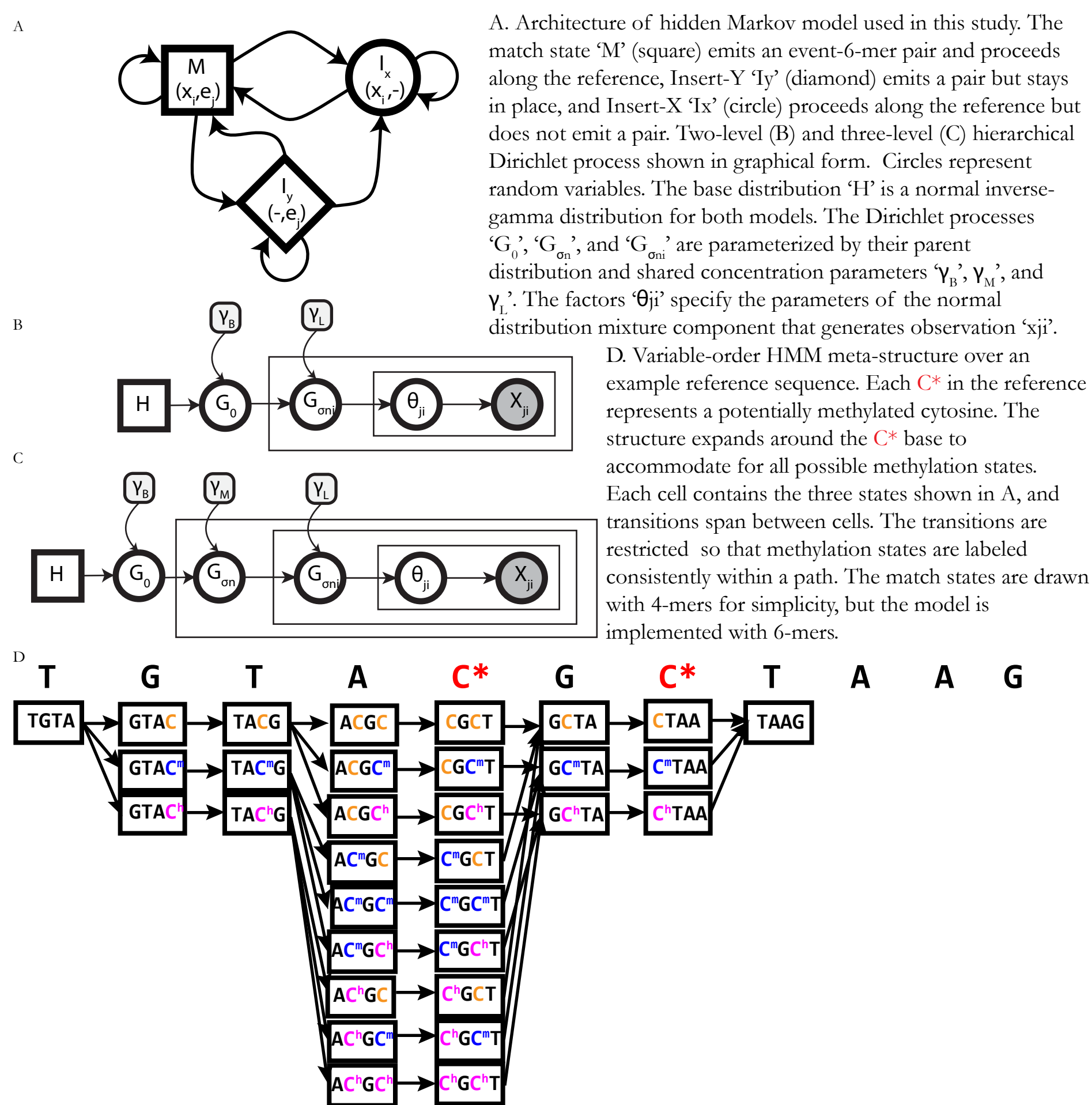
Abstract

Chemical modifications to DNA regulate cellular state and function. The Oxford Nanopore MinION is a portable single-molecule DNA sequencer that can sequence long fragments of genomic DNA. Here we show that the MinION can be used to detect and map three cytosine variants: cytosine, 5-methylcytosine, and 5-hydroxymethylcytosine. We present a probabilistic method that enables expansion of the nucleotide alphabet to include bases containing chemical modifications. Our results on synthetic DNA show that individual cytosine base modifications can be classified with accuracy up to 95% in a three-way comparison and 98% in a two-way comparison. We also demonstrate that 5-methylcytosine can be accurately mapped in *E. coli* genomic DNA

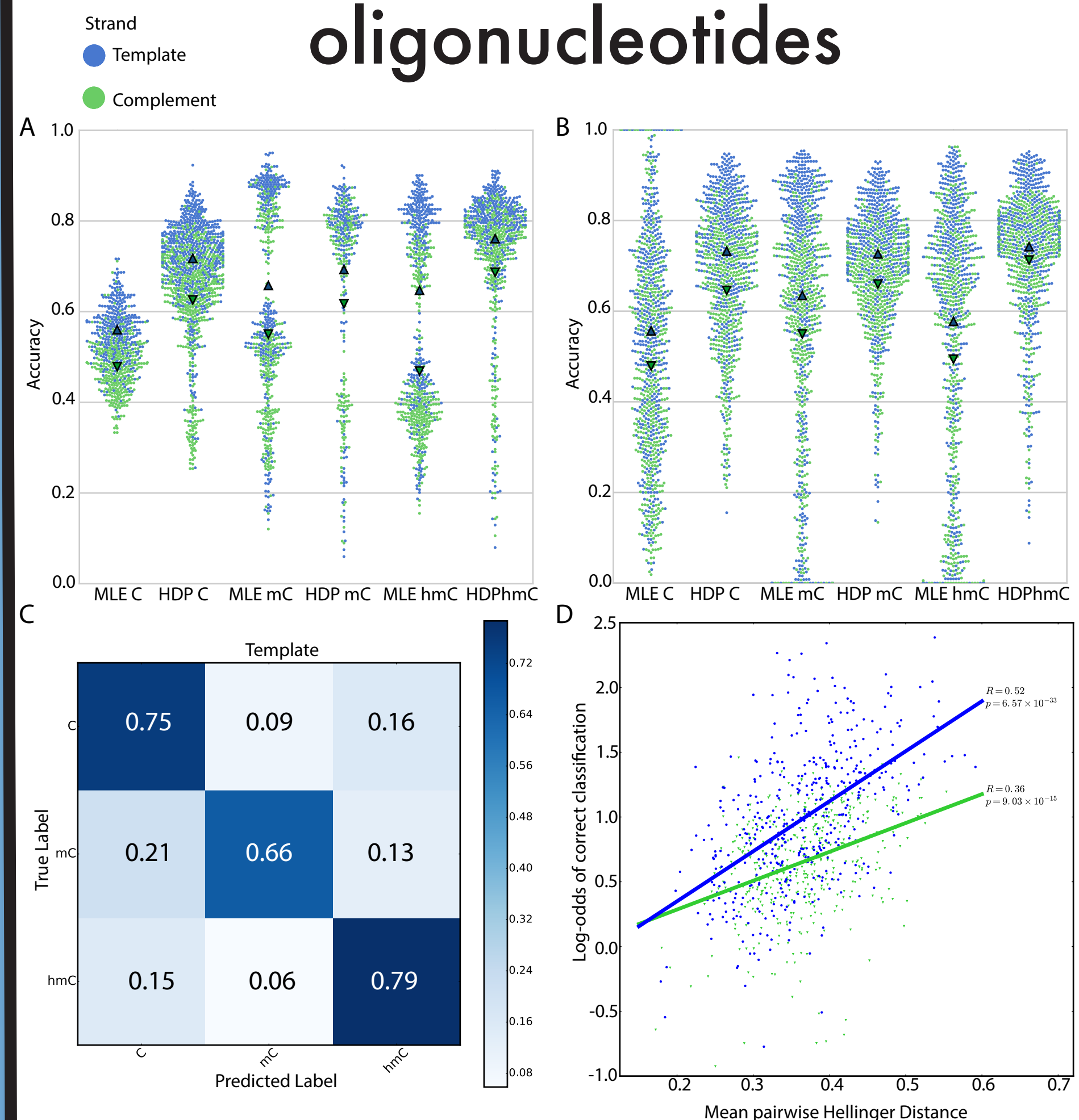
Nanopore Sequencing



Modeling Ionic Current with a hidden Markov model

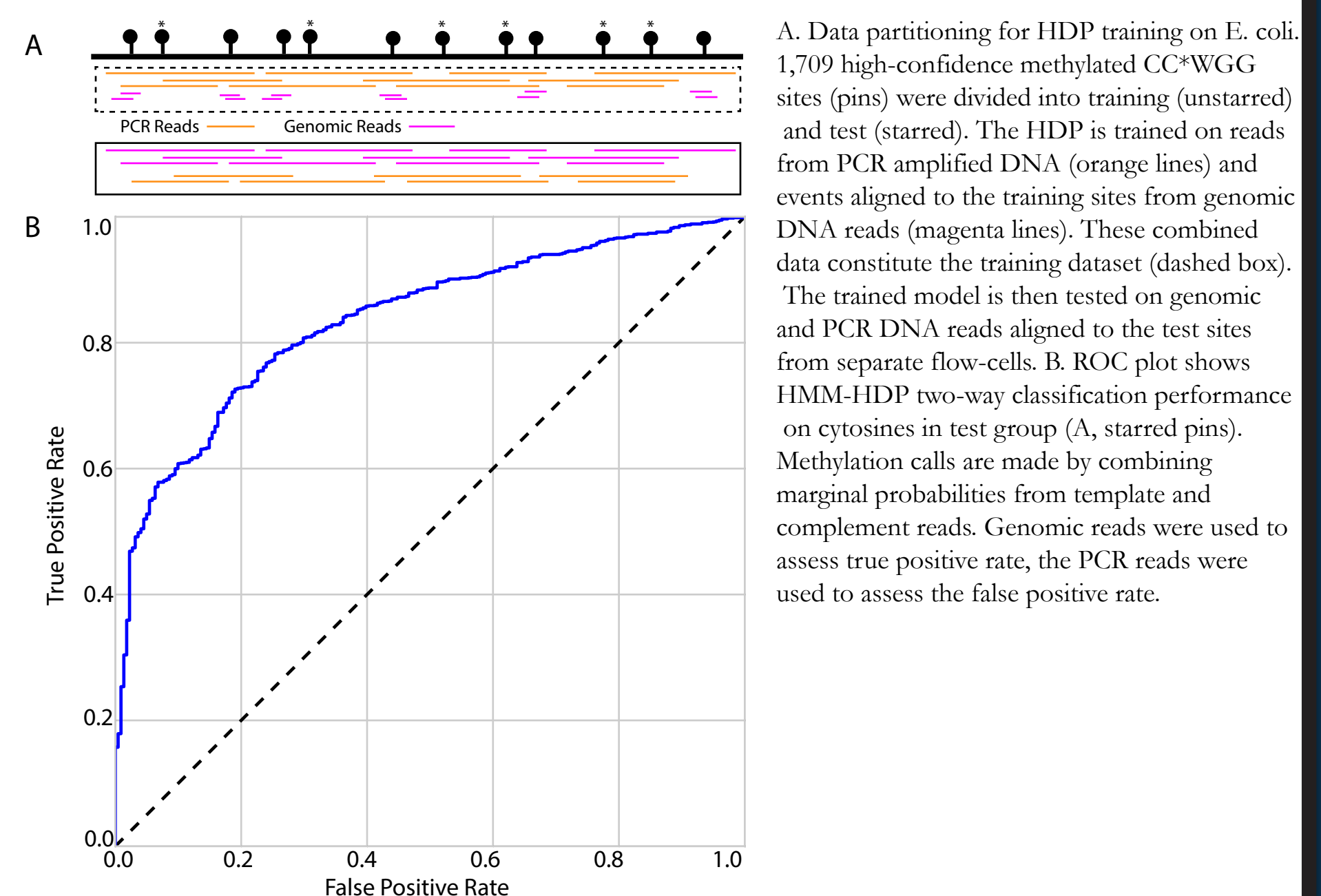


Base modification calling accuracy results on synthetic oligonucleotides

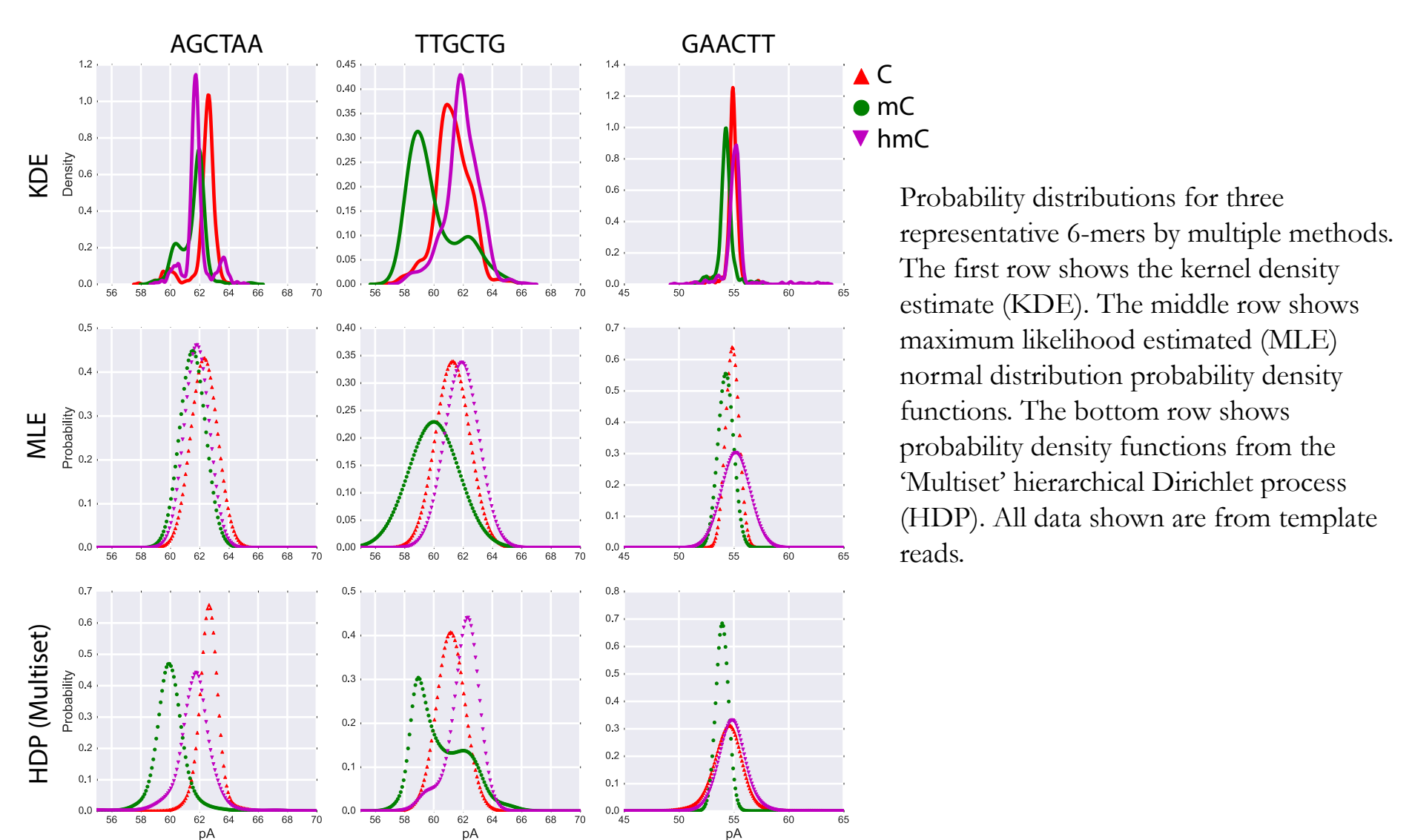


A and B. The accuracy distribution by read (A) and by context (B) is shown for the MLE emission distributions and the 'Multiset' HDP model on synthetic oligonucleotides. The triangles represent the mean of the distribution. C. Confusion matrix showing HMM-HDP three-way cytosine classification performance on template reads of synthetic oligonucleotides. D. Scatter plot shows the correlation between log-odds of correct classification and the mean pairwise Hellinger distance between the methylation statuses of the 6-mer distributions overlapping a cytosine.

Mapping 5-methylcytosine in *E. coli* genomic DNA



The HDP more realistically models ionic current distributions



Comparison of different HDP topologies

Model	Three-Way Accuracy			
	Mean Accuracy (read)	Median Accuracy (read)	Mean Accuracy (site)	Median Accuracy (site)
MLE	62% / 50%	58% / 47%	59% / 51%	66% / 57%
singlelevel	74% / 66%	79% / 72%	73% / 63%	76% / 69%
multiset	74% / 67%	80% / 76%	73% / 67%	76% / 70%
composition	73% / 66%	78% / 71%	73% / 66%	76% / 69%
middleNts	72% / 64%	76% / 69%	72% / 64%	75% / 67%
group	73% / 65%	78% / 71%	72% / 66%	75% / 69%
Model	Two-Way Accuracy			
	Mean Accuracy (read)	Median Accuracy (read)	Mean Accuracy (site)	Median Accuracy (site)
singlelevel	83% / 78%	86.5% / 84.5%	82% / 77%	83% / 78%
multiset	83% / 78%	86.5% / 84.5%	82% / 77%	83% / 78%

MLE is the maximum likelihood estimate of a normal distribution. 'Two-level' is an HDP model with no subgroupings of 6-mers, 'Multiset', 'Composition', 'MiddleNucleotides', and 'GroupMultiset' are three-level HDP models. Three-way classification was performed between cytosine, 5-methylcytosine, and 5-hydroxymethylcytosine. Two-way classifications were between cytosine and 5-methylcytosine.