



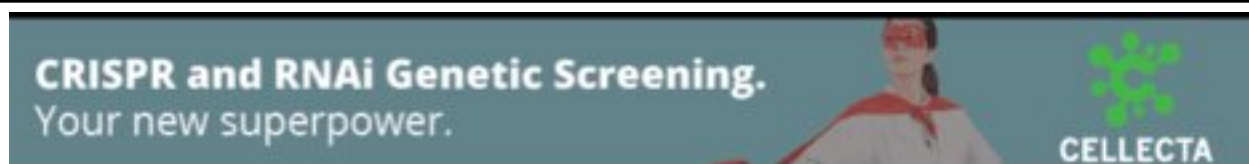
DNA-m6A calling and integrated long-read epigenetic and genetic analysis with fibertools

Anupama Jha, Stephanie C. Bohaczuk, Yizi Mao, et al.

Genome Res. published online June 7, 2024

Access the most recent version at doi:[10.1101/gr.279095.124](https://doi.org/10.1101/gr.279095.124)

P<P	Published online June 7, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 DNA-m6A calling and integrated long-read 2 epigenetic and genetic analysis with fibertools

3 Anupama Jha^{1,*}, Stephanie C. Bohaczuk^{2,*}, Yizi Mao², Jane Ranchalis², Benjamin J. Mallory¹,
4 Alan T. Min³, Morgan O. Hamm¹, Elliott Swanson¹, Danilo Dubocanin⁴, Connor Finkbeiner¹,
5 Tony Li¹, Dale Whittington⁵, William Stafford Noble^{1,6}, Andrew B. Stergachis^{1,2,7,†}, Mitchell R.
6 Vollger^{2,†}

7
8 1. Department of Genome Sciences, University of Washington, Seattle, WA, USA

9 2. Division of Medical Genetics, University of Washington School of Medicine, Seattle, WA, USA

10 3. Department of Statistics, University of Washington, Seattle, WA, USA

11 4. Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

12 5. Department of Medical Chemistry, University of Washington, Seattle, WA, USA

13 6. Paul G. Allen School of Computer Science and Engineering, University of Washington,
14 Seattle, WA, USA

15 7. Brotman Baty Institute for Precision Medicine, Seattle, WA, USA

16

17 † Corresponding author(s). absterga@uw.edu; mvollger@uw.edu

18 * These authors contributed equally to this work.

19 Abstract

20 Long-read DNA sequencing has recently emerged as a powerful tool for studying both genetic
21 and epigenetic architectures at single-molecule and single-nucleotide resolution. Long-read
22 epigenetic studies encompass both the direct identification of native cytosine methylation as
23 well as the identification of exogenously placed DNA *N*⁶-methyladenine (DNA-m6A). However,
24 detecting DNA-m6A modifications using single-molecule sequencing, as well as co-processing
25 single-molecule genetic and epigenetic architectures, is limited by computational demands and
26 a lack of supporting tools. Here, we introduce *fibertools*, a state-of-the-art toolkit that features a
27 semi-supervised convolutional neural network for fast and accurate identification of m6A-
28 marked bases using PacBio single-molecule long-read sequencing, as well as the co-
29 processing of long-read genetic and epigenetic data produced using either PacBio or Oxford

30 Nanopore sequencing platforms. We demonstrate accurate DNA-m6A identification (>90%
31 precision and recall) along >20 kilobase long DNA molecules with a ~1,000-fold improvement in
32 speed. In addition, we demonstrate that *fibertools* can readily integrate genetic and epigenetic
33 data at single-molecule resolution, including the seamless conversion between molecular and
34 reference coordinate systems, allowing for accurate genetic and epigenetic analyses of long-
35 read data within structurally and somatically variable genomic regions.

36 Introduction

37 Highly accurate long-read single-molecule DNA sequencing has revolutionized the
38 comprehensive assembly of phased genetic architectures, enabling the first complete human
39 genome assemblies (Wenger et al. 2019; Vollger et al. 2020; Nurk et al. 2022). In addition,
40 single-molecule long-read DNA sequencing natively identifies endogenously modified DNA
41 bases, such as m6A and 5-methylcytosine (5mC), permitting the co-analysis of both genetic and
42 DNA methylation features at single-molecule resolution (Marks et al. 2012; Clark et al. 2012;
43 Murray et al. 2012; Loman et al. 2015; Töpfer and Wenger 2023). Furthermore, using
44 exogenous DNA methyltransferases to add DNA base modifications, such as in the context of
45 single-molecule chromatin fiber sequencing (Stergachis et al. 2020; Lee et al. 2020; Abdulhay et
46 al. 2020; Shipony et al. 2020; Altemose et al. 2022; Cheetham et al. 2022), permits the co-
47 analysis of genetic, DNA methylation, and chromatin epigenetic features at single-molecule and
48 single-nucleotide resolution.

49
50 Specifically, single-molecule chromatin fiber sequencing leverages non-specific
51 methyltransferases to selectively stencil chromatin protein occupancy patterns directly onto their
52 underlying DNA molecules in the form of modified bases. Modified bases along individual DNA
53 molecules are then directly identified using PCR-free single-molecule sequencing. For example,
54 during single-molecule, real-time (SMRT) sequencing, the identity of each base is determined
55 by the fluorophore-labeled nucleotide that is incorporated as the polymerase replicates the
56 base. In contrast, the modification status of each base is determined by signature changes in
57 polymerase kinetics at and surrounding that base as it is replicated by the polymerase, such as
58 elongation of the interpulse duration (IPD) due to polymerase pausing at modified bases
59 (Flusberg et al. 2010) (**Fig. 1A, Supplemental Fig. S1**). Recently developed tools leverage
60 these polymerase kinetic parameters to identify 5mC within specific sequence contexts (Tse et
61 al. 2021; Töpfer and Wenger 2023), genomic positions with consistent m6A signal across

62 multiple sequencing reads (Marks et al. 2012; Clark et al. 2012; Murray et al. 2012), total
63 adenine methylation levels along very short sequencing reads (Kong et al. 2022), and m6A-
64 modified bases at single-molecule resolution along only short ~2 kilobase DNA molecules
65 (ipdSummary, SAMOSA-ChAAT) (Abdulhay et al. 2023; Clark et al. 2012). However, accurately
66 identifying DNA-m6A along multi-kilobase reads is largely unsolved, as existing approaches
67 either have poor sensitivity/specificity (**Supplemental Figs. S2, S3**), require excessive compute
68 and storage resources (**Supplemental Fig. S4**), and/or are reliant on input files (i.e., subreads)
69 no longer available with modern sequencing chemistries.

70
71 Furthermore, although extensive tooling exists for processing short-read epigenetic data relative
72 to reference coordinates (BEDOPS, BEDTools, etc.) (Neph et al. 2012; Quinlan 2014),
73 comparable tools for processing long-read epigenetic data are limited in their ability to leverage
74 the rich epigenetic and genetic data embedded within long-read sequencing data (Razaghi et al.
75 2022). Specifically, tools for processing long-read epigenetic data need to operate in four
76 dimensions: (1) they must process multiple types of genetic and epigenetic information present
77 on a single read (e.g., DNA base, mCpG, DNA-m6A, inferred epigenetic marks, etc.); (2) they
78 must capture this information across multiple reads mapping to a given reference position; (3)
79 they must capture how this information co-occurs along each read in both reference and
80 molecular coordinates (the original positions from the sequenced read, without adjustments for
81 alignment to the reference sequence) in order to accurately display and detect the impact of
82 structural or somatic variation on epigenetic marks; and (4) they must capture all of this
83 information across the various haplotypes mapping to the same position within a reference.

84
85 Here, we introduce a semi-supervised machine learning approach for accurately identifying
86 DNA-m6A in PacBio sequencing along multi-kilobase reads that permits the accurate learning of
87 modified DNA bases from noisy training data - a common occurrence with single-molecule
88 sequencing data owing to inherent biological heterogeneity in DNA methylation status between
89 individual DNA molecules. Furthermore, we introduce a comprehensive toolkit for co-processing
90 long-read genetic and epigenetic data designed for use across sequencing platforms.

91 Results

92 Building an accurate tool for DNA-m6A identification requires a training dataset of multi-kilobase
93 reads with both methylated and unmethylated adenines across diverse sequence contexts (i.e.,

94 all possible 7-mers containing a central adenine) and methylation density contexts (i.e., isolated
95 or clustered DNA-m6As). Because creating such a dataset is not achievable using synthetic
96 DNA or fully methylated and unmethylated samples, we leveraged DNA from single-molecule
97 chromatin fiber sequencing reactions (i.e., Fiber-seq) as the basis for training. Specifically,
98 Fiber-seq uses non-specific m6A-MTases to selectively mark sites of protein occupancy along
99 individual DNA molecules via m6A-marked bases. Because protein occupancy is highly
100 heterogeneous across chromatinized DNA (**Supplemental Fig. S5**), each DNA molecule
101 contains methylated adenines within diverse sequence and methylation density contexts
102 (**Supplemental Fig. S6**). Furthermore, we can employ chromatin features, such as nucleosome
103 occupancy, to bolster training and validation, making Fiber-seq well suited for training a general-
104 purpose DNA-m6A caller given labeled data.

105
106 To generate initial positive and negative labels, we used a previously published DNA-m6A caller
107 (referred to here as the “subread model”) (Dubocanin et al. 2022). The subread model improves
108 upon ipdSummary by using ipdSummary’s IPD normalization for sequence context (ipdRatios),
109 followed by a Gaussian mixture model (GMM) to identify adenines with ipdRatios that
110 significantly deviate from the expected distribution of unmethylated adenines (i.e., m6A-modified
111 bases) (Dubocanin et al. 2022) (**Supplemental Fig. S7**). These calls are then used to identify
112 m6A-modified bases (positive labels). Negative labels are drawn from regions with extended
113 stretches devoid of m6A corresponding to inferred nucleosome-occluded regions (**Fig. 1A**,
114 *Methods*).

115
116 All existing DNA-m6A callers (i.e., ipdSummary, SAMOSA-ChAAT, and the “subread model”)
117 require subreads. However, modern SMRT sequencing does not output subreads and only
118 produces summary kinetic information, making these existing tools largely obsolete.
119 Consequently, we designed *fibertools* using a two-staged training approach. In stage 1 of
120 training (**Fig. 1A,B**), we used a fully supervised training regime to validate that m6A calls can be
121 generated using only summary kinetics, bypassing the requirement of all other PacBio m6A
122 callers for individual subread kinetics. Using the dataset described above, we independently
123 trained two machine-learning models, XGBoost (Chen and Guestrin 2016) and a fully-
124 supervised convolutional neural network (fully-supervised CNN), and evaluated their
125 performance on a held-out dataset from a separate sequencing experiment (**Fig. 1A,B**, and
126 *Methods*). Both the fully-supervised CNN and XGBoost models maintained high precision and
127 recall, with average precision over 97% (**Supplemental Table S1**). As a comparison, we

128 benchmarked the models against ipdSummary (we excluded SAMOSA-ChAAT from
129 benchmarking as it was only optimized for a single deprecated polymerase chemistry). Both our
130 CNN and XGBoost models outperformed ipdSummary in terms of both AUPR and AUROC. At a
131 95% precision threshold, the fully-supervised CNN model has a recall of 89.6%, whereas the
132 recall for ipdSummary is only 47.4%. Thus, the fully-supervised CNN nearly doubles the number
133 of m6A identifications over ipdSummary at this precision threshold. Notably, the fully-supervised
134 CNN and XGBoost models are ~1,000× faster than the subread model used to generate the
135 training dataset. Overall, the CNN model had the best performance, and we used it as the basis
136 for stage 2 of training as well as subsequent improvements and validation.

137
138 In stage 2 of training, we sought to polish the architecture of our supervised CNN beyond the
139 capabilities and accuracy limitations of existing tools and training datasets using a semi-
140 supervised training regime inspired by well-established methods in the field of proteomics (Käll
141 et al. 2007; Fondrie and Noble 2021). To do this, we trained a new semi-supervised model
142 initiated with the architecture weights of our stage 1 model to overcome the imperfect training
143 labels derived from existing models. The resulting semi-supervised model (referred to onwards
144 as “*fibertools*”), allows for the possibility that the positive labels in the training dataset are
145 incorrect (**Fig. 1C**, *Methods*). Using evaluation metrics (e.g., AUPR/AUROC) based on
146 training/test labels in semi-supervised learning violates the inherent assumption that the labels
147 may be incorrect; therefore, to evaluate this model, we established a series of biological
148 validations to test the performance of the semi-supervised training.

149
150 First, we assessed the accuracy of *fibertools* for identifying nucleosome footprints along Fiber-
151 seq data. Using single-molecule m6A calls with a predicted precision of >95% (*Methods*), we
152 performed an autocorrelation analysis. Compared to the subread and fully supervised models,
153 *fibertools* more accurately recapitulated the exact length of nucleosomes (147 bp) (**Fig. 1D**)
154 (Luger et al. 1997) and showed an overall higher amplitude autocorrelation, consistent with
155 higher quality identification of m6A with nucleosome patterning characteristic of human
156 chromatin. In addition, comparing the distance between adjacent m6A methylation marks in
157 human Fiber-seq data demonstrated clear oscillatory patterns suggestive of nucleosome
158 breathing (Polach and Widom 1995; Anderson and Widom 2000; Hall et al. 2009), further
159 indicative of high-quality m6A identification (**Fig. 1E,G**).

160

161 Second, we evaluated the false positive rate (FPR) of *fibertools* using whole-genome amplified
162 (WGA) DNA that lacked m6A (**Fig. 2A**). Our findings indicated an FPR of 0.23% at the model-
163 predicted precision level of >95% (**Fig. 2B, Supplemental Fig. S8**). Notably, given this FPR,
164 *fibertools* is not suited for identifying non-specific genomic m6A events within species with low-
165 level endogenous m6A (Kong et al. 2022; Debo et al. 2023).

166
167 Third, we evaluated the ability of *fibertools* to accurately quantify the total amount of m6A within
168 a sample. Specifically, we spiked varying levels of m⁶ATP into a WGA reaction (**Fig. 2A**) and
169 employed ultra-high-performance liquid chromatography tandem mass spectrometry (UHPLC–
170 MS/MS) to determine the percentage of methylated adenines with respect to all adenines
171 (*Methods*). We then sequenced these samples, applied *fibertools*, and found a strong
172 correlation (Pearson=0.998, p-value=5.2e-6) between our method and mass spectrometry (**Fig.**
173 **2C**).

174
175 Fourth, since the above validations demonstrate that m6A calls from *fibertools* recapitulate bulk
176 chromatin features and total m6A content measured by UHPLC–MS/MS, we next evaluated the
177 precision of *fibertools* for identifying isolated m6A events, which is relevant for m6A calls within
178 small internucleosomal regions. Using *fibertools* to predict m6A on genomic DNA treated with
179 motif-specific methyltransferases, we found that m6A calls were enriched by 415-fold, 586-fold,
180 and 407-fold in motifs specific to Dam, EcoRI, and TaqI, respectively (**Fig. 2D**), consistent with
181 precise m6A calls. The ability to call single m6A events is biologically relevant for fiber-seq as
182 10-25% of all internucleosomal linker regions (i.e., methyltransferase-sensitive patches, MSPs)
183 within a Fiber-seq dataset contain only a single m6A separating adjacent nucleosomes
184 (**Supplemental Fig. S9**).

185
186 Fifth, the precision of *fibertools* in identifying single m6A events indicated it might be useful for
187 identifying endogenously m6A-modified bases within bacteria, a context that also provides a
188 true positive validation set as nearly 100% m6A methylation can be expected within
189 methyltransferase-specific motifs. Notably, DNA isolated from bacteria expressing both the Dam
190 and HsdM methyltransferases exhibited m6A at 94.1% of the target Dam sites for these
191 methyltransferases, indicating a false negative rate of less than 6% in this sequence context
192 (**Fig. 2E**). In contrast, DNA from bacteria lacking these methyltransferases (Anton et al. 2015)
193 exhibited m6A at <1% of adenines, consistent with our prior FPR estimate.

194

195 Sixth, we evaluated the accuracy of *fibertools* for quantifying m6A-marked chromatin
196 architectures along multi-kilobase reads. Current SMRT cell chemistries target sequencing
197 reads ~20 kb in length, yet traditional subread-based m6A models were designed for reads of
198 only ~2kb in length. We demonstrate that in comparison to existing subread-based models,
199 *fibertools* substantially reduce false-negative methylation calls along increasingly longer reads
200 (**Fig. 3, Supplemental Fig. S10**), enabling the accurate quantification of m6A-marked
201 chromatin architectures along reads >25 kb in length (**Fig. 3A,B**). The increased false-negative
202 rate of previous callers is due to a reliance on many subread passes which becomes less likely
203 with increasing insert size. This limitation is avoided in the *fibertools* model, which can call m6A
204 on HiFi reads with any number of subread passes.

205

206 Collectively, these biological validations provide strong evidence that *fibertools* is highly
207 accurate and specific in identifying m6A events using PacBio HiFi data. Importantly, the semi-
208 supervised training design of *fibertools* enables it to readily adapt to new sequencing
209 chemistries, which often contain updated polymerases that may differ in their kinetic values
210 (**Supplemental Fig. S1**). For example, we used calls from the model for the PacBio Sequel II
211 2.2 chemistry as the initialization point for training a semi-supervised model for the PacBio
212 Revio chemistry - demonstrating that this new Revio model is similarly highly accurate in
213 identifying m6A events (**Fig. 1D,E**).

214

215 Having established that *fibertools* can identify highly accurate m6A events using PacBio HiFi
216 data, we next sought to extend *fibertools* to enable the simultaneous processing of genetic,
217 cytosine methylation, and adenine methylation data. To accomplish this, we designed *fibertools*
218 to integrate m6A calls directly into the BAM format using the MM and ML tags (**Supplemental**
219 **Fig. S11**). Next, we optimized *fibertools* using a compiled language, which we provide as a
220 single binary (*ft*) accessible through bioconda (package "*fibertools-rs*"). Of note, *fibertools* can
221 process individual Revio SMRT cells in 15-24 CPU hours and Sequel II SMRT cells in 5-8 CPU
222 hours, a >1,000-fold increase in speed compared to the previous pipeline when using GPU
223 acceleration (>150-fold increase without GPU) (**Fig. 1F, Supplemental Table S2**).

224

225 We next extended the utility of *fibertools* to perform fundamental operations necessary for
226 processing single-molecule epigenetic and genetic data produced using either PacBio and ONT
227 sequencing platforms (i.e., *fibertools add-nucleosome*, *fibertools extract*, and *fibertools center*).

228

229 *Fibertools add-nucleosome* enables the identification of stretches of unmethylated adenine
230 bases, which are stored directly in the BAM file using custom flags. It processes 10 million
231 ~20kb reads in just 4.6 CPU hours. Importantly, *fibertools add-nucleosome* works seamlessly
232 with Fiber-seq data sequenced using either a PacBio or ONT instrument. For example, the
233 application of *fibertools add-nucleosome* to a Fiber-seq library sequenced on an ONT R10.4
234 flow cell and base called using Dorado v0.4.2 (Oxford Nanopore 2023) enabled the robust
235 identification of clear nucleosomal patterns (**Fig. 4A,B**), despite the substantial decrease in m6A
236 and DNA base identification accuracy with ONT sequencing (**Fig. 4C, Supplemental Fig. S3**).

237
238 *Fibertools extract* enables multithreaded conversion of genetic and epigenetic BAM features into
239 plain text formats regardless of upstream tooling with coordinates in either reference or
240 molecular space. For example, nucleosomes, m6A, and 5mC methylation can all be extracted
241 into BED12 format to visualize with the UCSC Genome Browser. Alternatively, all of these
242 features and more can be extracted into a unified table for custom downstream analysis or
243 visualization with both m6A and 5mC base modifications across multiple technologies (**Fig 5**).

244
245 *Fibertools center* enables the processing of single-molecule genetic and epigenetic data (i.e.,
246 DNA base, mCpG, DNA-m6A, inferred epigenetic marks, etc.) relative to a set of reference
247 genomic coordinates while maintaining how these features co-occur along each read using both
248 reference and molecular coordinate systems - addressing a need that is unique to long-read
249 epigenetic studies. To demonstrate the utility of *fibertools center*, we applied it to address two
250 fundamental biological questions that require the integration of long-read epigenetic and genetic
251 data: (1) the relationship between somatic DNA variability and overlying altered epigenetic
252 architecture along individual DNA molecules (**Fig. 5**); and (2) the co-occupancy of transcription
253 factor (TF) binding elements along individual DNA molecules (**Fig. 6**).

254
255 First, we applied *fibertools center* to resolve the relationship between somatic DNA variability
256 and overlying altered epigenetic architecture along telomeric and sub-telomeric regions using
257 Fiber-seq PacBio HiFi, standard PacBio HiFi, and standard ONT sequencing data. Telomere
258 repeat arrays exhibit substantial per-molecule somatic alterations in both their length and
259 sequence content (Dubocanin et al. 2022) and are known to be transcribed into Telomeric
260 Repeat-containing RNA (TERRA) via a CpG island promoter located within a TAR1 repeat
261 positioned adjacent to the majority of telomere repeats (Azzalin et al. 2007). The substantial
262 molecule-to-molecule heterogeneity in the DNA content of individual telomere repeats

263 originating from the same chromosome arm requires the use of a molecular coordinate system
264 for their appropriate analysis (i.e., the original positions from the sequenced read, without
265 adjustments for alignment to the reference sequence). Application of *fibertools center* to Fiber-
266 seq data from HG002 cells aligned to the HG002 reference genome enabled the resolution of
267 both genetic and epigenetic architectures of sub-telomeric regions in reference coordinates and
268 telomere repeat arrays in molecular coordinates by using the telomere-subtelomere boundary
269 as the centering point. Importantly, by using molecular coordinates to display the telomere
270 repeat array, we were able to identify multiple non-TTAGGG telomere variant repeats (TVRs)
271 (Baird et al. 1995; Allshire et al. 1989) absent from the HG002 reference sequence
272 (**Supplemental Fig. S12**). Overall, this analysis exposed that the sub-telomeric TAR1 element
273 was hyper-CpG methylated and lacked chromatin accessibility within HG002 cells (**Fig. 5**), in
274 contrast to CHM13 cells (Dubocanin et al. 2022). Notably, this region of hyper-CpG methylation
275 extended into TVRs harboring CpG dinucleotides within the telomere repeat array itself (Baird et
276 al. 1995; Allshire et al. 1989), which we validated by applying *fibertools* to additional genomic
277 ONT and PacBio HiFi sequencing data from HG002 (Zook et al. 2020). TVRs can undergo
278 somatic conversion back into TTAGGG repeats, likely via somatic shortening and subsequent
279 telomerase-mediated elongation (Dubocanin et al. 2022). Integrating the genetic and epigenetic
280 data along individual telomere fibers identified numerous fibers that had undergone such
281 somatic conversions, resulting in the loss of CpG methylation within that region of the telomere
282 repeat array. However, hyper-CpG methylation of the TAR1 element and the rest of the
283 telomere repeat array was maintained on these fibers, suggesting that CpG methylation of
284 TVRs may be a bystander effect of TAR1 hyper-CpG methylation within HG002 cells.

285

286 Second, we applied Fiber-seq and *fibertools center* to resolve the role of CTCF co-occupancy in
287 guiding higher-order chromatin architectures along the ~175 kbp Epstein-Barr virus (EBV)
288 genome. CTCF elements within the EBV OriP and *LMP* loci are known to form a cohesin-
289 dependent loop important for maintaining viral latent cycle gene expression (Arvey et al. 2012;
290 Morgan et al. 2022). Using *fibertools*, CTCF footprints were resolved on individual DNA
291 molecules (Yin et al. 2017), and we determined that the CTCF ChIP-seq peak within the EBV
292 *LMP* locus is predominantly mediated by CTCF co-occupying two immediately adjacent CTCF
293 binding elements along the same DNA molecule (**Fig. 6**). In addition, we observed that EBV
294 nucleoids bound by CTCF within the *LMP* locus are also preferentially co-bound by a single
295 CTCF element within the EBV OriP locus, which is located ~12 kbp away (**Fig. 6A**). However,
296 only 13.4% of EBV nucleoids are co-bound by CTCF at both the OriP and *LMP* loci (**Fig. 6B**),

297 setting an upper limit on the proportion of EBV nucleoids that are directly structured by CTCF-
298 mediated three-dimensional looping between these two loci in a stable manner. Notably,
299 dynamic CTCF-mediated looping has been observed along the nuclear genome (Gabriele et al.
300 2022), and it is possible that these CTCF co-bound EBV nucleoids may represent similarly
301 dynamic CTCF-mediated looping along the EBV genome.

302 Discussion

303 In summary, *fibertools* provides foundational tooling for processing long-read genetic and
304 epigenetic data. Specifically, *fibertools* enables highly accurate single-molecule DNA-m6A
305 identification using PacBio sequencing with a 1,000-fold improvement in speed (**Fig. 1F**),
306 enabling highly accurate single-molecule chromatin fiber sequencing and endogenous bacterial
307 adenine methylation identification (**Fig. 1D, 2E, Supplemental Fig. S13**). In addition, *fibertools*
308 enables the synchronous processing of multiple types of genetic and epigenetic information
309 present within BAM files produced using either PacBio or ONT sequencing platforms at single-
310 molecule and single-nucleotide resolution (**Fig. 4**). Furthermore, *fibertools* enables the seamless
311 conversion between molecular and reference coordinate systems, allowing for accurate genetic
312 and epigenetic analyses of long-read data within structurally and somatically variable genomic
313 regions (**Fig. 5**). Finally, *fibertools* is written in a compiled language and is available as a single
314 binary accessible through bioconda (package “*fibertools-rs*”) and PyPI (package “*pyft*”), enabling
315 its broad utility within existing computational pipelines for processing long-read data. Recent
316 advances in the cost, accuracy, and throughput of long-read sequencing have enabled the
317 broader adoption of this technology across the genetics community, and we anticipate that
318 *fibertools* will serve as an integral tool for processing long-read sequencing data.

319
320 Finally, the semi-supervised training structure that *fibertools* introduces for identifying m6A-
321 marked bases will likely prove useful for the accurate identification of other base modifications
322 using long-read sequencing (**Fig. 1C, Methods**). Specifically, this approach accounts for
323 imperfect training data, which is a core feature of long-read sequencing data owing to the
324 inherent biological per-molecule heterogeneity in the distribution of modified bases. As such, we
325 anticipate that this approach can both improve the accuracy of identifying core base
326 modifications in humans, as well as novel base modifications in bacteria. Furthermore, we
327 demonstrate that models trained using a semi-supervised structure can readily adapt to updated
328 sequencing chemistries, making this training approach highly advantageous for long-read

329 sequencing chemistries, which are undergoing frequent iterative development cycles (**Fig.**
330 **1D,E**).

331

332 In summary, we present *fibertools* as a toolkit to facilitate the identification and analysis of long-
333 read genetic and epigenetic sequencing data and establish semi-supervised machine learning
334 as an adaptable tool for accurately identifying base modifications using long-read sequencing.

335 Methods

336 Initial m6A calling with ipdSummary followed by filtering with a GMM

337 To initially identify m6A, we use a previously published protocol (Dubocanin et al. 2022) with the
338 following details and modifications. Raw PacBio subread BAM files were converted into CCS
339 (circular consensus sequence) reads using pbccs (v6.0.0) (<https://ccs.how/>) with average
340 kinetics information included. Then subreads were aligned to their respective CCS reads using
341 actc (<https://github.com/PacificBiosciences/actc>). The resulting alignments were passed to
342 ipdSummary (v3.0) (<https://github.com/PacificBiosciences/kineticsTools/>) to identify positions
343 with m6A modifications with the CCS-specific extracted genomic sequence as the reference and
344 these additional flags: --pvalue 0.001 --identify m6A. For each read, we then trained a two-
345 component Gaussian mixture model (GMM) from the ipdRatios generated by ipdSummary for all
346 adenine bases within the read ([push_m6a_to_bam.py](https://github.com/PacificBiosciences/kineticsTools/)). Adenine bases were then classified as
347 m6A modified if the probability that the ipdRatio came from the larger distribution was greater
348 than or equal to 0.99999999. All code to repeat these steps is made available as a single
349 Snakemake (Köster and Rahmann 2012) pipeline on GitHub
350 (<https://github.com/StergachisLab/fiberseq-smk>) and in Supplemental Code. K562 and CHM1
351 PacBio data generated as part of a previous study are available at GEO accession
352 GSE226394.

353 Hidden Markov model nucleosome calling

354 To identify nucleosomes initially, we use a previously published protocol (Dubocanin et al. 2022)
355 with the following details and modifications. We used m6A calls to train a two-state
356 (nucleosome, non-nucleosome) Hidden Markov model (HMM) with a standardized post-hoc
357 correction. The input data for the HMM skipped all G/C bases and considered only the A/T
358 sequence so that GC-rich regions without adenine would not falsely be called nucleosomes.
359 The HMM was then trained using the Baum-Welch algorithm (Baum et al. 1970) using a
360 subsample of 5,000 CCS reads per SMRT cell, and the trained model was subsequently applied
361 across all fibers from that SMRT cell using the Python package pomegranate (Schreiber 2017).
362 In tandem with the HMM delineation of nucleosomes, we also applied a simplified approach that
363 looked for unmethylated stretches larger than 85 bp in length (irrespective of A/T content). We
364 used these 'simple' calls to refine our HMM nucleosome calls by splitting large HMM
365 nucleosomes that contained multiple 'simple' calls. We also refined the terminal boundaries of

366 each nucleosome by bookending it to the nearest m6A call within 10 bp, if present. Nucleosome
367 calling is available as part of a Snakemake pipeline ([https://github.com/StergachisLab/fiberseq-](https://github.com/StergachisLab/fiberseq-smk)
368 [smk](https://github.com/StergachisLab/fiberseq-smk)) and encodes the resulting calls in the custom BAM flags ns (nucleosome start) and nl
369 (nucleosome length), which can be extracted using fibertools-rs.

370 Generating positive and negative labels for training and validation data

371 To generate positive labels, we used the previously existing calls made using the GMM-filtered
372 calls from ipdSummary in our Snakemake pipeline. Negative labels were drawn from
373 nucleosome regions as defined by the HMM since we had increased confidence that these
374 regions should be inaccessible due to the presence of a nucleosome. For each SMRT cell, we
375 selected ~350,000 CCS reads and randomly sampled 5% of available positive and negative
376 m6A positions within each read to reduce the number of adjacent and, therefore, non-
377 independent calls in our training data. This resulted in a dataset of ~100 million negative labels
378 and ~8 million positive labels for each SMRT cell. We note that by selecting negative labels
379 from within nucleosomes, we allow for the possibility for our models to outperform ipdSummary
380 and the GMM correction even though we draw positive labels from these tools. For each label in
381 our dataset, we included a 15 bp window centered around the adenine encoding the sequencing
382 information with one-hot encoding and the CCS kinetics information across the same 15 bases
383 for pulse width and interpulse duration resulting in a 6 × 15 matrix for each labeled position. The
384 15 bp window and inclusion of pulse width and interpulse duration were selected based on the
385 ablation study of the supervised model (**Supplemental Fig. S14**). Code to generate training
386 data from CCS reads is available on GitHub (<https://github.com/fiberseq/train-m6A-calling>) and
387 in Supplemental Code.

388 Training, validation, and testing datasets for machine learning

389 For the 2.2 PacBio sequencing chemistry, we established three separate datasets, each from a
390 different sequencing run, for developing our models. The three datasets were used for training,
391 validation, and a completely held-out testing dataset for determining final accuracies. Data for
392 v2.2 chemistry was generated from previously described samples (K562, CHM1) treated with
393 200 units of Hia5 for 10 minutes at 25°C (Dubocanin et al. 2022). Data for v3.2 chemistry was
394 generated from a K562 Fiber-seq sample (described below), and Revio data was provided by
395 PacBio.

396 Architectures and training of machine-learning models (XGBoost, CNN, and semi-
397 supervised CNN)

398 We trained an XGBoost model with a binary logistic objective function using our training and
399 validation datasets. We selected the following hyperparameters using threefold cross-validation;
400 learning rate of 1 (gamma), maximum tree depth of 8, minimum child weight of 100, and 150
401 estimators. The code to repeat this training is available on GitHub
402 (<https://github.com/fiberseq/train-m6A-calling>) and in Supplemental Code.

403 Convolutional neural network

404 We trained a CNN to predict m6A in Fiber-seq HiFi reads. Input to the CNN model is the 6×15
405 matrix described above. The model has three convolutional layers with 30, 10 and 5 filters of
406 size 5, 5, and 3, respectively. A dense layer of size 25×5 and an output layer of size 5×2
407 follow the convolutional layers. All internal layers have ReLU activation, and the output layer has
408 softmax activation. The output layer has two classes, one for m6A and one for unmethylated
409 adenines. Each class generates scores between 0 and 1 for each input matrix, where scores
410 close to 1 denote high confidence that the input belongs to that class. For example,
411 unmethylated adenines score close to 0 for the m6A class, and methylated adenines score
412 close to 1. The CNN optimizes a binary cross-entropy loss function using the Adam optimizer
413 (Kingma and Ba 2014). We trained the model iteratively for 30 epochs, where the training data
414 was input in random batches of 32 in each epoch. See **Supplemental Table S3** for the size of
415 training and validation datasets from different Fiber-seq chemistries for supervised training. The
416 code to repeat this training is available on GitHub (<https://github.com/fiberseq/train-m6A-calling>)
417 and in Supplemental Code.

418 Semi-supervised convolutional neural network

419 We derive m6A labels from GMM-filtered ipdRatios from ipdSummary. These labels can contain
420 false positives. The lack of clean m6A labels makes the supervised training approach less
421 suitable since it assumes that accurate labels are available. Therefore, we developed a semi-
422 supervised approach, which assumes that our m6A class has a mixed population of true and
423 false positives and our non-m6A class is a clean set. Our training approach is derived from the
424 Percolator and Mokapot proteomics tools for identifying peptides from tandem mass
425 spectrometry data (Käll et al. 2007; Fondrie and Noble 2021). This approach yields a classifier
426 with m6A calls at a target precision. The semi-supervised algorithm is outlined in *Supplementary*
427 *Methods* and *Supplementary Algorithm 1*. First, we split our dataset into training and validation

428 sets stratified by class labels (see **Supplemental Table S3** for the size of training and validation
429 datasets from different Fiber-seq experiment chemistries for semi-supervised training). Then our
430 method proceeds in two phases. In the first phase, we use the IPD score of the central base as
431 a classifier and generate an m6A classification score for all examples in the validation set. The
432 classification score ranks the validation examples, and precision is computed at every score
433 threshold. At the end of this phase, we select a score threshold to achieve the target precision.
434 In this work, we use a target precision of 95%. The second phase is iterative, and each iteration
435 consists of three steps. The first step is selecting a high-confidence m6A training set using the
436 current score threshold. The second step consists of training a CNN model on this training data.
437 In the final step, the validation data is rescored using the trained CNN model from the second
438 step, and a new score threshold is generated at 95% precision with the rescored validation data.
439 In the case of a successful second phase training, the number of positives in the validation data
440 identified at target precision increases with every iteration and plateaus when most m6A
441 examples in the validation data have been identified. We define two conditions for convergence,
442 both of which must be satisfied. First, more than 70% of putative m6A calls from the validation
443 set have been identified. Second, the number of additional m6A calls in a new iteration is less
444 than 1% of the total putative m6A calls. In practice, it took 12, 11, and 3 iterations of phase two
445 training to converge 2.2, 3.2, and Revio chemistry Fiber-seq experiments, respectively
446 (**Supplemental Fig. S15**). The code to repeat this training is available on GitHub
447 (<https://github.com/fiberseq/train-m6A-calling>) and in Supplemental Code.

448 Encoding and selecting precision levels for m6A calling

449 Using the approach outlined in the semi-supervised method, we calculated the empirical
450 precision using the validation data for every score output from the CNN model (see
451 *Supplementary Methods* and **Supplemental Table S4** for details). These precisions are then
452 multiplied by 256 and rounded into an 8-bit integer following the BAM specification for the ML
453 tag (<https://samtools.github.io/hts-specs/SAMtags.pdf>, (Li et al. 2009). We then chose the first
454 value (244) with a precision greater than 95% (244/256) as the threshold for calling positive
455 m6A events.

456 Base calling for Oxford Nanopore

457 HG002 DNA prepared with the standard fiber-seq protocol (below) was sequenced on a
458 PromethION flow cell (R 10.4.1) with a 5kHz sampling rate. Alignment to the HG002 or GRCh38
459 genome and DNA base and base modifications calling were performed with Dorado v.0.4.2 with

460 the basecalling model dna_r10.4.1_e8.2_400bps_sup@v4.2.0, 5mC model V2, and 6mA model
461 V3 (Oxford Nanopore 2023).

462 Preparation of calcium-competent ER2796 *E. coli*

463 A 10mL culture of the ER2796 *E. coli* strain (obtained from NEB) was grown overnight without
464 antibiotics. 1mL of the overnight culture was added to 99mL of fresh LB media without
465 antibiotics and incubated with shaking at 37°C and 200 rpm for 3-4 hours until OD reached 0.4.
466 The culture was separated into two 50mL falcon tubes and placed on ice for 20 minutes before
467 being pelleted by centrifugation at 4°C and 4000 rpm for 10 minutes. The supernatant was
468 discarded, and the cell pellets were resuspended with 20mL of ice-cold 0.1 M CaCl₂ and
469 incubated on ice for 30 minutes. The cells were then pelleted again by centrifugation and the
470 supernatant was discarded. The pellets were then combined by resuspending in 5mL of ice-cold
471 0.1M CaCl₂ with 15% glycerol. Cells were aliquoted in 50ul aliquots, frozen in liquid N₂, and
472 stored at -80C.

473 Plasmid DNA preparation

474 A pCS2+ plasmid containing flag-tagged mgfp5 cloned into the EcoRI/XhoI sites (a gift from Lea
475 Starita) was transformed into two strains of competent *E. coli*, ER2796 (from NEB, see
476 "Preparation of calcium competent ER2796 *E. coli*") and NEB 5-alpha (NEB, cat#C2987I), using
477 a standard heat shock method. 50 µL of the chemically competent cells were thawed on ice and
478 mixed with 50 ng of plasmid DNA. The mixture was placed on ice for 30 minutes before being
479 heat shocked at 42°C for 30 seconds. After heat shock, the cells were immediately placed back
480 on ice for 5 minutes. Following incubation, 950 µL of room temperature SOC media (NEB,
481 cat#B9020S) was added, and the cells were allowed to outgrow for 1 hour in the absence of
482 selection. 50 µL of the outgrown cells were diluted in 5 mL of selective media and grown
483 overnight with shaking at 37°C and 220 rpm. Specifically, the NEB 5-alpha cells were grown in
484 LB media with 100 µg/mL ampicillin, while the ER2796 cells were grown in LB media with 50
485 µg/mL kanamycin + 100 µg/mL ampicillin. The following day, plasmid DNA was extracted from
486 the bacterial cells using the Monarch Plasmid Miniprep Kit (NEB, cat#T1010L) following the
487 manufacturer's protocol. The elution of plasmid DNA was done with sterile water and the
488 concentration was measured using the Qubit 1X dsDNA HS Assay Kit (Invitrogen, cat#Q33231).

489 Cell culture

490 The following cell lines/DNA samples were obtained from the NIGMS Human Genetic Cell
491 Repository at the Coriell Institute for Medical Research: GM12878 and GM24385 (HG002).
492 K562 cells were a gift from the Stamatoyannopoulos lab. Cells were maintained in suspension in
493 IMDM media supplemented with 10% FBS (HyClone, cat#SH30396.03IH25-40) and antibiotic
494 (100 I.U/mL penicillin, 100ug/mL streptomycin, Gibco, cat#15140122) at 37°C and 5% CO₂ in T-
495 75 flasks. Cells were split 1:10 every 3-4 days.

496 gDNA preparation

497 Whole-genome DNA was extracted from K562 cells using HMW DNA extraction kit (Promega,
498 cat#A2920). gDNA was used as a template for all WGA samples. All DNA samples were
499 rehydrated with 10 mM Tris-HCl, pH 8 (ThermoScientific, cat#J22638.K2) or eluted with
500 ultrapure water unless specified.

501 WGA preparation

502 WGA was performed with REPLI-G Mini kit (Qiagen, cat#150023) according to the
503 manufacturer's protocol. 0, 10, 25, 64, 160, 400, or 1000 μM of N⁶-methyl-2-dATP (TriLink,
504 cat#N-2025) (m6dATP) final concentration was spiked into 50 μL of each WGA to generate
505 samples with 0, 2.56, 8.46, 14.45, 33.42, 52.27, and 67.94% m6A samples, as calculated using
506 mass spec according to an eleven point standard quantified every time with the samples during
507 the same run (see below). Samples were purified with ProNex® Size-Selective Chemistry
508 (Promega, cat#NG2001) by adding a 1:1.3 ratio by volume of sample to magnetic beads.
509 Following incubation for 15 minutes, the sample was placed on a magnetic rack for three
510 minutes, washed twice with 80% ethanol, and eluted in 50 μL of PacBio Elution buffer (PacBio,
511 cat#101-633-500). All m6dATP-spiked samples were purified three times with bead purification,
512 and residual free m6dATP quantified by MS of the 10 and 160 μM m6dATP samples treated
513 with Quick CIP only was near background levels.

514 gDNA methylation with sequence-specific methyltransferases

515 Samples labeled with the sequence-specific adenine methyltransferases EcoRI (GAATTC)
516 (NEB, M0211S), Dam (GATC) (NEB, cat#M0222S), or TaqI (TCGA) (NEB, cat#M0219S) were
517 generated as follows: 400 ng of purified gDNA was treated for one hour with 40U of EcoRI
518 (37°C), 10U of TaqI (65°C), or 8U of Dam (37°C) in supplied buffer and in the presence of 160

519 μ M SAM. All samples were purified once by ProNex® Size-Selective Chemistry as described
520 above but with a 1.25:1 ratio by volume of sample to magnetic beads.

521 Fiber-seq

522 Non-specific methyltransferase, Hia5, was purified, and the activity was quantified as previously
523 described (Stergachis et al. 2020). Fiber-seq samples were prepared as previously described
524 (Stergachis et al. 2020).

525 Library preparation

526 Samples were sheared in g-TUBEs (Covaris, cat# 520079) for four passes at 3200 RPM for 2-4
527 minutes in an Eppendorf 5424R centrifuge. Post-shear samples were quantified by Qubit
528 dsDNA high-sensitivity assay (Qubit, cat#Q32851) following the manufacturer's protocol.
529 Multiplexed library preparation was performed using the SMRTbell prep kit 3.0 (PacBio,
530 cat#102-141-700) and SMRTbell barcoded adapter plate 3.0 (bc2001-bc2019) (PacBio,
531 cat#102-009-200) according to the manufacturer's instructions but with the following
532 modifications: after barcoded adapter ligation, samples were incubated at 65°C for 10 minutes
533 to heat-inactivate the ligase. Barcoded samples were pooled and purified with ProNex® Size-
534 Selective Chemistry as described above, with a 1:1 ratio by volume of sample to magnetic
535 beads. Following nuclease treatment, the library was purified first with a 1:3.1 ratio of sample to
536 35% v/v Ampure PB beads (PacBio, cat#100-265-900)/PacBio elution buffer. A second
537 purification was performed with a 1:1 ratio of sample to ProNex® Size-Selective Chemistry, as
538 described. The sample was loaded onto a single Sequel II SMRT cell (v3.2 chemistry) and
539 sequenced by the University of Washington PacBio Sequencing Services. The full composition
540 and sample barcode IDs of the multiplexed library are listed in **Supplemental Table S5**.
541 Plasmid DNA was prepared in a separate multiplexed library. Plasmids were linearized with
542 KpnI prior to multiplexed library preparation.

543 Quantification of m6A/A by UHPLC-MS/MS

544 Samples for quantification were treated as previously described with minor modifications (Kong
545 et al. 2022). In brief, 30-50 ng of DNA from each sample was mixed with 0.02 U
546 phosphodiesterase I (Worthington, cat#LS003926), 1 U Benzonase (Millipore Sigma, cat#
547 E1014), and 2 U Quick CIP (NEB, cat#M0525S) in digestion buffer (10 mM Tris, 1 mM MgCl, pH
548 8 at RT) for 3 hours at 37°C. Single nucleotides were separated from the enzymes by collecting
549 the flow-through of a Nanosep centrifugal filter (MWCO 3 kDa, Pall, cat#OD003C33). The

550 UHPLC-MS/MS analysis of adenosine and m6A was performed on an ACQUITY Premier UPLC
551 System coupled with XEVO-TQ-XS triple quadrupole mass spectrometer. UPLC was performed
552 on a ZORBAX Eclipse Plus C18 column (2.1 × 50 mm I.D., 1.8 μm particle size) (Agilent, cat#
553 959757-902) using 10-90% linear gradient of solvent B (0.1% acetic acid in 100% methanol) in
554 solvent A (0.1% acetic acid in water) within 4 minutes and a flow rate of 0.3 ml/min. MS/MS
555 analysis was operated in positive ionization mode with 3000 V capillary voltage as well as 150
556 °C and 1000 L/Hour nitrogen drying gas. A multiple reaction monitoring (MRM) mode was
557 adopted with the following m/z transition: 252.10 → 136.09 for dA (collision energy, 14 eV), and
558 266.2 → 150.2 for m6A (collision energy, 15 eV). MassLynX was used to quantify the data.

559
560 A calibration curve was generated with 11 mixtures containing different ratios of 2-
561 deoxyadenosine (>99%, FisherScientific, cat#AAJ6388606)(A) to N⁶-methyl-2-deoxyadenine
562 (>99%, FisherScientific, cat#AAJ64961MD)(m6A). A new standard was measured and used for
563 each run. The standard was fit to a third-degree polynomial (equation 1) with y as m6A
564 percentage (%) and X as the quantified MS peak area of m6A over the sum of adenosine and
565 m6A peak area.

$$566 \quad y = a_1 + a_2 X + a_3 X^2 + a_4 X^3 \quad [1]$$

567 Data access

568 All raw and processed sequencing data generated in this study have been submitted to the
569 NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession
570 number PRJNA956114.

571 Software availability

572 Code for training (<https://github.com/fiberseq/train-m6A-calling>), evaluation
573 (<https://github.com/mrvollger/Fiber-m6A-figures-and-tables>), and execution
574 (<https://github.com/fiberseq/fibertools-rs>) of *fibertools* is publicly available on GitHub. These are
575 also included as a supplemental file called Supplemental Code.

576 Competing interests

577 All authors declare no competing interests.

578 Acknowledgments

579 The authors thank Sayeh Gorjifard for organizing the University of Washington Genome
580 Sciences Hackathon, which initiated this work, Michelle Noyes for designing the *fibertools*
581 logo art, and Tonia Brown for assistance in editing this manuscript. Cell lines obtained from
582 the NIGMS Human Genetic Cell Repository at the Coriell Institute for Medical Research
583 include GM12878 and GM24385.

584 Funding

585 A.B.S. holds a Career Award for Medical Scientists from the Burroughs Wellcome Fund and
586 is a Pew Biomedical Scholar. This study was supported by National Institutes of Health
587 (NIH) grants 1DP5OD029630, UM1DA058220, and OT2OD002748 to A.B.S and R01-
588 HG011466 to W.S.N. M.R.V. and S.C.B. were supported by a training grant (T32) from the
589 NIH (2T32GM007454-46).

590 Author contributions

591 Conceptualization and design: M.R.V., A.J., S.C.B., and A.B.S. Experimental design and
592 execution: S.C.B., Y.M., J.R., B.J.M., and A.B.S. Preliminary implementation: M.R.V., A.J.,
593 S.C.B., A.T.M., M.O.H., E.S., C.F., and T.L. Final implementation: M.R.V. and A.J.
594 Supplemental material organization: M.R.V., S.C.B., and A.J. Display items: M.R.V., A.J.,
595 S.C.B., and A.B.S. Manuscript writing: M.R.V., S.C.B., and A.B.S. with input from all
596 authors.

597 **References**

- 598 Abdulhay NJ, Hsieh LJ, McNally CP, Ostrowski MS, Moore CM, Ketavarapu M, Kasinathan S,
599 Nanda AS, Wu K, Chio US, et al. 2023. Nucleosome density shapes kilobase-scale
600 regulation by a mammalian chromatin remodeler. *Nat Struct Mol Biol* **30**: 1571–1581.
601 <http://dx.doi.org/10.1038/s41594-023-01093-6>.
- 602 Abdulhay NJ, McNally CP, Hsieh LJ, Kasinathan S, Keith A, Estes LS, Karimzadeh M,
603 Underwood JG, Goodarzi H, Narlikar GJ, et al. 2020. Massively multiplex single-molecule
604 oligonucleosome footprinting. *Elife* **9**. <http://dx.doi.org/10.7554/eLife.59404>.
- 605 Allshire RC, Dempster M, Hastie ND. 1989. Human telomeres contain at least three types of G-
606 rich repeat distributed non-randomly. *Nucleic Acids Res* **17**: 4611–4627.
607 <http://dx.doi.org/10.1093/nar/17.12.4611>.
- 608 Altemose N, Maslan A, Smith OK, Sundararajan K, Brown RR, Mishra R, Detweiler AM, Neff N,
609 Miga KH, Straight AF, et al. 2022. DiMeLo-seq: a long-read, single-molecule method for
610 mapping protein-DNA interactions genome wide. *Nat Methods* **19**: 711–723.
611 <http://dx.doi.org/10.1038/s41592-022-01475-6>.
- 612 Anderson JD, Widom J. 2000. Sequence and position-dependence of the equilibrium
613 accessibility of nucleosomal DNA target sites. *J Mol Biol* **296**: 979–987.
614 <http://dx.doi.org/10.1006/jmbi.2000.3531>.
- 615 Anton BP, Mongodin EF, Agrawal S, Fomenkov A, Byrd DR, Roberts RJ, Raleigh EA. 2015.
616 Complete Genome Sequence of ER2796, a DNA Methyltransferase-Deficient Strain of
617 *Escherichia coli* K-12. *PLoS One* **10**: e0127446.
618 <http://dx.doi.org/10.1371/journal.pone.0127446>.
- 619 Arvey A, Tempera I, Tsai K, Chen H-S, Tikhmyanova N, Klichinsky M, Leslie C, Lieberman PM.
620 2012. An atlas of the Epstein-Barr virus transcriptome and epigenome reveals host-virus
621 regulatory interactions. *Cell Host Microbe* **12**: 233–245.
622 <http://dx.doi.org/10.1016/j.chom.2012.06.008>.
- 623 Azzalin CM, Reichenbach P, Khoriauli L, Giulotto E, Lingner J. 2007. Telomeric repeat
624 containing RNA and RNA surveillance factors at mammalian chromosome ends. *Science*
625 **318**: 798–801. <http://dx.doi.org/10.1126/science.1147182>.
- 626 Baird DM, Jeffreys AJ, Royle NJ. 1995. Mechanisms underlying telomere repeat turnover,
627 revealed by hypervariable variant repeat distribution patterns in the human Xp/Yp telomere.
628 *EMBO J* **14**: 5433–5443. <http://dx.doi.org/10.1002/j.1460-2075.1995.tb00227.x>.
- 629 Baum LE, Petrie T, Soules G, Weiss N. 1970. A Maximization Technique Occurring in the
630 Statistical Analysis of Probabilistic Functions of Markov Chains. *Ann Math Stat* **41**: 164–
631 171. <http://www.jstor.org/stable/2239727>.
- 632 Cheetham SW, Jafrani YMA, Andersen SB, Jansz N, Kindlova M, Ewing AD, Faulkner GJ.
633 2022. Single-molecule simultaneous profiling of DNA methylation and DNA-protein
634 interactions with Nanopore-DamID. *bioRxiv* 2021.08.09.455753.
635 <https://www.biorxiv.org/content/10.1101/2021.08.09.455753v2> (Accessed April 16, 2024).

- 636 Chen T, Guestrin C. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the*
637 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining,*
638 *KDD '16*, pp. 785–794, Association for Computing Machinery, New York, NY, USA
639 <https://doi.org/10.1145/2939672.2939785> (Accessed April 4, 2023).
- 640 Clark TA, Murray IA, Morgan RD, Kislyuk AO, Spittle KE, Boitano M, Fomenkov A, Roberts RJ,
641 Korlach J. 2012. Characterization of DNA methyltransferase specificities using single-
642 molecule, real-time DNA sequencing. *Nucleic Acids Res* **40**: e29.
643 <http://dx.doi.org/10.1093/nar/gkr1146>.
- 644 Debo BM, Mallory BJ, Stergachis AB. 2023. Evaluation of N6-methyldeoxyadenosine antibody-
645 based genomic profiling in eukaryotes. *Genome Res*.
646 <http://genome.cshlp.org/content/early/2023/03/20/gr.276696.122.abstract>.
- 647 Dubocanin D, Sedeno Cortes AE, Ranchalis J, Real T, Mallory B, Stergachis AB. 2022. Single-
648 molecule architecture and heterogeneity of human telomeric DNA and chromatin. *bioRxiv*
649 2022.05.09.491186. <https://www.biorxiv.org/content/10.1101/2022.05.09.491186v1.full>
650 (Accessed November 22, 2022).
- 651 Flusberg BA, Webster DR, Lee JH, Travers KJ, Olivares EC, Clark TA, Korlach J, Turner SW.
652 2010. Direct detection of DNA methylation during single-molecule, real-time sequencing.
653 *Nat Methods* **7**: 461–465. <http://dx.doi.org/10.1038/nmeth.1459>.
- 654 Fondrie WE, Noble WS. 2021. mokapot: Fast and Flexible Semisupervised Learning for Peptide
655 Detection. *J Proteome Res* **20**: 1966–1971.
656 <http://dx.doi.org/10.1021/acs.jproteome.0c01010>.
- 657 Gabriele M, Brandão HB, Grosse-Holz S, Jha A, Dailey GM, Cattoglio C, Hsieh T-HS, Mirny L,
658 Zechner C, Hansen AS. 2022. Dynamics of CTCF- and cohesin-mediated chromatin
659 looping revealed by live-cell imaging. *Science* **376**: 496–501.
660 <http://dx.doi.org/10.1126/science.abn6583>.
- 661 Hall MA, Shundrovsky A, Bai L, Fulbright RM, Lis JT, Wang MD. 2009. High-resolution dynamic
662 mapping of histone-DNA interactions in a nucleosome. *Nat Struct Mol Biol* **16**: 124–129.
663 <http://dx.doi.org/10.1038/nsmb.1526>.
- 664 Käll L, Canterbury JD, Weston J, Noble WS, MacCoss MJ. 2007. Semi-supervised learning for
665 peptide identification from shotgun proteomics datasets. *Nat Methods* **4**: 923–925.
666 <http://dx.doi.org/10.1038/nmeth1113>.
- 667 Kingma DP, Ba J. 2014. Adam: A Method for Stochastic Optimization. *arXiv [cs.LG]*.
668 <http://arxiv.org/abs/1412.6980>.
- 669 Kong Y, Cao L, Deikus G, Fan Y, Mead EA, Lai W, Zhang Y, Yong R, Sebra R, Wang H, et al.
670 2022. Critical assessment of DNA adenine methylation in eukaryotes using quantitative
671 deconvolution. *Science* **375**: 515–522. <http://dx.doi.org/10.1126/science.abe7489>.
- 672 Köster J, Rahmann S. 2012. Snakemake—a scalable bioinformatics workflow engine.
673 *Bioinformatics* **28**: 2520–2522.
674 <https://academic.oup.com/bioinformatics/article/28/19/2520/290322> (Accessed May 23,
675 2021).

- 676 Lee I, Razaghi R, Gilpatrick T, Molnar M, Gershman A, Sadowski N, Sedlazeck FJ, Hansen KD,
677 Simpson JT, Timp W. 2020. Simultaneous profiling of chromatin accessibility and
678 methylation on human cell lines with nanopore sequencing. *Nat Methods* **17**: 1191–1199.
679 <http://dx.doi.org/10.1038/s41592-020-01000-7>.
- 680 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
681 1000 Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map
682 format and SAMtools. *Bioinformatics* **25**: 2078–2079.
683 <http://dx.doi.org/10.1093/bioinformatics/btp352>.
- 684 Loman NJ, Quick J, Simpson JT. 2015. A complete bacterial genome assembled de novo using
685 only nanopore sequencing data. *Nat Methods* **12**: 733–735.
686 <http://dx.doi.org/10.1038/nmeth.3444>.
- 687 Luger K, Mäder AW, Richmond RK, Sargent DF, Richmond TJ. 1997. Crystal structure of the
688 nucleosome core particle at 2.8 Å resolution. *Nature* **389**: 251–260.
689 <http://dx.doi.org/10.1038/38444>.
- 690 Marks P, Banerjee O, Alexander D. 2012. Detection and Identification of Base Modifications
691 with Single Molecule Real-Time Sequencing Data.
- 692 Morgan SM, Tanizawa H, Caruso LB, Hulse M, Kossenkov A, Madzo J, Keith K, Tan Y, Boyle S,
693 Lieberman PM, et al. 2022. The three-dimensional structure of Epstein-Barr virus genome
694 varies by latency type and is regulated by PARP1 enzymatic activity. *Nat Commun* **13**: 187.
695 <http://dx.doi.org/10.1038/s41467-021-27894-1>.
- 696 Murray IA, Clark TA, Morgan RD, Boitano M, Anton BP, Luong K, Fomenkov A, Turner SW,
697 Korlach J, Roberts RJ. 2012. The methylomes of six bacteria. *Nucleic Acids Res* **40**:
698 11450–11462. <http://dx.doi.org/10.1093/nar/gks891>.
- 699 Neph S, Kuehn MS, Reynolds AP, Haugen E, Thurman RE, Johnson AK, Rynes E, Maurano
700 MT, Vierstra J, Thomas S, et al. 2012. BEDOPS: high-performance genomic feature
701 operations. *Bioinformatics* **28**: 1919–1920. <http://dx.doi.org/10.1093/bioinformatics/bts277>.
- 702 Nurk S, Koren S, Rhie A, Rautiainen M, Bzikadze AV, Mikheenko A, Vollger MR, Altemose N,
703 Uralsky L, Gershman A, et al. 2022. The complete sequence of a human genome. *Science*
704 **376**: 44–53. <http://dx.doi.org/10.1126/science.abj6987>.
- 705 Oxford Nanopore. 2023. Dorado. *Dorado*. <https://github.com/nanoporetech/dorado>.
- 706 Polach KJ, Widom J. 1995. Mechanism of protein access to specific DNA sequences in
707 chromatin: a dynamic equilibrium model for gene regulation. *J Mol Biol* **254**: 130–149.
708 <http://dx.doi.org/10.1006/jmbi.1995.0606>.
- 709 Quinlan AR. 2014. BEDTools: The Swiss-army tool for genome feature analysis. *Curr Protoc*
710 *Bioinformatics* **47**: 11.12.1–34. <http://dx.doi.org/10.1002/0471250953.bi1112s47>.
- 711 Razaghi R, Hook PW, Ou S, Schatz MC, Hansen KD, Jain M, Timp W. 2022. Modbamtools:
712 Analysis of single-molecule epigenetic data for long-range profiling, heterogeneity, and
713 clustering. *bioRxiv* 2022.07.07.499188.
714 <https://www.biorxiv.org/content/10.1101/2022.07.07.499188> (Accessed November 18,
715 2023).

- 716 Schreiber J. 2017. Pomegranate: fast and flexible probabilistic modeling in python. *arXiv [csAI]*.
 717 <https://www.jmlr.org/papers/volume18/17-636/17-636.pdf?ref=https://githubhelp.com>
 718 (Accessed March 24, 2023).
- 719 Shipony Z, Marinov GK, Swaffer MP, Sinnott-Armstrong NA, Skotheim JM, Kundaje A,
 720 Greenleaf WJ. 2020. Long-range single-molecule mapping of chromatin accessibility in
 721 eukaryotes. *Nat Methods* **17**: 319–327. <http://dx.doi.org/10.1038/s41592-019-0730-2>.
- 722 Stergachis AB, Debo BM, Haugen E, Churchman LS, Stamatoyannopoulos JA. 2020. Single-
 723 molecule regulatory architectures captured by chromatin fiber sequencing. *Science* **368**:
 724 1449–1454. <http://dx.doi.org/10.1126/science.aaz1646>.
- 725 Töpfer A, Wenger A. 2023. Jasmine: Predict 5mC in PacBio HiFi reads.
 726 <https://github.com/PacificBiosciences/jasmine>.
- 727 Tse OYO, Jiang P, Cheng SH, Peng W, Shang H, Wong J, Chan SL, Poon LCY, Leung TY,
 728 Chan KCA, et al. 2021. Genome-wide detection of cytosine methylation by single molecule
 729 real-time sequencing. *Proc Natl Acad Sci U S A* **118**.
 730 <http://dx.doi.org/10.1073/pnas.2019768118>.
- 731 Vollger MR, Logsdon GA, Audano PA, Sulovari A, Porubsky D, Peluso P, Wenger AM,
 732 Concepcion GT, Kronenberg ZN, Munson KM, et al. 2020. Improved assembly and variant
 733 detection of a haploid human genome using single-molecule, high-fidelity long reads. *Ann*
 734 *Hum Genet* **84**: 125–140. <http://dx.doi.org/10.1111/ahg.12364>.
- 735 Wenger AM, Peluso P, Rowell WJ, Chang P-C, Hall RJ, Concepcion GT, Ebler J,
 736 Functammasan A, Kolesnikov A, Olson ND, et al. 2019. Accurate circular consensus long-
 737 read sequencing improves variant detection and assembly of a human genome. *Nat*
 738 *Biotechnol*. <http://dx.doi.org/10.1038/s41587-019-0217-9>.
- 739 Yin M, Wang J, Wang M, Li X, Zhang M, Wu Q, Wang Y. 2017. Molecular mechanism of
 740 directional CTCF recognition of a diverse range of genomic sites. *Cell Res* **27**: 1365–1377.
 741 <http://dx.doi.org/10.1038/cr.2017.131>.
- 742 Zook JM, Hansen NF, Olson ND, Chapman L, Mullikin JC, Xiao C, Sherry S, Koren S, Phillippy
 743 AM, Boutros PC, et al. 2020. A robust benchmark for detection of germline large deletions
 744 and insertions. *Nat Biotechnol* **38**: 1347–1355. [http://dx.doi.org/10.1038/s41587-020-0538-](http://dx.doi.org/10.1038/s41587-020-0538-8)
 745 [8](http://dx.doi.org/10.1038/s41587-020-0538-8).

746 Figure Legends

747 **Figure 1. Accurate identification of m6A with supervised machine learning (ML) and**
 748 **refinement with semi-supervised ML. A)** Methodology for generating training data and
 749 identifying m6A modifications using PacBio HiFi (see *Methods* for details). **B)** Receiver
 750 operating characteristic and precision-recall curves for the CNN (purple), XGBoost (orange),
 751 and ipdSummary (red) models. Dashed lines indicate the performance of a random classifier. **C)**
 752 Methodology for semi-supervised machine learning (see *Methods* for details). **D)** Autocorrelation
 753 between m6A calls made by the subread (red), semi-supervised (purple), and Revio (pink)

754 models. **E)** Density of the distance between adjacent m6A on the same chromatin fiber (10,000
 755 reads) for the same datasets and models as in d. **F)** CPU hours used by *fibertools* (purple) and
 756 subread-based GMM model (red) for individual SMRT cells. *Fibertools* was run with GPU
 757 acceleration (NVIDIA A40), which is unavailable for the GMM model. **G)** Visualization of m6A
 758 calls in the *HMBS* locus that are unique to *fibertools* (purple), unique to the subread GMM
 759 model (red), or shared by both (gray). Reads are sorted by the number of CCS passes (low to
 760 high). DNase-seq (ENCODE July 2012 Freeze) and CAGE signal is shown above.

761
 762 **Figure 2. Biological validation of the semi-supervised m6A caller.** **A)** Description of
 763 biological samples used for validation of *fibertools*. **B)** Percent of methylated adenines called by
 764 *fibertools* relative to all adenines in a whole-genome amplified (WGA) negative control as a
 765 function of the estimated precision reported by *fibertools*. This serves as an estimate of the false
 766 positive rate. The red line marks the default threshold used by *fibertools*. **C)** Percent m6A as
 767 determined by UHPLC–MS/MS (y) and *fibertools* (x) at the default precision level for WGA
 768 samples with varying levels of m6ATP spiked-in. The text (upper left) indicates the value of the
 769 Pearson correlation coefficient and the P value from a two-sided *t*-test without adjustment
 770 for multiple comparisons. **D)** Enrichment of m6A calls within targeted motifs of three motif-
 771 specific methyltransferases [*Dam* (blue), *EcoRI* (orange), and *TaqI* (green)] as a function of
 772 *fibertools* estimated precision. **E)** Methylation percent at recognition sites for *Dam* (purple),
 773 *HsdM* (orange), and other sites (green) among all sequencing reads of a plasmid grown in a
 774 *dam*⁺/*hsdM*⁺ *E. coli* strain (top) compared to a *dam*⁻/*hsdM*⁻ negative control (bottom). Dotted
 775 lines show the average across each category.

776
 777 **Figure 3. Increased m6A calling on long reads (>20kb) via fibertools.** **A)** Percent increase
 778 in *fibertools* m6A calls over the GMM model as a function of the minimum read length of the
 779 underlying sequencing data. The histogram below shows how many reads were used to
 780 calculate each percent increase. **B)** Comparison of *fibertools* and the subread model for m6A
 781 calling over CAGE-positive TSS in K562 cells across the genome, separated by read length.
 782 Reads are matched between *fibertools* and the subread model, and the number of Fiber-seq
 783 reads used in the calculation of each size range (n) is indicated.

784
 785 **Figure 4. Fibertools nucleosome calling with PacBio and ONT Fiber-seq.** **A)** Data
 786 processing pipelines for nucleosome calling with *fibertools*. **B)** Density of nucleosome lengths
 787 called by *fibertools* for PacBio (pink) and ONT (blue) Fiber-seq (n = 500,000 nucleosomes). **C)**

788 Visualization of the *NAPA* and *HMBS* representative loci for PacBio (top) and ONT (bottom)
789 Fiber-seq. m6A calls from *fibertools* (PacBio) or *Dorado* (ONT) are represented by vertical
790 purple dashes, along with nucleosome (gray) and MTase-sensitive patch (MSP) (orange) calls
791 from *fibertools*.

792

793 **Figure 5. Organization of the HG002 telomere. (Top)** HG002 PacBio Fiber-seq with genetic
794 variants, m6A methylation (purple), and CpG methylation (brown) overlaid for fibers overlapping
795 the maternal telomere of Chromosome 13q (telomeric boundaries were determined with seqtk
796 v1.3). Telomeric sequences (blue), telomeric variants (red), and non-telomeric sequences (gray)
797 are highlighted to show telomeric genetic variation. **(Middle)** ONT and **(bottom)** PacBio
798 standard sequencing of HG002 telomeres (Zook et al. 2020) with CpG methylation overlaid on
799 chromatin architecture.

800

801 **Figure 6. CTCF co-occupancy along the EBV genome. A)** CTCF ChIP-seq (green),
802 significant (red) or insignificant (gray) HiC loops (Morgan et al. 2022), and significant CTCF site
803 co-occupancy by Fiber-seq (purple) along the EBV genome. The significance of CTCF site co-
804 occupancy was determined by comparing the expected number of co-occupied fibers to the
805 observed number using Fisher's Exact test (see **Supplemental Table S6** for exact counts and
806 p-values). **B)** Zoom-in of the indicated CTCF peak, which contains two CTCF binding elements.
807 Single-molecule occupancy and co-occupancy from Fiber-seq are shown below.

Figure 1

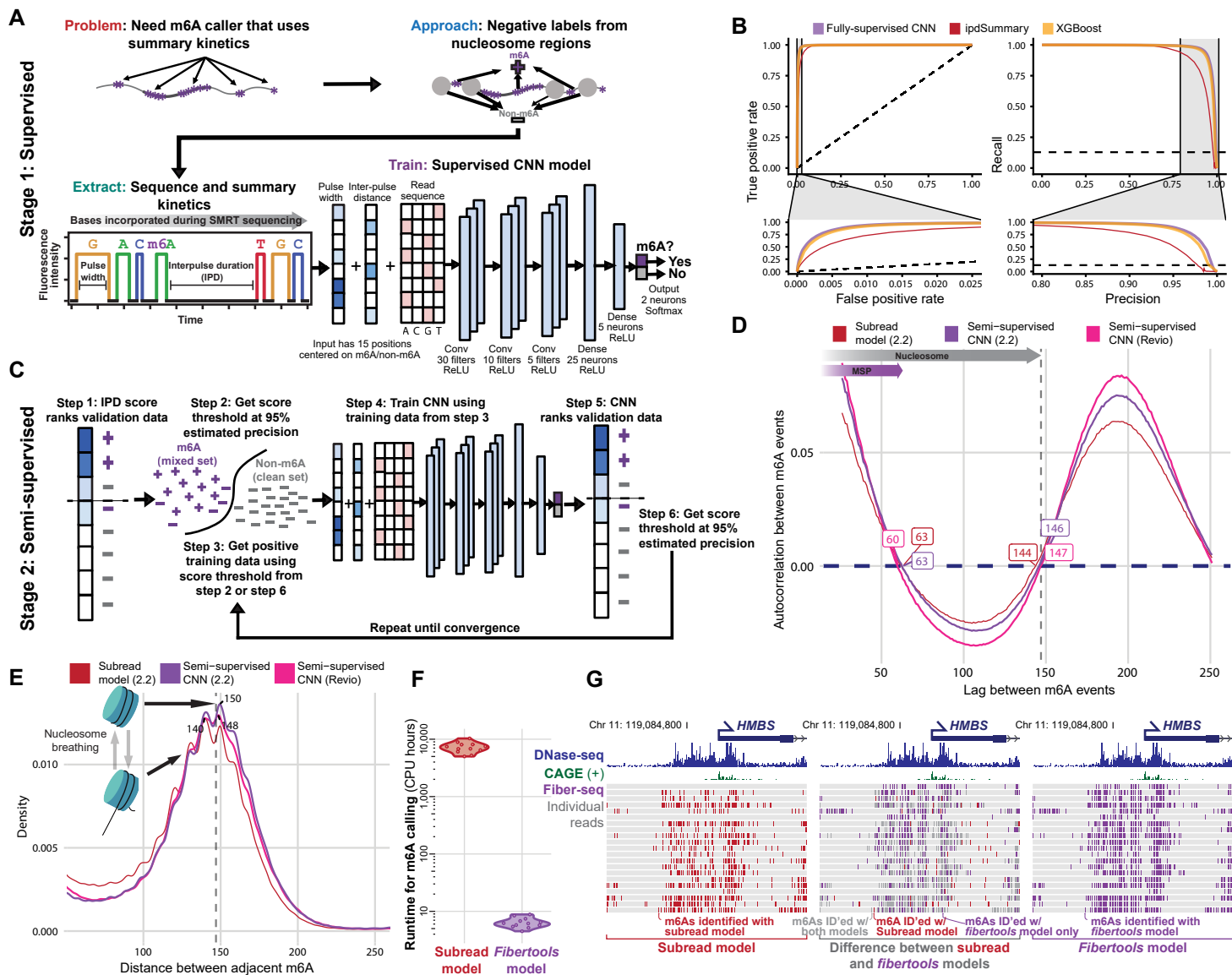


Figure 2

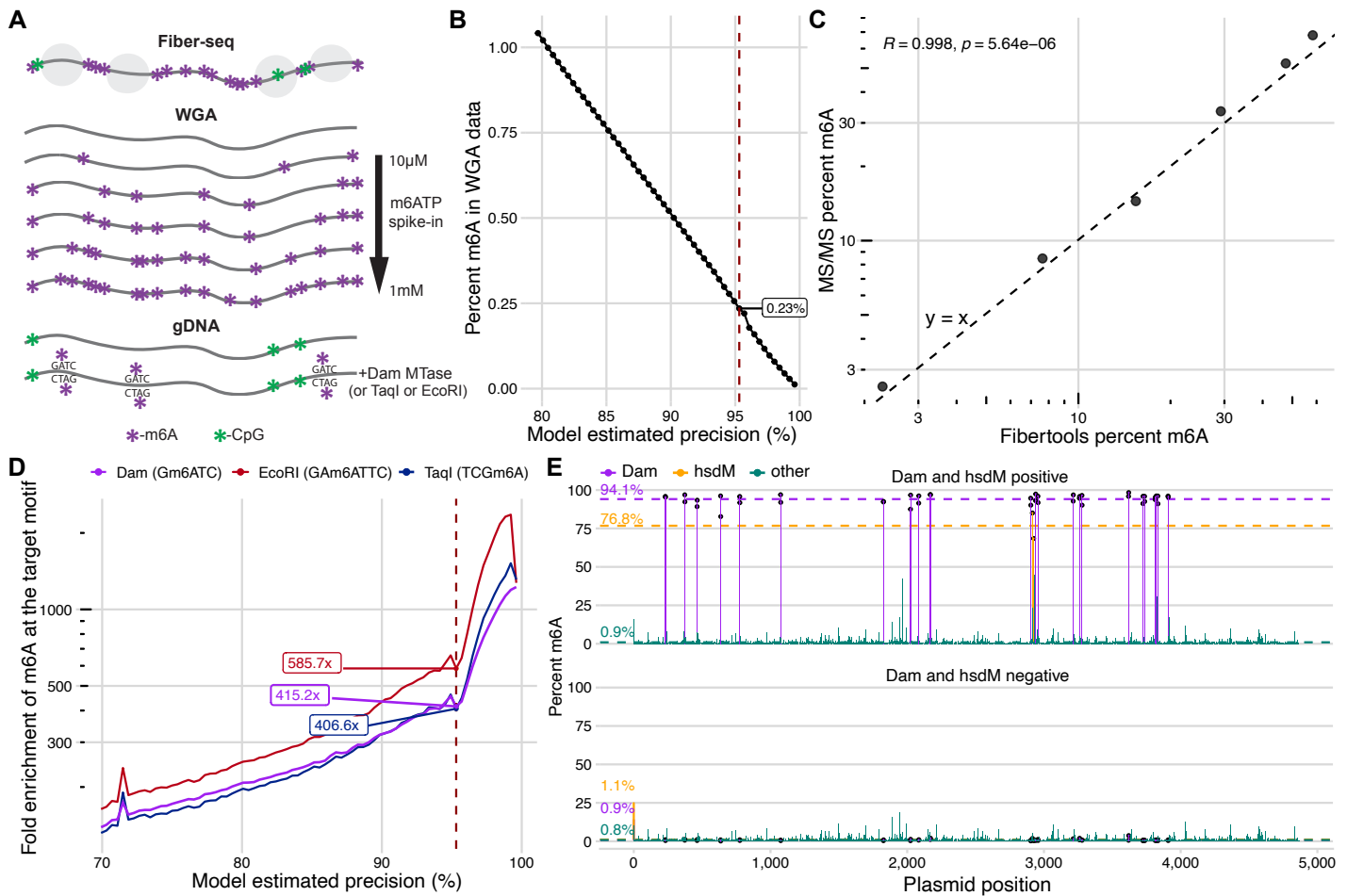
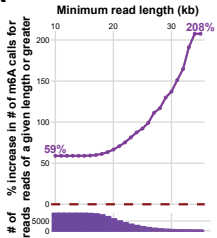


Figure 3

A



B

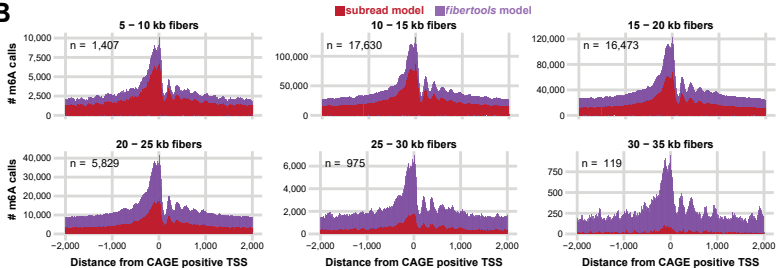


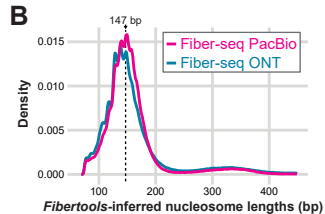
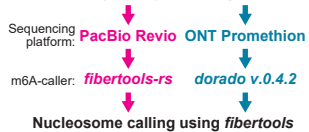
Figure 4**A** GM24385 (HG002) Fiber-seq reaction**C**

Figure 5

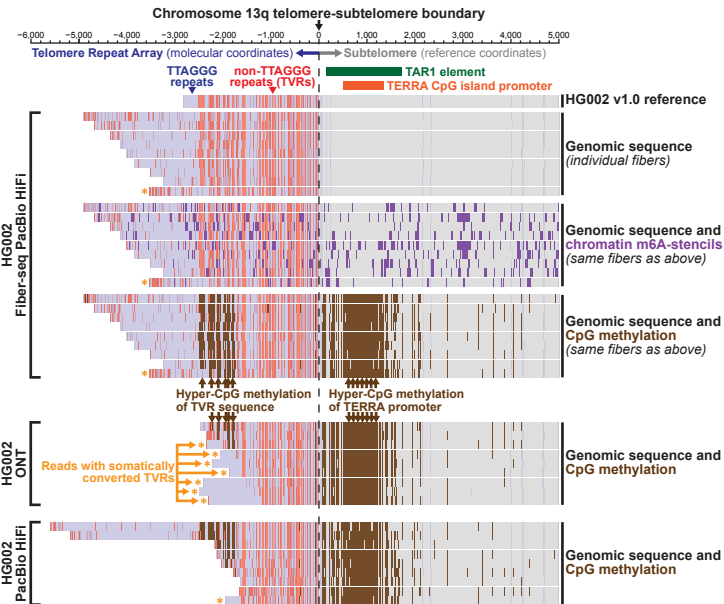
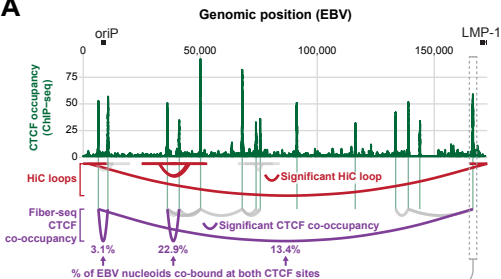


Figure 6**A****B**