

---

# Detecting Sparse Colorectal Cancer Signals from Multi-Modal Cell-Free DNA Representations Using Modern Hopfield Attention

---

Michael Widrich<sup>1</sup> Anooj Patel<sup>1</sup> Elisabeth Rumetshofer<sup>1</sup> Peter Ulz<sup>1</sup> Kaitlyn Coil<sup>1</sup> Thomas Royce<sup>1</sup>  
Cheng-Ho Jimmy Lin<sup>1</sup> Richard Bourgon<sup>1</sup> Anindita Dutta<sup>1</sup>

## Abstract

Next-generation sequencing-based early cancer detection from cell-free DNA (cfDNA) presents an extreme-scale multiple instance learning challenge: identifying rare tumor signals amidst millions of instances per sample (with witness rates as low as  $<0.0001\%$ ). It also provides a testbed for translating DNA foundation model embeddings to a clinically important supervised learning task. We propose Fragment-Level Deep Learning (FLDL), an end-to-end deep learning framework that combines multi-modal fragment features, including a HyenaDNA-derived sequence embedding, with Modern Hopfield Networks to perform dense associative retrieval over the massive cfDNA instance space. Using held-out real-world clinical and challenging contrived test sets, we compare FLDL’s performance to a state-of-the-art machine learning model and to a deep learning model without attention (max pooling). Our results demonstrate that only the attention-based FLDL model outperforms the machine learning model, in spite of a modest training set size ( $n = 4,394$ ). FLDL also scales effectively with sample size and with the number of instances per sample while offering useful biological insights via attention weights and learned sample representations. This work establishes a new frontier for DNA foundation-model-augmented cfDNA representations and highly scalable attention-based deep learning in clinical cfDNA diagnostics.

## 1. Introduction

Foundation models are rapidly expanding the scope of machine learning in the life sciences, particularly for genomic

---

<sup>1</sup>Freenome, Brisbane, CA, USA. Correspondence to: Anindita Dutta <anindita.dutta@freenome.com>.

*Proceedings of the ICML 2026 3rd Workshop on Multi-modal Foundation Models and Large Language Models for Life Sciences*, Seoul, South Korea. 2026. Copyright 2026 by the author(s).

sequence analysis. Recently, genomic foundation models have emerged, including DNA-BERT, HyenaDNA, Nucleotide Transformer, Evo 2, and related models, which learn sequence representations from large genomic corpora and have shown strong performance on regulatory and functional genomics benchmarks (Ji et al., 2021; Nguyen et al., 2023; Dalla-Torre et al., 2025; Brixi et al., 2026). Despite the proliferation of DNA foundation models, there are few demonstrations of clinical diagnostic workflows leveraging their learned genomic representations to improve predictions from noisy biomedical data. **Colorectal cancer (CRC)** detection from blood provides a particularly stringent clinical test case for this question. CRC remains the second most common cause of cancer-related death in the US (Siegel et al., 2025). Despite the availability of multiple screening modalities, only an estimated 63% of eligible individuals aged 45 years and older were up to date with guideline-recommended screening in 2023 (Bandi et al., 2025), motivating the development of more convenient, non-invasive blood-based tests.

**Cell-free DNA (cfDNA)** analysis is central to non-invasive blood-based tests. Next-generation sequencing (NGS) based cfDNA assays generate millions of short DNA fragments per sample yielding a high-dimensional, naturally multi-resolution signal spanning nucleotide-level, fragment-level, and sample-level structure. In this work, we focus on **DNA methylation**, the conversion of cytosine to 5-methylcytosine (mC) at CpG dinucleotides. Aberrant methylation patterns in cfDNA are a well-established non-invasive cancer biomarker (Yamaguchi et al., 2003), and NGS-based methylation blood tests for CRC and advanced precancerous lesions (APLs) have recently achieved or are approaching FDA approval (Chung et al., 2024; Shaukat et al., 2025). Thus, individual CpG methylation states and fragment-level methylation patterns constitute additional relevant modalities. Unlike approaches profiling cancer tissue, which benefit from relatively high tumor DNA content (sometimes exceeding 50%), a hallmark of cfDNA based cancer detection is that cancer signal is usually very sparse, especially for early-stage disease, with some “low-shedding” cancers releasing little to no detectable circulating tumor DNA (ctDNA) (Bettegowda et al., 2014; Luo et al., 2021).

Thus, cfDNA-based early cancer detection tests present a formidable computational challenge: identifying a minute number of ctDNA “needles” within an enormous haystack of millions of healthy cfDNA fragments.

**Traditional approaches** to ctDNA detection largely rely on biological prior knowledge and summarize methylation measurements through aggregation at the per-CpG ( $\beta$ ) or per-fragment ( $\alpha$ ) level (Li et al., 2018). While effective, these methods depend on manual definitions rather than learning optimal representations directly from raw data. Recently proposed multi-step deep learning approaches for methylation sequencing data (Jeong et al., 2025; Niki et al., 2025; Deng et al., 2023) demonstrate the value of fragment-level representation learning. However, these methods are typically not trained end-to-end for final patient-level classification and instead rely on separate representation learning and downstream aggregation. We hypothesize that robust detection of early-stage cancer benefits from a) the ability to identify and upweight rare informative fragments, and b) representations jointly optimized with the classification objective to capture subtle task-specific signatures.

To enable multi-modal fragment representations and multi-resolution learning across nucleotides/CpG, fragment, genomic context of the fragment and sample levels, we propose the **Fragment-Level Deep Learning (FLDL) model**, a method that formulates early cancer detection as a multiple instance learning (MIL) problem (Carbonneau et al., 2018; Ilse et al., 2018). FLDL incorporates a specialized attention module based on Modern Hopfield Networks (MHN) (Ramsauer et al., 2020), which have been shown to be successful for immune repertoire classification (Widrich et al., 2020). By adapting this mechanism to ctDNA detection, our method retains the clinically relevant extreme scale “needle-in-a-haystack” structure of cfDNA: rare fragments can receive high attention without requiring quadratic self-attention across all fragments.

In this work, we show that FLDL for blood-based detection of CRC a) operates on billions of nucleotides from millions of cfDNA fragments per subject, b) shows utility of a multi-modal fragment representation including DNA foundation model embeddings in a clinical setting, c) directionally outperforms a state-of-the-art machine learning baseline on contrived and real-world clinical test sets, d) is capable of implicit denoising of large and sparse input spaces, e) supports biological interpretability via attention weights and generalizability assessment via learned sample embeddings, and f) scales consistently with training data volume. Together, these results show that multi-modal fragment-level learning with DNA foundation-model sequence representations and scalable attention can be combined in an end-to-end model for a clinically meaningful cfDNA-based diagnostic task.

## 2. Fragment-Level Deep Learning model

We base our FLDL model on the attention mechanism of continuous MHNs (Ramsauer et al., 2020), a generalization of Hopfield Networks (Hopfield, 1982; 1984), that provides a trainable associative memory for deep learning architectures (Ramsauer et al., 2020; Hu et al., 2023). Due to their exponential storage capacity, they excel in low signal-to-noise problems, as demonstrated by the DeepRC model in Widrich et al. (2020). FLDL extends DeepRC to the prediction of CRC from cfDNA.

**Problem formulation.** We follow the formulation of the MIL problem in Widrich et al. (2020), where a bag  $\mathcal{X} = \{s_1, \dots, s_N\}$  of  $N$  instances constitutes a sample. Assuming binary classification, each such instance  $s_i$  is associated with an inaccessible label  $y_i \in \{0, 1\}$ . Only a sample-level label  $y = \max_i y_i$  is observed per bag. Correct classification of a positive sample therefore requires identifying the instances responsible for the positive label (Foulds & Frank, 2010). Here, each cfDNA fragment constitutes an instance and each patient sample a bag, with  $N$  typically ranging from 1 to 10 million and witness rates (fraction of instances indicating a positive label) as low as  $< 0.0001\%$  in low shedding cases (compared to  $N \approx 300,000$  and  $0.01\%$  in DeepRC).

**FLDL model.** We design FLDL as an end-to-end trainable deep learning architecture (Fig. 1). A dedicated sub-network  $\phi(\cdot)$  (*Fragment Embedding*) maps each fragment  $s_i$  independently to a fixed-size vector  $\mathbf{h}_i = \phi(s_i) \in \mathbb{R}^m$ . FLDL then aggregates the embedded fragments into a sample representation using a MHN (*Sample Embedding*) and predicts CRC status using a fully connected *Output Network*. Each fragment  $s_i$  is represented by five modalities: a) sequence and CpG methylation, b) methylation statistics, c) HyenaDNA foundation model embedding (Nguyen et al., 2023), d) genomic position, e) strand. The HyenaDNA embedding, computed from the human reference genome, brings broad genomic context priors for the fragment while other modalities encode assay-specific methylation, position and sequence information that are not captured by the DNA foundation model. Each modality is projected to a latent embedding by a dedicated sub-network. The individual modality embeddings are concatenated and further processed to form a unified fragment representation  $\mathbf{h}_i \in \mathbb{R}^m$ . See also Appendix Section A.1 and Fig. A1.

To aggregate these fragment representations into a sample-level embedding, we employ a Hopfield Pooling layer following DeepRC. As shown in Fig. 1, we learn a set of  $K$  cancer-indicative fragment prototypes (*state patterns* or *queries*)  $\mathbf{Q} \in \mathbb{R}^{K \times d}$  in a high-dimensional association space. A fully-connected self-normalizing neural network (Klambauer et al., 2017),  $\text{NN} : \mathbb{R}^m \rightarrow \mathbb{R}^d$ , is applied row-wise to project the  $N$  embedded fragments

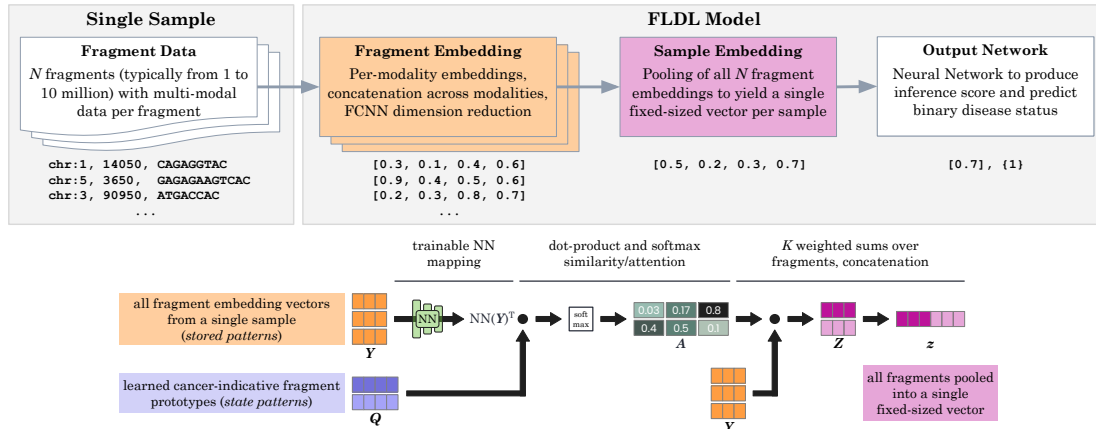


Figure 1. Overview of the FLDL architecture. **Top:** The end-to-end pipeline embeds multi-modal data from millions of cfDNA fragments and aggregates them into a sample-level representation. **Bottom:** Hopfield Pooling uses an attention matrix ( $A$ ) to aggregate fragment embeddings ( $Y$ ) into a sample representation ( $Z$ ) based on similarity to learned cancer-indicative prototypes ( $Q$ ).

$Y = [h_1, \dots, h_N] \in \mathbb{R}^{N \times m}$  (stored patterns) to the same high-dimensional association space  $\mathbb{R}^{N \times d}$  (yielding keys). The softmax-normalized dot-product similarity (attention values)  $A = \text{softmax}(\beta Q \text{NN}(Y)^T) \in \mathbb{R}^{K \times N}$  between these queries and the keys is then used to aggregate the embedded fragments  $Y$  into a sample embedding  $Z = AY \in \mathbb{R}^{K \times m}$ , where  $\beta$  controls the attention distribution sharpness.

Note that the fragment embeddings  $Y$  serve two distinct roles: they are projected via NN to form keys for similarity matching and also serve as values (Vaswani et al., 2017) for aggregation. Thus, the model can learn to assign high attention weights to fragments resembling tumor-derived signal by optimizing the association space and learned state patterns. Finally,  $Z$  is flattened into a vector  $z \in \mathbb{R}^{mK}$  and passed to an MLP to predict the probability of CRC and for use in biologically relevant auxiliary tasks. See Appendix Section A.2 for details on FLDL training.

### 3. Experimental setup and results

We evaluate three FLDL configurations, each differing by the subset of the targeted capture panel used: ER-FLDL uses the full capture panel to test performance in a larger, noisier input space, C-FLDL uses the CRC panel, a section of the full panel most relevant for CRC, and PD-FLDL uses biologically informed pre-filtering on the CRC panel. We further compare against MaxPool variants that replace Hopfield pooling with max-activation-based fragment aggregation, and against a state-of-the-art machine learning baseline (ML baseline) (Shaukat et al., 2025) developed on the CRC panel. These methods are described in detail in Appendix Section B. Datasets used for training the FLDL models, hold-out test sets for performance comparison and model interpretability analysis are described in Appendix

Section C. Note that PD-FLDL regions are a subset of C-FLDL regions, which are a subset of ER-FLDL regions. Using a larger panel with more regions results in more fragments per sample, highlighted for the training and hold-out test sets in Appendix Section D and Appendix Fig. A3.

#### 3.1. Comparison of predictive performance

**Predictive performance, real-world clinical sample test set.** One primary indicator of the model’s efficacy is its predictive performance on the real-world clinical sample test set. For blood-based CRC screening, a specificity at or near 90% is considered to be the most clinically appropriate (Chung et al., 2024; Shaukat et al., 2025). As a result, the most relevant performance metric is not the area under the receiver operating characteristic curve (AUROC), but rather, the sensitivity at a fixed specificity of 90%. We report the achieved APL and CRC sensitivities at a specificity of 90%, along with Wilson’s method 95% confidence intervals (CI), in Table 1. For each method, a classification threshold is selected to achieve the desired specificity among the negative samples in the real-world clinical sample set. As shown

Table 1. Real-world clinical sample test set sensitivity at 90% specificity, with Wilson’s method 95% CIs, for all models.

|             | APL                      | CRC                      |
|-------------|--------------------------|--------------------------|
| ML Baseline | 27.1 (22.8, 31.9)        | 88.2 (82.9, 92.1)        |
| PD-FLDL     | <b>30.2 (25.7, 35.0)</b> | <b>89.6 (84.5, 93.2)</b> |
| C-FLDL      | 28.9 (24.5, 33.7)        | 89.1 (83.9, 92.8)        |
| ER-FLDL     | 28.4 (24.0, 33.2)        | 88.2 (82.9, 92.1)        |
| PD-MaxPool  | 19.9 (16.1, 24.3)        | 79.2 (73.0, 84.3)        |
| C-MaxPool   | 17.3 (13.8, 21.5)        | 76.3 (69.9, 81.8)        |

in Table 1, PD-FLDL outperforms all competing methods, achieving sensitivities of 30.2% for APLs and 89.6% for

CRCs. This improvement over the ML Baseline is more pronounced in the challenging APL group: a gain of 3.1 points over the ML Baseline vs. a gain of 1.4 points for CRCs. Further, PD-FLDL outperforms C-FLDL and ER-FLDL, which suggests that leveraging explicit biological priors to denoise background signals before FLDL training is an effective mechanism for improving model detection sensitivity, at least at the current training set sample size.

It is noteworthy that the C-FLDL model, which operates on the CRC panel but does not use an additional negative sample set for biologically informed filtering, achieves a sensitivity of 28.9% for APL and 89.1% for CRC, also outperforming the ML Baseline. Additionally, even when the input space is significantly expanded to include more task-irrelevant genomic regions (ER-FLDL), the FLDL architecture performs competitively against the ML Baseline model. In contrast, models employing max pooling (PD-MaxPool, C-MaxPool) show significantly degraded performance, with sensitivities dropping below 20% for APLs and below 80% for CRCs. This evidence supports the hypothesis that a simple max pooling aggregation is insufficient for identifying sparse ctDNA signals; instead, the ability to dynamically attend to rare informative fragments (here via MHN) contributes to high-sensitivity cancer detection.

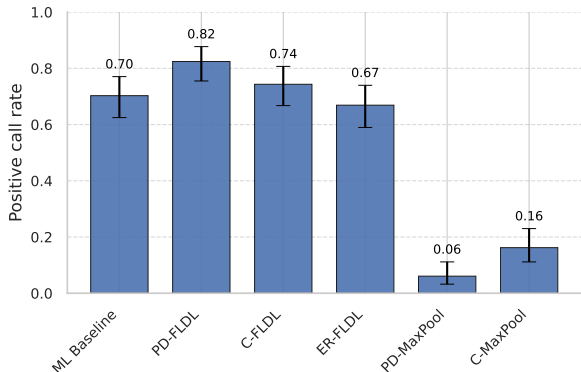


Figure 2. Positive call rates among 148 replicates of the challenging contrived test set, at 90% specificity. Whiskers indicate Wilson’s method 95% CIs.

**Predictive performance, challenging contrived test set.**

The challenging contrived test set provides control over the true ctDNA level, and the blending process provides enough plasma volume for substantially more replication than is possible with a conventional clinical sample. Using the same model-specific 90% specificity classification thresholds as in Table 1, we report each model’s proportion of positive predictions among the 148 replicates.

Performance on the challenging contrived test set (Fig. 2) recapitulates the trends observed in the real-world clinical sample set. PD-FLDL achieves a positive call rate of 84.2%, substantially outperforming ML Baseline (70.2%).

PD-FLDL again outperforms C-FLDL, which in turn outperforms ER-FLDL. C-FLDL (74.3%) again exceeds the ML Baseline, but in contrast to the real-world clinical sample set, here ER-FLDL (66.9%) underperforms relative to the ML Baseline model. Finally, both PD-MaxPool and C-MaxPool struggle in this task and fail to achieve a competitive detection rate on cases with challenging ctDNA levels. These results validate the importance of the MHN attention mechanism for low-signal regimes and demonstrate that while biological prior denoising yields improved FLDL performance (at the current training set sample size), the architecture is capable of effective implicit denoising.

**Ablation: DNA foundation model embedding.** We performed a targeted ablation of the HyenaDNA embedding using the PD-FLDL input space and evaluated the ablated model’s performance on the real-world clinical and challenging contrived hold-out sets at the same 90% specificity operating point shown in Table 1 and Fig. 2. Relative to full-feature PD-FLDL, the variant without the HyenaDNA embedding had 1.3 and 1.0 point lower APL and CRC sensitivities, respectively, and a 2.0 point lower positive call rate on the challenging contrived set. These modest but consistent declines suggest that HyenaDNA embeddings contribute incremental, clinically relevant signal to the FLDL fragment representation.

Taken together, these results show that the FLDL architecture with its multi-modal fragment representation and MHN attention mechanism can a) scale to larger numbers of fragments per sample, b) reduce the need for explicit biological priors, and c) operate on broad genomic panels (the full capture panel), which targets genomic regions selected for use with other cancers, and is thus suitable for multi-cancer detection beyond CRC.

**3.2. Scaling behavior: training dataset size**

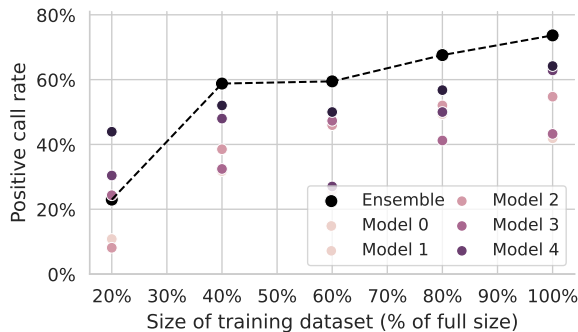


Figure 3. Positive call rates among 148 replicates of the challenging contrived test set, at 90% specificity, for C-FLDL models trained on subsets (20% to 100%) of the full training dataset. Results are shown for the ensemble of 5 models as well as for individual ensemble members (Model 0 to Model 4).

In order to understand the scale of data needed for effective training for the CRC detection task and to identify the onset of potential performance plateaus, we explore the scaling behavior of the FLDL architecture when trained on different dataset sizes. For this, we train the C-FLDL model on nested subsets of our full training dataset and report the resulting models’ positive call rate on the challenging contrived test samples. To preserve the ratio of positives to negatives, we subsample the positive cases (CRCs and APLs) and negative controls separately, to fractions ranging from 20% to 100% in 20% increments. To reduce computational effort during training, for this analysis we only ensemble 5 models from a single round of 5-fold cross validation.

As shown in Fig. 3, at all but the smallest subsampling fraction, the ensemble outperforms all of its 5 individual members, suggesting that ensembling is an effective strategy for smoothing noisy performance of individual members. We also observe a clear improvement in model performance as the training dataset size increases, and we do not observe plateauing, even when approaching the full training dataset size. The increasing performance trend indicates that the FLDL model, as currently parameterized, is likely to continue to benefit from even more training samples, and that there may be potential for further performance improvements by scaling the complexity of the model along with increased training data volume in the future.

### 3.3. Latent space analysis and biological interpretability

For FLDL interpretability at the fragment level, we analyze the model’s attention weights, evaluate genomic localization of high-attention fragments, and compare to prioritized genomic loci of the ML Baseline in Appendix Section E and Appendix Fig. A4.

At the sample level, embeddings after MHN aggregation from C-FLDL are used to visualize sample distributions and qualitatively assess the model’s generalizability to unseen test data. In Fig. 4 we project both training and test set embeddings into a two-dimensional space using UMAP (McInnes et al., 2018) to inspect the latent structure for potential distribution shifts and biological relevance. The UMAP visualization reveals a coherent, biologically relevant continuum: we observe a smooth gradient transitioning from no-ctDNA negative clinical samples (blue) to challenging low-ctDNA clinical blends (orange), to high-ctDNA clinical blends (purple) and real-world clinical APL and CRC samples (red). The artificial no-methylation samples, which lack all CpG methylation and are hence biologically implausible, form a distinct cluster that is well separated from real-world clinical samples and clinical blends. These results demonstrate that the FLDL model learns a robust, biologically meaningful representation that effectively distinguishes among biologically plausible samples based on

their ctDNA content, and clearly separates biologically implausible inputs.

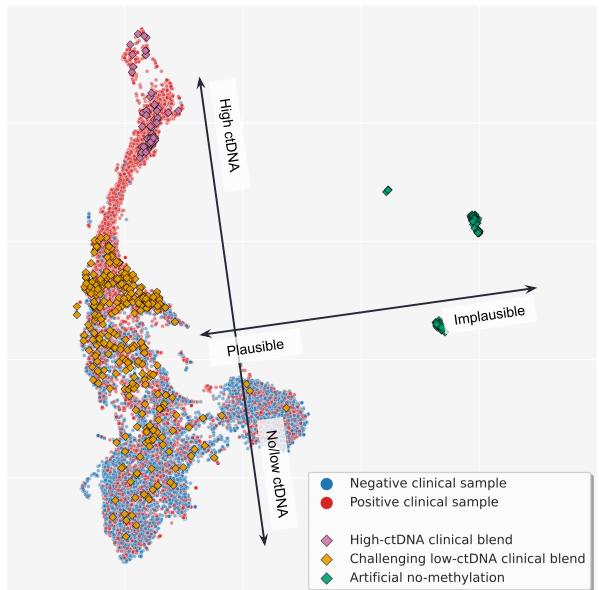


Figure 4. UMAP projection of learned FLDL sample embeddings from the C-FLDL model, applied to real-world clinical, contrived blend, and artificial no-methylation samples.

## 4. Conclusion

In this work, we present Fragment-Level Deep Learning (FLDL), an end-to-end multiple instance learning architecture designed to address the extreme needle-in-a-haystack challenge of early cancer detection from cfDNA in blood. Leveraging Modern Hopfield Networks to attend to rare ctDNA signals amidst millions of uninformative background fragments, FLDL directionally outperforms a state-of-the-art machine learning approach and a max pooling deep learning alternative, particularly at identifying low-signal APLs and challenging contrived samples, with modest but relevant contributions from DNA foundation model embeddings. These results suggest a path for translating DNA foundation model embeddings from genomic representation learning to a clinically meaningful patient-level prediction task. Beyond predictive performance, FLDL also provides interpretable and clinically relevant representations at both fragment and sample levels. In addition, the finding that FLDL can implicitly denoise large genomic search spaces suggests suitability for future multi-cancer detection tasks using a capture panel designed for more than one cancer. Finally, scaling behavior analysis reveals that the model’s predictive performance improves consistently with increasing training data volume, without yet reaching a plateau. As cfDNA-based blood tests are deployed at a larger scale, FLDL is well positioned to benefit from larger future datasets for even greater cancer detection sensitivity.

## Impact Statement

**Impact on deep learning and on diagnostics.** By successfully addressing the extreme-scale multiple instance learning challenge, the FLDL architecture may encourage evaluation of Modern Hopfield Networks and attention mechanisms in a broader class of biological problems with extremely low signal-to-noise ratios. The architecture’s ability to implicitly denoise large genomic search spaces suggests it could be effectively repurposed for cfDNA-based detection of other cancers or diseases. The interpretability of the attention mechanism allows for the identification of high-attention genomic regions, potentially aiding in the discovery or better understanding of disease biomarkers.

**Broader impact on society.** Colorectal cancer remains a critical public health challenge and is the second most common cause of cancer-related death in the US despite being preventable through screening (Siegel et al., 2025). While existing screening methods like colonoscopy and stool-based tests are available, adherence rates remain sub-optimal (Bandi et al., 2025). Blood-based tests offer a non-invasive alternative that can significantly increase screening adherence. A blood test with better sensitivity will be useful in this context, particularly for harder-to-detect cases with low tumor fraction, which is the stage at which the disease is most treatable. The FLDL model improves upon the current state-of-the-art machine learning model for this application, and it demonstrates the potential for further performance gains through training data scaling, pointing towards even greater cancer detection sensitivity in the future.

**Data considerations for bias.** Confounding factors such as age, sex, race/ethnicity, and comorbidities could inadvertently be used by the model for classification if they correlate with the target label in the training data. Our current training set consists of 4,394 samples, a size that is substantial for the clinical diagnostics space but small compared to image or large language model training datasets. Techniques for controlling confounding of this type (e.g., dynamic minibatch balancing) may be less effective in smaller training datasets. Further, there may be unknown confounders. We are encouraged by FLDL’s performance on the fully independent real-world clinical test set. However, continued monitoring and validation in additional diverse, representative datasets will be essential for confirming generalizability. It will also be important for these expanded datasets to increase sample counts for currently underrepresented subpopulations.

**Privacy, data availability.** The clinical data used in this study cannot be shared or distributed because they contain sensitive patient information and high-resolution genomic sequences that are subject to strict privacy regulations. Written informed consent was obtained from each subject who met eligibility criteria and contributed a biosample to this study. All personally identifiable information (PII) in the

source data was fully redacted prior to use in the model development process, to protect patient privacy and adhere strictly to all relevant regulations and guidelines. The de-identification process was approved and monitored by our Data Governance Committee.

## References

- Bandi, P., Star, J., Mazzitelli, N., Nargis, N., Islami, F., Siegel, R. L., Yabroff, K. R., and Jemal, A. Prevalence and review of major modifiable cancer risk factors, hpv vaccination, and cancer screenings in the united states: 2025 update. *Cancer Epidemiology, Biomarkers & Prevention*, 34(6):836–849, 2025.
- Betgegowda, C., Sausen, M., Leary, R. J., Kinde, I., Wang, Y., Agrawal, N., Bartlett, B. R., Wang, H., Luber, B., Alani, R. M., et al. Detection of circulating tumor dna in early-and late-stage human malignancies. *Science translational medicine*, 6(224):224ra24–224ra24, 2014.
- Brixi, G., Durrant, M. G., Ku, J., Naghipourfar, M., Poli, M., Sun, G., Brockman, G., Chang, D., Fanton, A., Gonzalez, G. A., et al. Genome modelling and design across all domains of life with evo 2. *Nature*, pp. 1–13, 2026.
- Carbonneau, M.-A., Cheplygina, V., Granger, E., and Gagnon, G. Multiple instance learning: A survey of problem characteristics and applications. *Pattern recognition*, 77:329–353, 2018.
- Chung, D. C., Gray, D. M., Singh, H., Issaka, R. B., Raymond, V. M., Eagle, C., Hu, S., Chudova, D. I., Talasaz, A., Greenson, J. K., et al. A cell-free dna blood-based test for colorectal cancer screening. *New England Journal of Medicine*, 390(11):973–983, 2024.
- Dalla-Torre, H., Gonzalez, L., Mendoza-Revilla, J., Carranza, N. L., Grzywaczewski, A. H., Oteri, F., Dallago, C., Trop, E., de Almeida, B. P., Sirelkhatim, H., et al. Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods*, 22:287–297, 2025.
- Deng, Z., Ji, Y., Han, B., Tan, Z., Ren, Y., Gao, J., Chen, N., Ma, C., Zhang, Y., Yao, Y., et al. Early detection of hepatocellular carcinoma via no end-repair enzymatic methylation sequencing of cell-free dna and pre-trained neural network. *Genome Medicine*, 15(1):93, 2023.
- Foulds, J. and Frank, E. A review of multi-instance learning assumptions. *The Knowledge Engineering Review*, 25(1): 1–25, 2010.
- Hopfield, J. J. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences*, 79(8):2554–2558, 1982.

- Hopfield, J. J. Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences*, 81(10): 3088–3092, 1984. doi: 10.1073/pnas.81.10.3088.
- Hu, J. Y.-C., Yang, D., Wu, D., Xu, C., Chen, B.-Y., and Liu, H. On sparse modern hopfield model. *Advances in neural information processing systems*, 36:27594–27608, 2023.
- Ilse, M., Tomczak, J., and Welling, M. Attention-based deep multiple instance learning. In *International conference on machine learning*, pp. 2127–2136. PMLR, 2018.
- Jeong, Y., Gerhäuser, C., Sauter, G., Schlomm, T., Rohr, K., and Lutsik, P. Methylbert enables read-level dna methylation pattern identification and tumour deconvolution using a transformer-based model. *Nature Communications*, 16(1):788, 2025.
- Ji, Y., Zhou, Z., Liu, H., and Davuluri, R. V. Dnabert: pre-trained bidirectional encoder representations from transformers for dna-language model. *Bioinformatics*, 37(15):2112–2120, 2021.
- Klambauer, G., Unterthiner, T., Mayr, A., and Hochreiter, S. Self-normalizing neural networks. In *Advances in Neural Information Processing Systems*, pp. 971–980, 2017.
- Li, W., Li, Q., Kang, S., Same, M., Zhou, Y., Sun, C., Liu, C.-C., Matsuoka, L., Sher, L., Wong, W. H., et al. Cancerdetector: ultrasensitive and non-invasive cancer detection at the resolution of individual reads using cell-free dna methylation sequencing data. *Nucleic Acids Research*, 46(15):e89, 2018.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Luo, H., Wei, W., Ye, Z., Zheng, J., and Xu, R. Liquid biopsy of methylation biomarkers in cell-free dna. *Trends in Molecular Medicine*, 27(5):482–500, 2021.
- McInnes, L., Healy, J., and Melville, J. Umap: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018.
- Nguyen, E., Poli, M., Faizi, M., Thomas, A., Wornow, M., Birch-Sykes, C., Massaroli, S., Patel, A., Rabideau, C., Bengio, Y., et al. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. *Advances in neural information processing systems*, 36: 43177–43201, 2023.
- Niki, P., Nalmpantis, C., Ganbat, J.-O., Byrne, D., Babikir, H., Jhutti, A., Rowe, W., Liu, T., Loyfer, N., Toniolo, S., et al. Human whole epigenome modelling for clinical applications with pleiades. *bioRxiv*, pp. 2025–07, 2025.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: an imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pp. 8024–8035, 2019.
- Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlović, M., Sandve, G. K., Greiff, V., Kreil, D., Kopp, M., Klambauer, G., Brandstetter, J., and Hochreiter, S. Hopfield networks is all you need. *ArXiv*, 2008.02217, 2020.
- Shaukat, A., Burke, C. A., Chan, A. T., Grady, W. M., Gupta, S., Katona, B. W., Ladabaum, U., Liang, P. S., Liu, J. J., Putcha, G., et al. Clinical validation of a circulating tumor dna-based blood test to screen for colorectal cancer. *JAMA*, 334(1):56–63, 2025.
- Siegel, R. L., Kratzer, T. B., Giaquinto, A. N., Sung, H., and Jemal, A. Cancer statistics, 2025. *Ca*, 75(1):10, 2025.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Widrich, M. Long short-term memory and convolutional neural networks for SNV-based phenotype prediction. Master’s thesis, JOHANNES KEPLER UNIVERSITY LINZ, 2016.
- Widrich, M., Schäfl, B., Pavlović, M., Ramsauer, H., Gruber, L., Holzleitner, M., Brandstetter, J., Sandve, G. K., Greiff, V., Hochreiter, S., and Klambauer, G. Modern Hopfield networks and attention for immune repertoire classification. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M. F., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 18832–18845. Curran Associates, Inc., 2020.
- Yamaguchi, S., Asao, T., Marcet-Palacios, M., Al-Kasspooles, M., Lee, J., Weitz, J., Ambrosini, G., Ju, J., Wiley, E., and Bland, K. High frequency of dap-kinase gene promoter methylation in colorectal cancer specimens and its identification in serum. *Cancer Letters*, 194(1):99–105, 2003.

## A. Implementation details

### A.1. Details on Fragment Representation and Embedding

Each cfDNA fragment  $s_i$  is represented using five input modalities. **Sequence and CpG Methylation:** the nucleotide sequence and per-CpG methylation status of the fragment, represented as a one-hot encoded matrix with channels for sequenced nucleotides, reference nucleotides, and CpG methylation status. **Methylation Statistics:** a vector of per-fragment summarized methylation statistics, including the sequence length and length-normalized counts of methylated, unmethylated, and total CpGs. **Foundation Model Embedding:** a learned representation based on the HyenaDNA foundation model (Nguyen et al., 2023). This embedding is derived from large-scale training on the human reference genome and captures long-range genomic dependencies and high-order sequence motifs. **Genomic Position:** the numerical location of the fragment start position in a concatenated version of the human reference genome. **Strand:** a boolean indicator for the DNA strand. To ensure a balanced contribution from the heterogeneous data modalities, each modality is projected to a latent embedding of uniform dimension  $m$  using a dedicated modality-specific encoder. Sequence features are processed via a 1D Convolutional Neural Network (CNN) followed by max pooling and a Multi-Layer Perceptron (MLP). The foundation model embeddings are projected to  $m$  dimensions using an MLP. Scalar features such as methylation statistics and genomic position are encoded using triangular encoding (Widrich, 2016) followed by MLPs, while the boolean strand feature is processed directly by an MLP. Lastly, the individual modality embeddings, each of size  $m$ , are concatenated and processed by an MLP to form a unified fragment representation  $h_i \in \mathbb{R}^m$ . This is illustrated in Fig. A1.

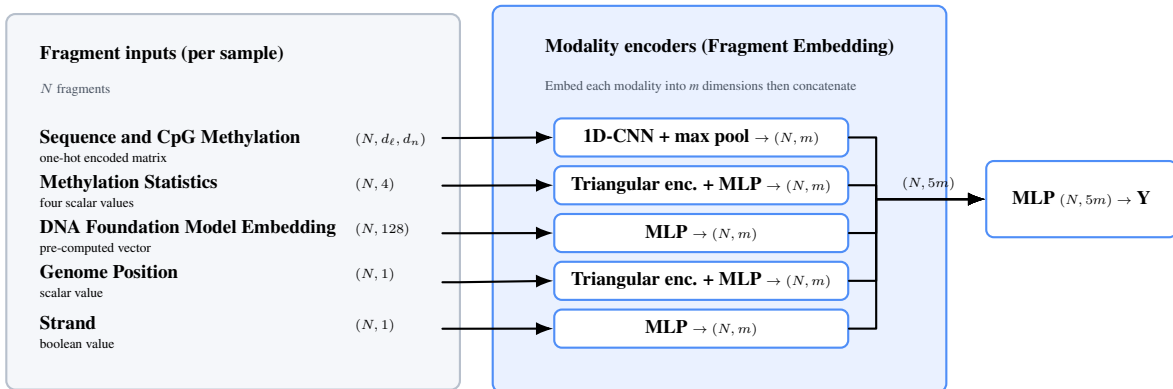


Figure A1. FLDL Fragment Embedding: Each input modality is embedded by a dedicated encoder. After initial embedding, the per-modality embedded vectors are concatenated and processed using an MLP.

### A.2. Implementation, training, ensembling

Training on millions of fragments per sample presents significant computational challenges. We implement a two-stage fragment subsampling strategy to mitigate runtime and GPU memory consumption. During training and inference for early stopping, random dropout of fragments is followed by further attention-based subsampling that retains the fragments with the highest current attention scores. For test-time inference, we omit random dropout and only apply the attention-based subsampling. We train the model end-to-end in PyTorch (Paszke et al., 2019) using the AdamW optimizer (Loshchilov & Hutter, 2017). To prevent exploitation of confounding effects such as collection batch or patient age, we employ a dynamic minibatch balancing scheme that pairs samples with opposite class labels but similar confounding characteristics. To improve generalization, the model is also trained with auxiliary tasks: predicting relevant clinical data and biological characteristics alongside the primary binary disease status.

Due to the relatively low number of training samples and low signal-to-noise ratio (as a result of low witness rate), the representation learned may vary among trained models based on the random weight initialization and the order in which samples are selected for use during training. To address this, we ensemble 25 models obtained from 5 random restarts of 5-fold cross validation (CV). Specifically, we apply 5-fold cross validation with 5 random restarts to obtain 25 members for the final ensemble model (Fig. A2). In each case, 4 different hyperparameter settings are considered, and the optimal hyperparameter choice for that iteration is the one that maximizes the minimum of training and tuning set accuracy. (While uncommon, this maximin metric provides stronger regularization, which is helpful given the small tuning set sample size.) The final model score for a sample is the average of the scores from the 25 models in the ensemble.

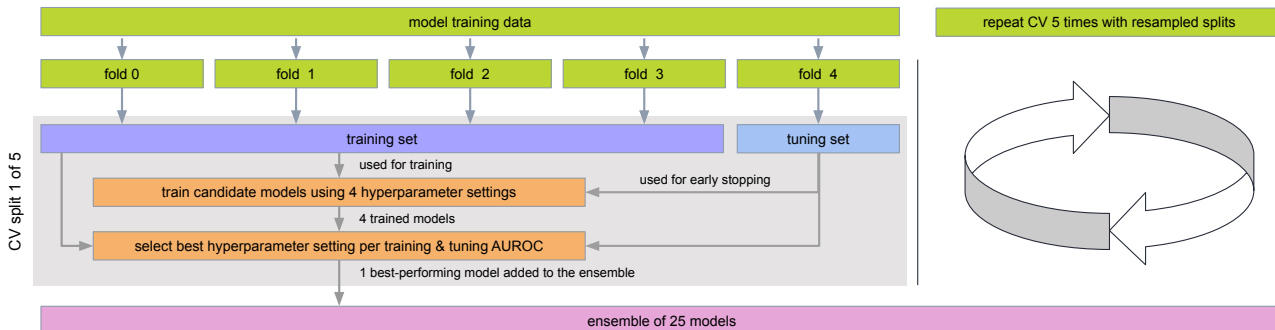


Figure A2. Model training and ensemble construction methodology using 5-fold cross validation with 5 random restarts.

## B. Compared Methods

**FLDL variants.** The wet-lab sample processing pipeline’s full capture panel targets >500 kilobases (kb) of differentially methylated genomic regions, identified through an iterative process leveraging public and internal DNA methylation data. While the full panel captures signals relevant to multiple cancer types, a subset is tailored to maximize CRC signal detection and is referred to as the CRC panel. In order to detect sparse signals necessary for early cancer detection, both the full and CRC panels retain certain loci that can exhibit sporadic background methylation in healthy individuals, which requires models to denoise, i.e., ignore such loci so that non-ctDNA fragments from noisy regions do not lead to false positives. While the machine learning baseline incorporates an explicit denoising step during training (described below), the FLDL architecture does not. FLDL can rely on attention for implicit denoising, and biologically informed pre-filtering can assist in denoising when training data are limited.

To investigate the effect of input space denoising on model performance at the current training dataset size, we evaluate three distinct FLDL configurations:

**Prior Denoised (PD-FLDL):** A variant incorporating a pre-filter step outside of model training to reduce noise in input data. We first profile methylation signal in the CRC panel using a hold-out cohort of 825 healthy individuals. Regions exhibiting elevated background hypermethylation in these controls are aggressively filtered, yielding a reduced, higher-signal-to-noise input space. We refer to this data-driven pre-filtration of the feature space as denoising with *explicit biological priors*.

**CRC Panel (C-FLDL):** The FLDL model operating on the CRC panel without prior filtering. This configuration tests the architecture’s capacity to implicitly denoise the input space and identify cancer-specific signals within the CRC panel, which is smaller than the full panel but still includes noisy regions.

**Expanded Region (ER-FLDL):** A scalability benchmark extending the input space to the full capture panel. This configuration further challenges the model to isolate CRC signals within a substantially larger search space that includes a higher fraction of noisy features.

**Deep learning via max pooling.** To isolate the contribution of the MHN attention mechanism, we evaluate an alternative deep learning model that retains the FLDL fragment feature extractor but replaces Hopfield Pooling with a fragment aggregation strategy based on maximum embedded feature activations. Starting with  $Y = [h_1, \dots, h_N]$  as defined in Section 2, the MaxPool approach constructs a sample embedding  $z \in \mathbb{R}^{m^2}$  by identifying  $m$  representative fragments, one for each fragment embedding dimension. Specifically, for each dimension  $j$ , we determine the fragment  $h_{i^*(j)}$  that maximizes the activation of that dimension:  $i^*(j) = \arg \max_i h_i^j$ . We then define  $r_j = h_{i^*(j)} \in \mathbb{R}^m$ . The final sample embedding  $z$  is obtained by concatenating  $r_1, r_2, \dots, r_m$  to yield  $z \in \mathbb{R}^{m^2}$ . This pooling strategy is motivated by two key considerations: a) it acts as a discrete analogue to attention-based pooling, where the maximum embedded feature activation serves as a proxy for importance; and b) by preserving the intact embedded fragment  $h_{i^*(j)}$  rather than creating a pseudo-fragment from dimension-wise scalar maxima, the model maintains the co-occurrence of features within the representative fragments. This in turn ensures that the sample representation is also constructed from intact embedded fragments. We evaluate this architecture on the Prior Denoised (PD-MaxPool) and CRC (C-MaxPool) input spaces, omitting the Expanded Region configuration due to poor performance in preliminary studies.

**State-of-the-art machine learning model (ML Baseline).** This model operates by identifying hypermethylated fragments

(HMFs) in the CRC panel, where HMFs are defined as cfDNA fragments showing significantly more CpG methylation than typically seen in similar fragments obtained from individuals without disease (Shaukat et al., 2025). In practice, the model first subdivides the CRC panel into small, similarly sized bins. For each bin  $b$ , fragments derived from healthy samples are compared to fragments derived from cases, and a per-bin methylated CpG threshold  $t_b$  is learned. Fragments intersecting bin  $b$  and with a methylated CpG count at or above  $t_b$  are deemed to be hypermethylated. The model aggregates HMF counts observed throughout the regions of interest in order to compute an overall score for the sample. This approach includes an explicit denoising step in the training process: only fragments with methylated CpG counts exceeding that seen in training set control samples are allowed to contribute to the classification score. As for FLDL, an ensemble of 25 such models, each trained on different subsets of the training data, produces the final binary classification.

### C. Dataset details

Our **training set** consists of 4,394 samples, of which 925 are positive cases, comprising 377 CRCs and 548 of the more difficult APLs, and 3,469 are negative controls. To assess classification accuracy on data that are representative of future use cases, we utilize two complementary independent hold-out test sets: a) a **real-world clinical sample set** with 930 samples, containing 331 negative and 599 positive samples (388 APLs, 211 CRCs) collected independently of the training dataset; and b) a **challenging contrived set**, to probe the detection of samples with low witness rates in a controlled setting and with much higher replication. The contrived set consists of 148 replicates created by mixing material from a single advanced CRC donor into plasma from a healthy donor pool to yield a ctDNA level just above the state-of-the-art machine learning model’s detection limit. To inspect the learned representations of the FLDL model’s sample-level embeddings, we leverage two additional hold-out test sets: c) a **high-signal contrived set** consisting of 36 replicates of a blend similar to the challenging contrived set but with ctDNA level well above the machine learning model’s detection limit, and d) an **artificial no-methylation set** consisting of 92 synthetically unmethylated samples, where any residual methylation signal can be attributed solely to technical noise.

### D. Fragment count distributions for different input space denoising strategies

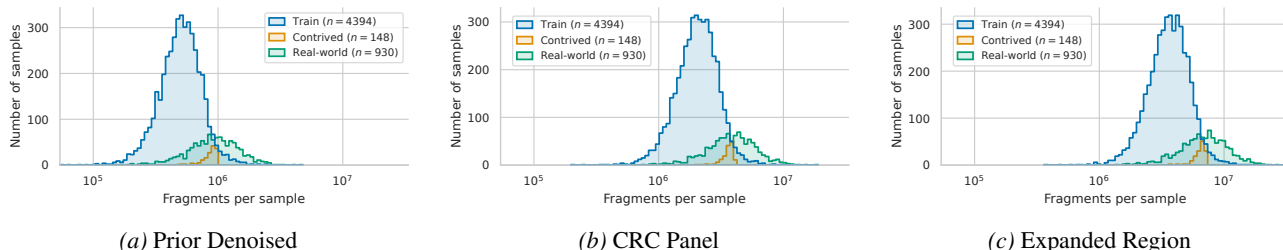


Figure A3. Fragment count distributions for the Prior Denoised (PD), CRC Panel (C), and Expanded Region (ER) training sets vs. the challenging contrived and real-world clinical test sets. Input space size increases from (a) to (b) to (c).

Fig. A3 illustrates the distribution of fragment counts per sample, for training and held-out test datasets, across the three input space configurations described in Sec. 3 and Sec. B. As the model’s input space expands, more fragments are in scope, and we observe the expected rightward shift in the fragment count distributions. The PD-FLDL configuration, which applies aggressive biological filtering to remove noisy regions, yields the smallest input space and the lowest fragment counts. The C-FLDL configuration uses the CRC panel without biological prior filtering, resulting in an intermediate increase in input space and typical fragment count. Finally, the ER-FLDL configuration, which uses the full capture panel, leads to the largest number of instances, with fragment counts almost always exceeding  $10^6$  per sample and sometimes reaching  $10^7$  fragments. Of note, both the real-world clinical and challenging contrived hold-out test sets are generated using an improved wet-lab sample processing pipeline; as a result, they show a marked increase in fragment count per sample compared to the training set. The pipeline improvements are intended to increase accuracy of future blood test versions, but as a side effect, the distribution shift also provides a rigorous testbed for evaluating model generalization in the face of an evolving input distribution.

## E. Attention value analysis for fragment interpretability

For FLDL interpretability at the fragment level, we analyze the model’s attention values. To identify the genomic loci driving the C-FLDL model’s predictions, we analyze the distribution of high-attention fragments within the CRC panel. The starting points for this analysis are the fragments retained after the attention-based subsampling of the Hopfield Pooling layer. To focus on those fragments with the most meaningful contributions to the prediction, we apply a sample-specific filtering step: Given the attention matrix  $\mathbf{A} \in \mathbb{R}^{K \times N}$ , for each of the  $K$  state patterns we retain only those fragments with an attention value above threshold  $t_i = \tau \cdot \max_j(A_i^j)$  where  $i \in \{1, \dots, K\}$ ,  $j \in \{1, \dots, N\}$ , and  $\tau = 0.1$ . Next, we partition the genomic regions interrogated by the CRC panel into  $b$  non-overlapping bins, maintaining consistency with the bins used by the ML Baseline model. The filtered fragments are aligned to these bins, and the relevance count for a bin is incremented for each fragment that at least partially overlaps the bin. This aggregation identifies the genomic regions prioritized by the C-FLDL model for its classification decisions. We compare this to the ML Baseline model’s per-bin HMF counts.

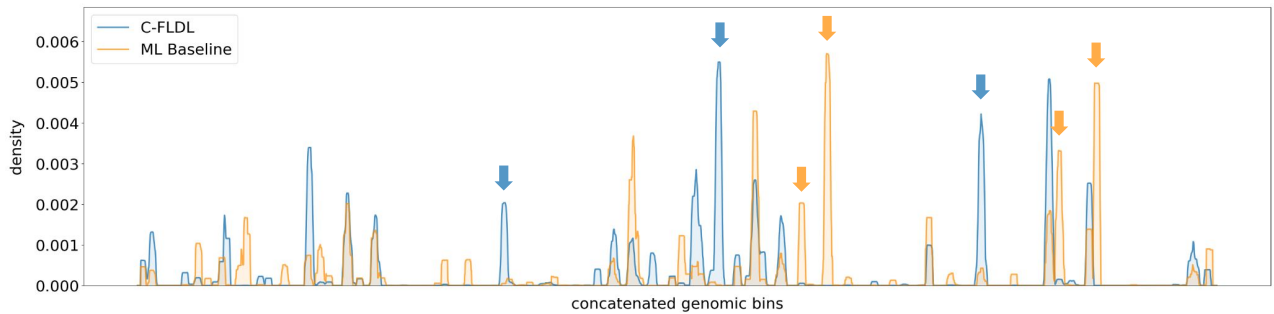


Figure A4. Classification relevance for genomic bins comprising the CRC panel, in 13 correctly predicted late-stage CRC samples from the real-world clinical test set. Blue curve: density of high-attention fragments from C-FLDL. Orange curve: density of HMFs identified by ML Baseline. Arrows indicate regions prioritized by one model but not the other. Both curves are smoothed with a moving average.

Fig. A4 illustrates this comparison using data from 13 correctly classified late-stage CRC samples from the real-world clinical test set. While the models exhibit high-frequency variation at the individual bin level, they demonstrate substantial agreement after smoothing with a moving average, as seen by the overlaps between the orange and the blue curves. Note that a subset of each model’s prioritized genomic loci is ignored by the other model (arrows). This complementarity may suggest that neither model has fully saturated the available signal, or that the models have made different but largely equivalent sparsification choices when presented with sets of loci showing correlated signal.