

---

# Multimodal Cell-Free DNA Embeddings are Informative for Early Cancer Detection

---

## Abstract

1 Cell-free DNA is a promising biomarker for early cancer detection, as it circulates  
2 in the blood and can be extracted non-invasively. However, methods of analysing  
3 the genetic and epigenetic patterns present in cell-free DNA are outdated, and  
4 fail to fully capture the wealth of biological information contained within these  
5 molecules. We present a Transformer based deep learning model that combines the  
6 three distinct modalities contained within cell-free DNA: epigenetic information  
7 in the form of DNA methylation patterns, genetic sequence, and cell-free DNA  
8 fragment length. After training on publicly available data, we demonstrate our  
9 model can accurately distinguish liver cancer patients using cell-free DNA samples  
10 alone. We demonstrate model generalisability by accurate classification of liver  
11 cancer patients from entirely distinct patient cohorts. Finally, we show that the  
12 vector embeddings of cell-free DNA learnt by this multimodal deep-learning model  
13 are biologically informative, and may help shed light on the origins and aetiology  
14 of this elusive bio-molecule.

## 15 1 Introduction

16 DNA is released from somatic cells undergoing apoptosis, and circulates in the peripheral bloodstream,  
17 contributing to what is known as the cell-free DNA (cfDNA) pool. Cell-free DNA molecules carry  
18 three distinct information modalities, all of which vary with cancer status: genomic sequence, a  
19 methylation pattern and the fragment length (as in Figure 1). This molecule can be extracted non-  
20 invasively from blood plasma, following routine blood draws. In cancer, cell cycle dysregulation  
21 leads to increased rates of cell turnover, which causes downstream fluctuations to tumour-specific  
22 cfDNA concentrations detected in the blood.

23 These changes are subtle, but if reliably detected they provide a route to early and non-invasive cancer  
24 diagnosis. Methylation sequencing methods are needed to detect these tissue-specific changes in  
25 cfDNA release, because genomic sequence is identical in all somatic cells whereas DNA methylation  
26 is tissue-specific. Recently, methods have been developed to simultaneously capture methylation,  
27 sequence and fragment length information of cfDNA with high fidelity and depth of coverage,  
28 resulting in a rich dataset describing the genome-wide cell-free DNA state of cancer patients.

29 Current clinical methods of analysing cfDNA are outdated in three important ways. Firstly, they were  
30 designed for low quality, low-depth genomic sequencing data, so tend to aggregate information across  
31 all cfDNA fragments at each genomic locus, disregarding the unique information each individual  
32 bio molecule may contain. Secondly, classification models of cfDNA have relied on manual feature  
33 extraction at known risk loci, which limits research to current genomic hypotheses and scales poorly.  
34 Automated feature extraction methods are more suitable to explore this relatively unknown bio  
35 molecule and its role in cancer aetiology at scale.

36 Finally, there is no way current way to combine the distinct data modalities of cfDNA sequence,  
37 methylation and fragment length. Clinical biomarker tests that rely on cfDNA tend to focus on just  
38 one of these modalities. In doing so, they neglect not only the additional information that could

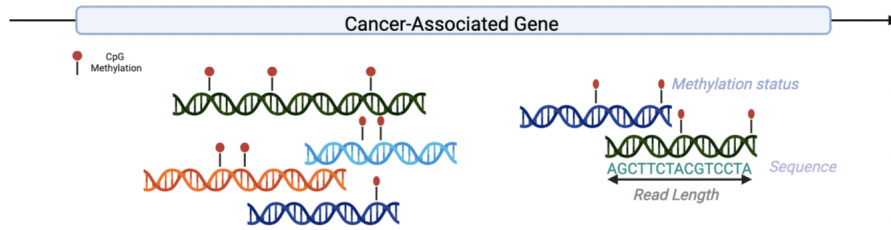


Figure 1: A schematic of cell-free DNA fragments found in peripheral blood, aligned to a gene of interest. These fragments contain three distinct modalities of information: DNA sequence, methylation pattern (shown here on CpGs) and their variable fragment length.

39 be gleaned from this molecule, but also the potential to better understand a key question in cfDNA  
 40 biogenesis: what is the interplay between genomic sequence, methylation status and fragment length?  
 41 Are certain sequence motifs prone more to aberrant methylation during tumourigenesis, and are more  
 42 highly methylated fragments typically shorter? As a result, despite garnering considerable interest as  
 43 an early cancer biomarker, the dynamics and mechanisms of cfDNA fragmentation and release into  
 44 the blood are poorly understood.

45 Here we present a Transformer model that classifies each individual cfDNA fragment, incorporating  
 46 methylation state, genomic sequence and fragment length simultaneously into a single, multimodal  
 47 classifier. We use this model to accurately distinguish liver cancer patients from healthy controls  
 48 across patient cohorts.

## 49 2 Results

50 We encode each individual cfDNA molecule as a variable length vector, where numbers 1-4 represent  
 51 each of the possible nucleotide bases at each position. Methylated cytosine is encoded as a 5, to  
 52 distinguish it from unmethylated cytosine. These encoded molecules are then passed to a Transformer  
 53 deep learning model (standard architecture (Vaswani *et al* 2017; Tay *et al* 2022) with a binary  
 54 classification final layer, whose task is to predict fragment origin: Healthy or Cancer. This model was  
 55 trained on public data, which consisted of billions of individual reads from either Healthy cfDNA  
 56 samples or Hepatocellular Carcinoma tumour samples.

57 For any given cancer patient or healthy control, each individual cfDNA fragment is assigned a  
 58 probability of originating from cancer, as shown in Figure 2. As shown, the probability distribution  
 59 over all fragments is largely similar between healthy and cancer cfDNA samples. This is to be  
 60 expected: cell-free DNA is thought to be released from most healthy somatic tissues, with the  
 61 majority contribution coming from healthy lymphocytes. In cancer patients, only a small subset of  
 62 cfDNA fragments actually originate from the tumour, and even this circulating tumour DNA fraction  
 63 varies considerably with cancer stage and cancer type. A small bump in the probability distribution  
 64 around  $P=0.9$  is seen in cancer cfDNA samples (red) but not in healthy cfDNA samples (green), and  
 65 can be attributed to cfDNA originating from the tumour.

66 We aggregate fragment scores to develop a patient-level risk score, which we then use to classify  
 67 patients. This model is still in active development, and we are currently seeking to improve model  
 68 architecture and hyperparameter selection. Patient cfDNA data collection is also ongoing, but we  
 69 have thus far evaluated the model on the following three datasets:

- 70 1. The held-out test set of unseen cfDNA samples, but taken from the same cohort as the  
 71 training dataset. All patients in the test cohort of this dataset ( $n=24$ ) are correctly classified.
- 72 2. A separate publicly available dataset of cfDNA samples from liver cancer patients and  
 73 healthy controls ( $n=8$ ). All of these patients are also correctly classified by the model,  
 74 demonstrating its generalisability to new patient cohorts.
- 75 3. A third dataset of cfDNA samples from liver cancer patients and healthy controls generated  
 76 in-house ( $n=53$ ). A subset of patients in this cohort have early stage liver cancer, and are

77  
78  
79

incorrectly classified as healthy by our model, most likely due to insufficient cfDNA material originating from cancer. Taking these misclassifications into account, we observe an overall AUC of 0.81 for this patient cohort, and late-stage disease is still detected.

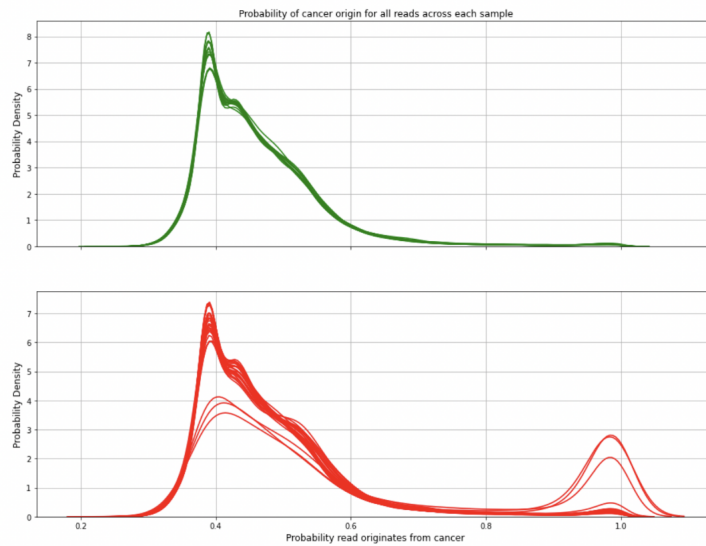


Figure 2: The distribution of cfDNA fragment scores across healthy controls (green) and liver cancer patients (red) in the test set. Each line represents the distribution over all scores for a single patient's cfDNA sample, where the score for a single fragment is the probability that fragment originated from cancer. This data is from patient cohort 1 and 2.

80 By obtaining classifications for each individual fragment in a patient's cfDNA population, we can  
81 begin to untangle the relationship between methylation state, genomic sequence and fragmentation  
82 length of individual cfDNA fragments, and how this affects the predicted origin of the fragment. For  
83 example in Figure 3, we can see that shorter fragments are assigned a lower cancer score on average.  
84 We can also see that as cfDNA fragment methylation rate decreases, the assigned probability that  
85 fragment originates from cancer increases.

86 Population-level cancer screening tests must be non-invasive, and liquid biopsy is the most promising  
87 pan-cancer non-invasive test. As cfDNA-based liquid biopsy diagnostics become clinically adopted,  
88 we desperately need new methods to make sense of the richly informative cfDNA molecules that  
89 circulate in our blood. This method details the initial results of an ongoing data collection and  
90 modelling effort towards this end.

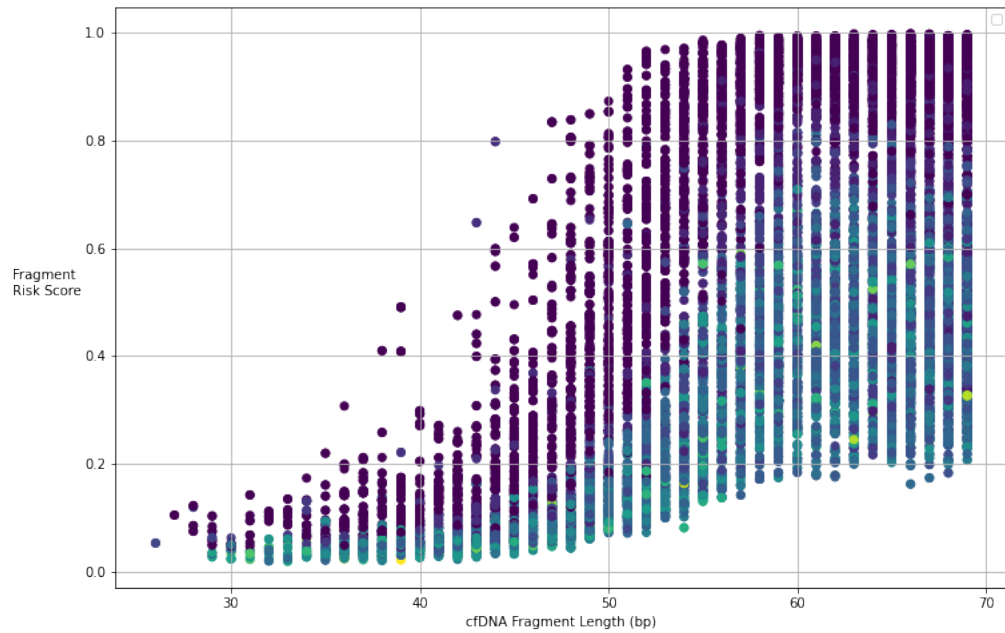


Figure 3: The distribution of cfDNA fragment risk scores stratified by fragment length, for a single liver cancer patient (downsampled for legibility). Methylation state (normalised by fragment length) is represented by marker color, where a lighter color means higher methylation.

91 **References**

- 92 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz  
 93 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing*  
 94 *systems*, pages 5998–6008, 2017
- 95 Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. 2022. Efficient Transformers: A Survey.  
 96 *ACM Comput.Surv.* (apr 2022). Just Accepted.