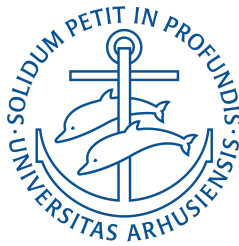


# Generative AI for Predicting Gene Expression from cfDNA Fragmentation Patterns



**Christina Espeseth Grand**

**Supervisor:** Jakob Skou Pedersen

Department of Bioinformatics  
Aarhus University

This thesis is submitted for the degree of  
*Master of Science in Bioinformatics*

June 2025



## Acknowledgements

I would like to express my gratitude to Professor Jakob Skou Pedersen for his invaluable guidance and support throughout this project. His expertise, insightful feedback, and encouragement have played a big role in shaping the direction of my research and in helping me navigate the various challenges and frustrations encountered along the way.

I am also deeply thankful to my co-supervisors, Farhad Zamani and Mathilde Diekema, for their support. Farhad's assistance, especially his insights and the use of code from his master's thesis, has been important in the development of this project. Mathilde's thoughtful advice, support, and willingness to help whenever needed have been essential to my progress.

Finally, I would like to thank everyone in the research group for good discussions and conversations. Their support has made my time here a truly positive and enriching experience.



# Abstract

## Motivation

Nucleosomes regulate gene expression by blocking DNA access to proteins important for transcription, requiring the displacement of nucleosomes at promoter regions for transcription to occur. Their positioning, especially around transcription start sites (TSSs), reflects gene activity. This positioning can be inferred from cell-free DNA (cfDNA) in blood, as nucleosomes protect DNA from cleavage during apoptosis. This protection results in nucleosome-bound DNA being more likely to appear in cfDNA, whereas DNA from actively transcribed genes, particularly around the TSS, are more exposed and thus underrepresented due to degradation. In cfDNA sequencing, this creates an inverse relationship between gene expression and the read depth around the TSS of the corresponding gene. This relationship suggests the possibility of modeling gene expression from cfDNA sequencing, offering potential for noninvasive disease detection based on expression changes.

This study primarily explores deep generative models to predict gene expression from cfDNA sequencing data. XGBoost is used as a baseline model for comparison. Two variants of a Deep Generative Decoder (DGD) are applied: a multimodal DGD and a newly developed factor multimodal DGD. The latter model explicitly separates the contributions of genes and samples to both cfDNA read depths and gene expression. These models aim to find the patterns linking the cfDNA signal to transcriptional activity. This approach can potentially enable noninvasive disease detection by leveraging differences in gene expression between diseased and normal cells.

## Results

Analysis of the data suggested a potential relationship between cfDNA read depth patterns around the TSS and gene expression levels. However, the models struggled to effectively learn this pattern. One of the reasons for this could be significant noise within the data. A key factor contributing to the poor performance of the models is that the cfDNA sequencing data and gene expression data were obtained from different individuals, which means that the gene expression associated with the cfDNA pattern may not always be correct, due to variance between patients.



# Table of contents

<b>Abbreviations</b>	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 Thesis Objectives . . . . .	2
1.3 Scope and Limitation . . . . .	3
1.4 Thesis Outline . . . . .	3
<b>2 Background Theory</b>	<b>5</b>
2.1 DNA, Histones, and Chromatin . . . . .	5
2.1.1 Packaging of DNA . . . . .	5
2.1.2 Regulation of Gene Expression . . . . .	6
2.2 Cell-free DNA . . . . .	8
2.3 Nucleosome Positioning in Relation to cfDNA Fragments . . . . .	9
2.3.1 CfDNA Fragmentation Pattern to Infer Expression . . . . .	10
2.3.2 CfDNA Fragmentation Pattern to Detect Cell Types . . . . .	11
2.4 Machine Learning and Deep Learning . . . . .	11
2.4.1 XGBoost . . . . .	11
2.4.2 Deep Learning and Feed Forward Network . . . . .	12
2.4.3 Training of Neural Networks . . . . .	13
2.4.4 Optimization of Neural Networks . . . . .	14
2.4.5 Variational Autoencoders . . . . .	15
2.4.6 Deep Generative Decoder . . . . .	16
<b>3 Methods and Data</b>	<b>19</b>
3.1 Datasets . . . . .	19
3.1.1 RNA-Sequencing Data . . . . .	19
3.1.2 CfDNA Data . . . . .	21
3.1.3 Train, validation, test split . . . . .	25
3.2 Model Overviews . . . . .	25
3.2.1 Multimodal DGD . . . . .	25

3.2.2	Factor Multi-DGD . . . . .	28
3.2.3	Evaluation of the Models . . . . .	30
3.2.4	Prediction of New Genes for the Multi-DGD . . . . .	31
3.2.5	Prediction on New Samples for Factor multi-DGD . . . . .	31
3.2.6	Hyperparameter tuning and training of XGBoost . . . . .	32
3.3	Measures of Prediction . . . . .	32
3.4	Weights and Biases for Optimization . . . . .	34
3.5	Code Availability . . . . .	35
<b>4</b>	<b>Results</b>	<b>37</b>
4.1	Overview of RNA-seq Data . . . . .	37
4.2	Overview of cfDNA Data . . . . .	40
4.3	XGBoost . . . . .	42
4.4	Multi-DGD . . . . .	44
4.4.1	Evaluation of Model Training . . . . .	45
4.4.2	Predicting of Gene Expression on Test Set . . . . .	46
4.5	Factor Multi-DGD . . . . .	48
4.5.1	Evaluation of Model Training . . . . .	50
4.5.2	Prediction of Gene Expression from cfDNA Read Depth Features on Test Set . . . . .	53
4.6	Evaluating the Effect of Noise on the Model Performance . . . . .	54
<b>5</b>	<b>Discussion</b>	<b>57</b>
5.1	Limitations Due to Uncoupled Data and Dataset Structure . . . . .	57
5.2	Technical Considerations . . . . .	59
5.3	Future Ideas . . . . .	59
<b>6</b>	<b>Conclusion</b>	<b>61</b>
	<b>References</b>	<b>63</b>



# Abbreviations

**cfDNA** cell-free DNA.

**ctDNA** circulating tumor DNA.

**CV** Coefficient of Variation.

**DGD** Deep Generative Decoder.

**FFT** Fast Fourier Transformation.

**FPKM** Fragments Per Kilobase of transcript per Million mapped reads.

**FPM** Fragments Per Million.

**GMM** Gaussian Mixture Model.

**mRNA** messenger RNA.

**NDR** nucleosome-depleted region.

**NFR** nucleosome-free region.

**PBMC** Peripheral Blood Mononuclear Cell.

**PFE** promoter fragmentation entropy.

**PIC** pre-initiation complex.

**RNA-seq** RNA-sequencing.

**rRNA** ribosomal RNA.

**SGD** stochastic gradient descent.

**SVM** support vector machine.

**TER** Transcription End Region.

**TF** transcription factor.

**TSS** Transcription Start Site.

**VAE** Variational Autoencoder.

**WGS** Whole Genome Sequencing.

**WPS** Window Protection Score.

**XGBoost** Extreme Gradient Boosting.