

Master's Thesis

Characterization of ctDNA and Classification Using Machine Learning Methods

Students

VICTORIA HAGELSKJÆR VILSTRØM
Student number 202007866
Exam number 186463

EFTYCHIA DASKALAKI
Student number 202402146
Exam number 227492

Supervisors

JAKOB SKOU PEDERSEN
MATHILDE HARTVIG DIEKEMA

JUNE 2026

Section for Bioinformatics and
Computational Biology
Aarhus University

Contents

Abstract	i
Acknowledgements	ii
Abbreviations	iii
1 Introduction	1
1.1 Background	1
1.2 Thesis objectives	2
1.3 Delimitations	2
1.4 Overview	3
1.5 Contributions	3
2 Background and related work	4
2.1 Fragmentomics	4
2.1.1 Fragment length	4
2.1.2 End motifs	6
2.1.3 Mismatch profile	6
2.2 Machine learning	7
2.2.1 Logistic regression	9
2.2.2 XGBoost	10
2.2.3 SHAP values	11
3 Extraction	13
3.1 Overview	13
3.2 Methods and Materials	13
3.2.1 Data	13
3.2.2 Quality control	13
3.2.3 Extraction pipeline	14
3.2.4 Tumor Probability Model	17
3.3 Results	20
4 Characterization	23
4.1 Overview	23
4.2 Methods	23
4.2.1 Fragment Length	24
4.2.2 End motifs	24
4.2.3 Mismatch rates	25
4.3 Results	32
4.3.1 Fragment Length	32
4.3.2 End motifs	33
4.3.3 Mismatch rate	35

5	Modelling	39
5.1	Overview	39
5.2	Methods	39
5.2.1	Fragment level	40
5.2.2	Sample level	42
5.2.3	Experiment 1: Statistical reliability	43
5.2.4	Experiment 2: Feature necessity	44
5.2.5	Experiment 3: Signal attribution	44
5.2.6	Experiment 4: Prevalence sensitivity	45
5.2.7	Experiment 5: Signal aggregation	45
5.2.8	Experiment 6: Hierarchical generalization	46
5.3	Results	47
5.3.1	Experiment 1: Statistical reliability	47
5.3.2	Experiment 2: Feature necessity	49
5.3.3	Experiment 3: Signal attribution	49
5.3.4	Experiment 4: Prevalence sensitivity	53
5.3.5	Experiment 5: Signal aggregation	56
5.3.6	Experiment 6: Hierarchical generalization	57
6	Discussion	58
6.1	Extraction	58
6.2	Fragment features	58
6.2.1	Fragment length and end motifs	58
6.2.2	Mismatch rate	59
6.3	Modelling	60
6.3.1	Experiments	60
6.3.2	Additional model choices	62
6.4	Limitations	63
6.5	Future work	64
7	Conclusion	65
8	Materials and code availability	66
A	Appendix	67
A.1	Mismatch rate plots before filtering	67
A.2	XGBoost hyperparameters	68
A.3	SHapley Additive exPlanations (SHAP) values for fourth position of end motifs in the XGBoost model with positional end motif encoding	69
A.4	Sample level model cross-validation performance	70

Abstract

Liquid biopsies provide a minimally invasive approach for analyzing tumor derived material in blood plasma. These samples contain circulating tumor DNA (ctDNA), which reflects tumor biology and has potential applications in cancer detection and disease monitoring when distinguished from healthy cell free DNA (cfDNA). However, ctDNA is typically present at low abundance, which limits reliable detection and poses significant analytical challenges. Fragmentomic features have emerged as informative signals of fragment origin that may improve discrimination between ctDNA and healthy cfDNA. These features capture patterns related to DNA fragmentation, providing complementary information to conventional mutation based approaches. In this thesis, we investigated multiple fragmentomic features, including fragment length, end motifs and mismatch rates, to characterize differences between ctDNA and healthy cfDNA fragments and determine predictive power. The results demonstrated consistent differences across several feature types. The ctDNA fragments exhibited shorter fragment lengths across nucleosome associated peaks, altered end motif frequencies and elevated overall mismatch rates. Notably, mismatch related features emerged as a previously underexplored source of fragmentomic signal, highlighting their potential relevance for future liquid biopsy research. A predictive model further confirmed that these features contain information relevant for fragment origin classification, with fragment length and mismatch rates identified as the most important contributors. However, the overall predictive performance indicated that these features alone are insufficient for robust classification.

Acknowledgements

Eftychia

I would like to thank my supervisor Jakob Skou Pedersen for giving me the opportunity to work on this thesis and for his valuable guidance during the semester. His ideas and proposals were really inspiring and motivating. I would also like to thank my daily supervisor, Mathilde Hartvig Diekema for her willingness to provide support to any question during the process. Furthermore, I would like to express my gratitude to my colleague, Victoria Hagelskjær Vilstrøm, for considering me as a thesis collaborator and for our excellent cooperation. Working with her was a fruitful experience, from which I evolved not only scientifically as a bioinformatician but also personally. Lastly, I want to thank my friends and family, who have always been by my side during my journey.

Victoria

I would like to thank my supervisor Jakob Skou Pedersen, for his guidance, constructive feedback and insightful ideas throughout this thesis, which were inspiring and shaped the direction of this work. My appreciation also extends to my daily supervisor Mathilde Hartvig Diekema, for her academic guidance, support and valuable feedback on the report. Furthermore, I am deeply grateful to my collaborator Eftychia Daskalaki. Without her great support, understanding and partnership through difficult times, finishing this thesis would not have been possible. I truly enjoyed working with her. Lastly, a special thanks goes to family for the amazing encouragement and patience they have shown me. Having them by my side has been a constant source of strength. Above all, I want to thank my mother for her kindness, support and encouragement, which has deeply influenced my life. She continues to motivate me every day and this thesis is a tribute to her.