

1     **MitoChime: A Machine Learning Pipeline for**  
2             **Detecting PCR-Induced Chimeras in**  
3             **Mitochondrial Illumina Reads**

4                     A Special Project Proposal  
5                     Presented to  
6     the Faculty of the Division of Physical Sciences and Mathematics  
7                     College of Arts and Sciences  
8                     University of the Philippines Visayas  
9                     Miagao, Iloilo

10                    In Partial Fulfillment  
11                    of the Requirements for the Degree of  
12     Bachelor of Science in Computer Science

13                    by

14                    Duranne Duran  
15                    Yvonne Lin  
16                    Daniella Pailden

17                    Adviser  
18     Francis D. Dimzon, Ph.D.

19                    December 5, 2025

## Abstract

21 Next-generation sequencing (NGS) platforms have advanced research but re-  
22 main susceptible to artifacts such as PCR-induced chimeras that compromise  
23 mitochondrial genome assembly. These artificial hybrid sequences are prob-  
24 lematic for small, circular, and repetitive mitochondrial genomes, where they  
25 can generate fragmented contigs and false junctions. Existing detection tools,  
26 such as UCHIME, are optimized for amplicon-based microbial community ana-  
27 lysis and depend on reference databases or abundance assumptions unsuitable  
28 for organellar assembly. To address this gap, this study presents MitoChime,  
29 a machine learning pipeline for detecting PCR-induced chimeric reads in *Sar-*  
30 *dinella lemuru* Illumina paired-end data without relying on external reference  
31 databases.

32 Using simulated datasets containing clean and chimeric reads, we extracted  
33 a feature set combining alignment-based metrics (e.g., supplementary align-  
34 ments, soft-clipping) with sequence-derived statistics (e.g., k-mer composition,  
35 microhomology). A comparative evaluation of supervised learning models  
36 identified tree-based ensembles CatBoost and Gradient Boosting as top per-  
37 formers, achieving an F1-score of 0.77 and an ROC-AUC of 0.84 on held-out  
38 test data. Feature importance analysis highlighted soft-clipping and k-mer  
39 compositional shifts as the strongest predictors of chimerism, whereas micro-  
40 homology contributed minimally. Integrating MitoChime as a pre-assembly  
41 step can aid in streamlining mitochondrial reconstruction pipelines.

42 **Keywords:** Chimera detection, Mitochondrial genome,  
Assembly, Machine learning

# Contents

44	<b>1 Introduction</b>	<b>1</b>
45	1.1 Overview . . . . .	1
46	1.2 Problem Statement . . . . .	3
47	1.3 Research Objectives . . . . .	4
48	1.3.1 General Objective . . . . .	4
49	1.3.2 Specific Objectives . . . . .	4
50	1.4 Scope and Limitations of the Research . . . . .	5
51	1.5 Significance of the Research . . . . .	6
52	<b>2 Review of Related Literature</b>	<b>7</b>
53	2.1 The Mitochondrial Genome . . . . .	7
54	2.1.1 Mitochondrial Genome Assembly . . . . .	8

55	2.2	PCR Amplification and Chimera Formation . . . . .	9
56	2.3	Existing Traditional Approaches for Chimera Detection . . . . .	10
57	2.3.1	UCHIME . . . . .	11
58	2.3.2	UCHIME2 . . . . .	12
59	2.3.3	CATch . . . . .	13
60	2.3.4	ChimPipe . . . . .	14
61	2.4	Machine Learning Approaches for Chimera and Sequence Quality	
62		Detection . . . . .	15
63	2.4.1	Feature-Based Representations of Genomic Sequences . . .	15
64	2.5	Synthesis of Chimera Detection Approaches . . . . .	16
65	<b>3</b>	<b>Research Methodology</b>	<b>19</b>
66	3.1	Research Activities . . . . .	19
67	3.1.1	Data Collection . . . . .	20
68	3.1.2	Feature Extraction Pipeline . . . . .	24
69	3.1.3	Machine Learning Model Development . . . . .	26
70	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
71		Evaluation . . . . .	28
72	3.1.5	Feature Importance and Interpretation . . . . .	29

73	3.1.6	Validation and Testing . . . . .	30
74	3.1.7	Documentation . . . . .	31
75	3.2	Calendar of Activities . . . . .	32
76	<b>4</b>	<b>Results and Discussion</b>	<b>33</b>
77	4.1	Descriptive Analysis of Features . . . . .	33
78	4.1.1	Exploratory Data Analysis . . . . .	34
79	4.2	Baseline Classification Performance . . . . .	35
80	4.3	Effect of Hyperparameter Tuning . . . . .	37
81	4.4	Detailed Evaluation of Representative Models . . . . .	38
82	4.4.1	Confusion Matrices and Error Patterns . . . . .	39
83	4.4.2	ROC and Precision–Recall Curves . . . . .	40
84	4.5	Feature Importance and Biological Interpretation . . . . .	42
85	4.5.1	Permutation Importance of Individual Features . . . . .	42
86	4.5.2	Feature Family Importance . . . . .	43
87	4.6	Summary of Findings . . . . .	45
88	<b>A</b>	<b>Exploratory Data Analysis</b>	<b>47</b>
89	A.1	Histograms of Key Features . . . . .	47

# 90 List of Figures

91	3.1	Process Diagram of Special Project . . . . .	20
92	4.1	Feature correlation heatmap showing relationships among alignment-	
93		derived and sequence-derived variables. . . . .	35
94	4.2	Test F1 of all baseline classifiers, showing that no single model	
95		clearly dominates and several achieve comparable performance. . .	36
96	4.3	Comparison of test F1 (left) and ROC–AUC (right) for baseline and	
97		tuned models. Hyperparameter tuning yields small but consistent	
98		gains, particularly for tree-based ensembles. . . . .	38
99	4.4	Confusion matrices for the four representative models on the held-	
100		out test set. All models show more false negatives (chimeric reads	
101		called clean) than false positives. . . . .	40
102	4.5	ROC (left) and precision–recall (right) curves for the four represen-	
103		tative models on the held-out test set. Tree-based ensembles cluster	
104		closely, with logistic regression performing slightly but consistently	
105		worse. . . . .	41

106	4.6	Permutation-based feature importance for four representative clas-	
107		sifiers. Clipping and k-mer composition features are generally the	
108		strongest predictors, whereas microhomology and other alignment	
109		metrics contribute minimally. . . . .	43
110	4.7	Aggregated feature family importance across four models. Clipping	
111		and k-mer compositional shifts are consistently the dominant con-	
112		tributors, while SA_structure, Micro_homology, and other features	
113		contribute minimally. . . . .	45
114	A.1	Histogram plots of six key features comparing clean and chimeric	
115		reads. . . . .	48

# 116 List of Tables

<small>117</small>	2.1 Comparison of Chimera Detection Approaches and Tools . . . . .	17
<small>118</small>	3.1 Timetable of Activities . . . . .	32
<small>119</small>	4.1 Performance of baseline classifiers on the held-out test set. . . . .	36
<small>120</small>	4.2 Performance of tuned classifiers on the held-out test set. . . . .	37



# Chapter 1

## Introduction

### 1.1 Overview

The rapid advancement of next-generation sequencing (NGS) technologies has transformed genomic research by enabling high-throughput and cost-effective DNA analysis (Metzker, 2010). Among current platforms, Illumina sequencing remains the most widely adopted, capable of producing millions of short reads that can be assembled into reference genomes or analyzed for genetic variation (Bentley et al., 2008; Glenn, 2011). Despite its high base-calling accuracy, Illumina sequencing is prone to artifacts introduced during library preparation, particularly polymerase chain reaction (PCR)-induced chimeras, which are artificial hybrid sequences that do not exist in the true genome (Judo, Wedel, & Wilson, 1998).

PCR chimeras form when incomplete extension products from one template

anneal to an unrelated DNA fragment and are extended, creating recombinant reads (Qiu et al., 2001). In mitochondrial genome assembly, such artifacts are especially problematic because the mitochondrial genome is small, circular, and often repetitive (Boore, 1999; Cameron, 2014). Even a small number of chimeric or misjoined reads can reduce assembly contiguity and introduce false junctions during organelle genome reconstruction (Dierckxsens, Mardulyn, & Smits, 2017; Hahn, Bachmann, & Chevreux, 2013; Jin et al., 2020). Existing assembly tools such as GetOrganelle and MITObim assume that input reads are largely free of such artifacts (Hahn et al., 2013; Jin et al., 2020). Consequently, undetected chimeras may produce fragmented assemblies or misidentified organellar boundaries. To ensure accurate reconstruction of mitochondrial genomes, a reliable method for detecting and filtering PCR-induced chimeras before assembly is essential.

This study focuses on mitochondrial sequencing data from the genus *Sardinella*, a group of small pelagic fishes widely distributed in Philippine waters. Among them, *Sardinella lemuru* (Bali sardinella) is one of the country’s most abundant and economically important species, providing protein and livelihood to coastal communities (Labrador, Agmata, Palermo, Ravago-Gotanco, & Pante, 2021; Willette, Bognot, Mutia, & Santos, 2011). Accurate mitochondrial assemblies are critical for understanding its population genetics, stock structure, and evolutionary history. However, assembly pipelines often encounter errors or fail to complete due to undetected chimeric reads. To address this gap, this research introduces MitoChime, a machine learning pipeline designed to detect and filter PCR-induced chimeric reads using both alignment-based and sequence-derived statistical features. The tool aims to provide bioinformatics laboratories, partic-

160 ularly the Philippine Genome Center Visayas (PGC Visayas), with an efficient  
161 solution for improving mitochondrial genome reconstruction.

## 162 1.2 Problem Statement

163 While NGS technologies have revolutionized genomic data acquisition, the ac-  
164 curacy of mitochondrial genome assembly remains limited by artifacts produced  
165 during PCR amplification. These chimeric reads can distort assembly graphs and  
166 cause misassemblies, with particularly severe effects in small, circular mitochon-  
167 drial genomes (Boore, 1999; Cameron, 2014). Existing assembly pipelines such  
168 as GetOrganelle, MITObim, and NOVOPlasty assume that sequencing reads are  
169 free of such artifacts (Dierckxsens et al., 2017; Hahn et al., 2013; Jin et al., 2020).  
170 At PGC Visayas, several mitochondrial assemblies have failed or yielded incom-  
171 plete contigs despite sufficient coverage, suggesting that undetected chimeric reads  
172 compromise assembly reliability. Meanwhile, existing chimera detection tools such  
173 as UCHIME and VSEARCH were developed primarily for amplicon-based com-  
174 munity analysis and rely heavily on reference or taxonomic comparisons (Edgar,  
175 Haas, Clemente, Quince, & Knight, 2011; Rognes, Flouri, Nichols, Quince, &  
176 Mahé, 2016). These approaches are unsuitable for single-species organellar data,  
177 where complete reference genomes are often unavailable. Therefore, there is a  
178 pressing need for a reference-independent, data-driven tool capable of detecting  
179 and filtering PCR-induced chimeras in mitochondrial sequencing datasets.

## 180 1.3 Research Objectives

### 181 1.3.1 General Objective

182 This study aims to develop and evaluate a machine learning-based pipeline (Mi-  
183 toChime) that detects PCR-induced chimeric reads in *Sardinella lemuru* mito-  
184 chondrial sequencing data in order to improve the quality and reliability of down-  
185 stream mitochondrial genome assemblies.

### 186 1.3.2 Specific Objectives

187 Specifically, the study aims to:

- 188 1. construct simulated *Sardinella lemuru* Illumina paired-end datasets contain-  
189 ing both clean and PCR-induced chimeric reads,
- 190 2. extract alignment-based and sequence-based features such as k-mer compo-  
191 sition, junction complexity, and split-alignment counts from both clean and  
192 chimeric reads,
- 193 3. train, validate, and compare supervised machine learning models for classi-  
194 fying reads as clean or chimeric,
- 195 4. determine feature importance and identify indicators of PCR-induced  
196 chimerism,
- 197 5. integrate the optimized classifier into a modular and interpretable pipeline  
198 deployable on standard computing environments at PGC Visayas.

## 1.4 Scope and Limitations of the Research

This study focuses on detecting PCR-induced chimeric reads in Illumina paired-end mitochondrial sequencing data from *Sardinella lemuru*. The decision to restrict the taxonomic scope to a single species is based on four considerations: to limit interspecific variation in mitochondrial genome size, GC content, and repetitive regions so that differences in read patterns can be attributed more directly to PCR-induced chimerism; to align the analysis with relevant *S. lemuru* sequencing projects at PGC Visayas; to take advantage of the availability of *S. lemuru* mitochondrial assemblies and raw datasets in public repositories such as the National Center for Biotechnology Information (NCBI), which facilitates reference selection and benchmarking; and to develop a tool that directly supports local studies on *S. lemuru* population structure and fisheries management.

The study emphasizes `wgsim`-based simulations and selected empirical mitochondrial datasets from *S. lemuru*. It excludes naturally occurring chimeras, nuclear mitochondrial pseudogenes (NUMTs), and large-scale assembly rearrangements in nuclear genomes. Feature extraction is restricted to low-dimensional alignment and sequence statistics, such as k-mer frequency profiles, GC content, read length, soft and hard clipping metrics, split-alignment counts, and mapping quality, rather than high-dimensional deep learning embeddings. This design keeps model behaviour interpretable and ensures that the pipeline can be run on standard workstations at PGC Visayas. Testing on long-read platforms (e.g., Nanopore, PacBio) and other taxa is outside the scope of this project; the implemented pipeline is evaluated only on short-read *S. lemuru* datasets.

Other limitations in this study include the following: simulations with varying

error rates were not performed, so the effect of different sequencing errors on model performance remains unexplored; alternative parameter settings, including k-mer lengths and microhomology window sizes, were not systematically tested, which could affect the sensitivity of both k-mer and microhomology feature detection as well as the identification of chimeric junctions; and the machine learning models rely on supervised training with labeled examples, which may limit their ability to detect novel or unexpected chimeric patterns.

## 1.5 Significance of the Research

This research provides both methodological and practical contributions to mitochondrial genomics and bioinformatics. First, MitoChime detects PCR-induced chimeric reads prior to genome assembly, with the goal of improving the contiguity and correctness of *Sardinella lemuru* mitochondrial assemblies. Second, it replaces informal manual curation with a documented workflow, improving automation and reproducibility. Third, the pipeline is designed to run on computing infrastructures commonly available in regional laboratories, enabling routine use at facilities such as PGC Visayas. Finally, more reliable mitochondrial assemblies for *S. lemuru* provide a stronger basis for downstream applications in the field of fisheries and genomics.

## Chapter 2

### Review of Related Literature

This chapter presents an overview of the literature relevant to the study. It discusses the biological and computational foundations underlying mitochondrial genome analysis and assembly, as well as existing tools, algorithms, and techniques related to chimera detection and genome quality assessment. The chapter aims to highlight the strengths, limitations, and research gaps in current approaches that motivate the development of the present study.

#### 2.1 The Mitochondrial Genome

Mitochondrial genome (mtDNA) is a small, typically circular molecule found in most eukaryotes. It encodes essential genes involved in oxidative phosphorylation and energy metabolism. Because of its conserved structure, mtDNA has become a valuable genetic marker for studies in population genetics and phylogenetics (Anderson et al., 1981; Boore, 1999). In animal species, the mitochondrial genome

255 ranges from 15–20 kilobase and contains 13 protein-coding genes, 22 tRNAs, and  
256 two rRNAs arranged compactly without introns (Gray, 2012). In comparison to  
257 nuclear DNA, the ratio of the number of copies of mtDNA is higher and has  
258 simple organization which make it particularly suitable for genome sequencing  
259 and assembly studies (Dierckxsens et al., 2017).

### 260 **2.1.1 Mitochondrial Genome Assembly**

261 Mitochondrial genome assembly refers to the reconstruction of the complete mito-  
262 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is  
263 conducted to obtain high-quality, continuous representations of the mitochondrial  
264 genome that can be used for a wide range of analyses, including species identi-  
265 fication, phylogenetic reconstruction, evolutionary studies, and investigations of  
266 mitochondrial diseases. Because mtDNA evolves rapidly, its assembled sequence  
267 provides valuable insights into population structure, lineage divergence, and adap-  
268 tive evolution across taxa (Boore, 1999). Compared to nuclear genome assembly,  
269 assembling the mitochondrial genome is often considered more straightforward but  
270 still encounters technical challenges such as the formation of chimeric reads. Com-  
271 monly used tools for mitogenome assembly such as GetOrganelle and MITObim  
272 operate under the assumption of organelle genome circularity, and are vulnerable  
273 when chimeric reads disrupt this circular structure, resulting in assembly errors  
274 (Hahn et al., 2013; Jin et al., 2020).



## 2.2 PCR Amplification and Chimera Formation

PCR plays an important role in NGS library preparation, as it amplifies target DNA fragments for downstream analysis. However as previously mentioned, the amplification process can also introduce chimeric reads which compromises the quality of the input reads supplied to sequencing or assembly workflows. Chimeras typically arise when incomplete extension occurs during a PCR cycle. This causes the DNA polymerase to switch from one template to another and generate hybrid recombinant molecules (Judo et al., 1998). Artificial chimeras are produced through such amplification errors, whereas biological chimeras occur naturally through genomic rearrangements or transcriptional events.

In the context of amplicon-based sequencing, the presence of chimeras can inflate estimates of genetic or microbial diversity and may cause misassemblies during genome reconstruction. Qin et al. (2023) has reported that chimeric sequences may account for more than 10% of raw reads in amplicon datasets. This artifact tends to be most prominent among rare operational taxonomic units (OTUs) or singletons, which are sometimes misinterpreted as novel diversity, further causing the complication of microbial diversity analyses (Gonzalez, Zimmermann, & Saiz-Jimenez, 2004). As such, determining and minimizing PCR-induced chimera formation is vital for improving the quality of mitochondrial genome assemblies, and ensuring the reliability of amplicon sequencing data.

## 2.3 Existing Traditional Approaches for Chimera Detection

Several computational tools have been developed to identify chimeric sequences in NGS datasets. These tools generally fall into two categories: reference-based and de novo approaches. Reference-based chimera detection, also known as database-dependent detection, is one of the earliest and most widely used computational strategies for identifying chimeric sequences in amplicon-based community studies. These methods rely on the comparison of each query sequence against a curated, high-quality database of known, non-chimeric reference sequences (Edgar et al., 2011).

On the other hand, the de novo chimera detection, also referred to as reference-free detection, represents an alternative computational paradigm that identifies chimeric sequences without reliance on external reference databases. This method infer chimeras based on internal relationships among the sequences present within the dataset itself, making it particularly advantageous in studies of under explored or taxonomically diverse communities where comprehensive reference databases are unavailable or incomplete (Edgar, 2016; Edgar et al., 2011). The underlying assumption on this method is that during PCR, true biological sequences are generally more abundant as they are amplified early and dominate the read pool, whereas chimeric sequences appear later and are generally less abundant. The de novo approach leverage this abundance hierarchy, treating the most abundant sequences as supposed parents and testing whether less abundant sequences can be reconstructed as mosaics of these templates. Compositional and structural similarity are also evaluated to check whether different regions of a candidate

319 sequence correspond to distinct high-abundance sequences.

320 In practice, many modern bioinformatics pipelines combine both paradigms  
321 sequentially: an initial de novo step identifies dataset-specific chimeras, followed  
322 by a reference-based pass that removes remaining artifacts relative to established  
323 databases (Edgar, 2016). These two methods of detection form the foundation of  
324 tools such as UCHIME and later UCHIME2.

### 325 **2.3.1 UCHIME**

326 UCHIME is one of the most widely used tools for detecting chimeric sequences in  
327 amplicon-based studies and remains a standard quality-control step in microbial  
328 community analysis. Its core strategy is to test whether a query sequence ( $Q$ ) can  
329 be explained as a mosaic of two parent sequences, ( $A$  and  $B$ ), and to score this  
330 relationship using a structured alignment model (Edgar et al., 2011).

331 In reference mode, UCHIME divides the query into several segments and maps  
332 them against a curated database of non-chimeric sequences. Candidate parents  
333 are identified, and a three-way alignment is constructed. The algorithm assigns  
334 “Yes” votes when different segments of the query match different parents and  
335 “No” votes when the alignment contradicts a chimeric pattern. The final score  
336 reflects the balance of these votes. In de novo mode, UCHIME operationalizes the  
337 abundance-skew principle described earlier: high-abundance sequences are treated  
338 as candidate parents, and lower-abundance sequences are evaluated as potential  
339 mosaics. This makes the method especially useful when no reliable reference  
340 database exists.

341 Although UCHIME is highly sensitive, it faces key constraints. Chimeras  
342 formed from parents with very low divergence (below 0.8%) are difficult to detect  
343 because they are nearly indistinguishable from sequencing errors. Accuracy in ref-  
344 erence mode depends strongly on database completeness, while de novo detection  
345 assumes that true parents are both present and sufficiently more abundant, such  
346 conditions are not always met.

### 347 **2.3.2 UCHIME2**

348 UCHIME2 extends the original algorithm with refinements tailored for high-  
349 resolution sequencing data. One of its major contributions is a re-evaluation  
350 of benchmarking practices. Edgar (2016) demonstrated that earlier accuracy es-  
351 timates for chimera detection were overly optimistic because they relied on un-  
352 realistic scenarios where all true parent sequences were assumed to be present.  
353 Using the more rigorous CHSIMA benchmark, UCHIME2 showed the prevalence  
354 of “fake models” or real biological sequences that can be perfectly reconstructed  
355 as apparent chimeras of other sequences, which suggests that perfect chimera de-  
356 tection is theoretically unattainable. UCHIME2 also introduces several preset  
357 modes (e.g., denoised, balanced, sensitive, specific, high-confidence) designed to  
358 tune sensitivity and specificity depending on dataset characteristics. These modes  
359 allow users to adjust the algorithm to the expected noise level or analytical goals.

360 Despite these improvements, UCHIME2 must be applied with caution. The  
361 author’s website manual (Edgar, n.d) explicitly advises against using UCHIME2  
362 as a standalone chimera-filtering step in OTU clustering or denoising workflows  
363 because doing so can inflate both false positives and false negatives.

### 364 2.3.3 CATCh

365 As previously mentioned, UCHIME (Edgar et al., 2011) relied on alignment-based  
366 sequences in amplicon data. However, researchers soon observed that different al-  
367 gorithms often produced inconsistent predictions. A sequence might be identified  
368 as chimeric by one tool but classified as non-chimeric by another, resulting in  
369 unreliable filtering outcomes across studies.

370 To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieus  
371 (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-  
372 resents the first ensemble machine learning system designed for chimera detection  
373 in 16S rRNA amplicon sequencing. Rather than depending on a single detec-  
374 tion strategy, CATCh integrates the outputs of several established tools, includ-  
375 ing UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual  
376 scores and binary decisions generated by these tools are used as input features for  
377 a supervised learning model. The algorithm employs a Support Vector Machine  
378 (SVM) with a Pearson VII Universal Kernel (PUK) to determine optimal weight-  
379 ings among the input features and to assign each sequence a probability of being  
380 chimeric.

381 Benchmarking in both reference-based and de novo modes demonstrated signif-  
382 icant performance improvements. CATCh achieved sensitivities of approximately  
383 85 percent in reference-based mode and 92 percent in de novo mode, with corre-  
384 sponding specificities of approximately 96 percent and 95 percent. These results  
385 indicate that CATCh detected 7 to 12 percent more chimeras than any individual  
386 algorithm while maintaining high precision.

### 387 2.3.4 ChimPipe

388 Among the available tools for chimera detection, ChimPipe is a pipeline developed  
389 to identify chimeric sequences such as biological chimeras. It uses both discordant  
390 paired-end reads and split-read alignments to improve the accuracy and sensitivity  
391 of detecting biological chimeras (Rodriguez-Martin et al., 2017). By combining  
392 these two sources of information, ChimPipe achieves better precision than meth-  
393 ods that depend on a single type of indicator.

394 The pipeline works with many eukaryotic species that have available genome  
395 and annotation data (Rodriguez-Martin et al., 2017). It can also predict multiple  
396 isoforms for each gene pair and identify breakpoint coordinates that are useful  
397 for reconstructing and verifying chimeric transcripts. Tests using both simulated  
398 and real datasets have shown that ChimPipe maintains high accuracy and reliable  
399 performance.

400 ChimPipe lets users adjust parameters to fit different sequencing protocols or  
401 organism characteristics. Experimental results have confirmed that many chimeric  
402 transcripts detected by the tool correspond to functional fusion proteins, demon-  
403 strating its utility for understanding chimera biology and its potential applications  
404 in disease research (Rodriguez-Martin et al., 2017).

## 405 **2.4 Machine Learning Approaches for Chimera** 406 **and Sequence Quality Detection**

407 Traditional chimera detection tools rely primarily on heuristic or alignment-based  
408 rules. Recent advances in machine learning (ML) have demonstrated that models  
409 trained on sequence-derived features can effectively capture compositional and  
410 structural patterns in biological sequences. Although most existing ML systems  
411 such as those used for antibiotic resistance prediction, taxonomic classification,  
412 or viral identification are not specifically designed for chimera detection, they  
413 highlight how data-driven models can outperform similarity-based heuristics by  
414 learning intrinsic sequence signatures. In principle, ML frameworks can integrate  
415 indicators such as k-mer frequencies, GC-content variation and split-alignment  
416 metrics to identify subtle anomalies that may indicate a chimeric origin (Arango  
417 et al., 2018; Liang, Bible, Liu, Zou, & Wei, 2020; Ren et al., 2020).

### 418 **2.4.1 Feature-Based Representations of Genomic Se-** 419 **quences**

420 Feature extraction converts DNA sequences into numerical representations suit-  
421 able for machine learning models. One approach is k-mer frequency analysis,  
422 which counts short nucleotide sequences within a read (Vervier, Mahé, Tournoud,  
423 Veyrieras, & Vert, 2015). High-frequency k-mers, including simple repeats such  
424 as “AAAAAA,” can highlight repetitive or unusual regions that may occur near  
425 chimeric junctions. Comparing k-mer patterns across adjacent parts of a read can  
426 help identify such regions, while GC content provides an additional descriptor of

427 local sequence composition (Ren et al., 2020).

428 Alignment-derived features further inform junction detection. Long-read tools  
429 such as Sniffles (Sedlazeck et al., 2018) use split alignments to locate breakpoints  
430 across extended sequences, whereas short-read aligners like Minimap2 (Li, 2018)  
431 report supplementary and secondary alignments that indicate local discontinu-  
432 ities. Split alignments, where parts of a read map to different regions, can reveal  
433 template-switching events. These features complement k-mer profiles and en-  
434 hance detection of potentially chimeric reads, even in datasets with incomplete  
435 references.

436 Microhomology, or short sequences shared between adjacent segments, is an-  
437 other biologically meaningful feature. Its length, typically a few to tens of base  
438 pairs, has been linked to microhomology-mediated repair and template-switching  
439 mechanisms (Sfeir & Symington, 2015). In PCR-induced chimeras, short iden-  
440 tical sequences at junctions provide a clear signature of chimerism. Measuring  
441 the longest exact overlap at each breakpoint complements k-mer and alignment  
442 features and helps identify reads that are potentially chimeric.

## 443 2.5 Synthesis of Chimera Detection Approaches

444 To provide an integrated overview of the literature discussed in this chapter, Ta-  
445 ble 2.1 summarizes the major chimera detection studies, their methodological  
446 approaches, and their known limitations.



Table 2.1: Comparison of Chimera Detection Approaches and Tools

Method / Tool	Core Approach	Key Limitations
<b>Reference-based Detection</b>	Compares each query sequence against curated databases of verified, non-chimeric sequences; evaluates segment similarity to identify mosaic patterns.	Accuracy depends on database completeness; performs poorly for novel taxa or missing parents; limited sensitivity for low-divergence chimeras.
<b>De novo Detection</b>	Identifies chimeras using only internal dataset structure; leverages abundance hierarchy and compositional similarity to infer whether low-abundance sequences can be reconstructed from abundant parents.	Assumes true sequences are more abundant; fails when amplification bias distorts abundances; struggles when parental sequences are similarly abundant or highly similar.
<b>UCHIME</b>	Alignment-based model that partitions the query into segments, identifies parent candidates, and computes a chimera score via a three-way alignment; supports reference and de novo modes.	Reduced accuracy for very closely related parents (<0.8% divergence); sensitive to incomplete databases; de novo mode fails if parents are absent or not sufficiently more abundant.
<b>UCHIME2</b>	Updated UCHIME with improved benchmarking (CHSIMA) and multiple sensitivity/specificity presets; better handles incomplete references and dataset variability.	“Fake models” limit theoretical accuracy; genuine variants may mimic chimeras; not recommended as a standalone step in OTU or denoising pipelines due to increased false positives/negatives.
<b>CATCh</b>	First ensemble ML model for 16S chimera detection; integrates outputs of UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus using an SVM to boost overall prediction accuracy.	Performance constrained by underlying tools; ML model cannot capture features not present in component algorithms; may misclassify in highly novel or low-coverage datasets.
<b>ChimPipe</b>	Pipeline for detecting biological chimeras in RNA-seq using discordant paired-end reads and split-read alignments; identifies isoforms and breakpoint coordinates.	Requires high-quality genome and annotation; tailored to RNA-seq rather than amplicons; computationally intensive; limited to organisms with available reference genomes.

447 Across existing studies, no single approach reliably detects all forms of chimeric  
448 sequences, and the reviewed literature consistently shows that chimeras remain a  
449 persistent challenge in genomics and bioinformatics. Although the surveyed tools  
450 are not designed specifically for organelle genome assembly, they provide valu-  
451 able insights into which methodological strategies are effective and where current  
452 approaches fall short. These limitations collectively define a clear research gap:  
453 the need for a specialized, feature-driven detection framework tailored to PCR-  
454 induced mitochondrial chimeras. Addressing this gap aligns with the research  
455 objective outlined in Section 1.3, which is to develop and evaluate a machine  
456 learning-based pipeline (MitoChime) that improves the quality of downstream  
457 mitochondrial genome assembly. In support of this aim, the subsequent chapters  
458 describe the design, implementation, and evaluation of the proposed tool.

## 459 Chapter 3

# 460 Research Methodology

461 This chapter outlines the steps involved in completing the study, including data  
462 gathering, generating simulated mitochondrial Illumina reads, preprocessing and  
463 indexing the data, developing a feature extraction pipeline to extract key features,  
464 applying machine learning algorithms for chimera detection, and validating and  
465 comparing model performance.

## 466 3.1 Research Activities

467 As illustrated in Figure 3.1, this study carried out a sequence of procedures to  
468 detect PCR-induced chimeric reads in mitochondrial genomes. The process began  
469 with collecting a mitochondrial reference sequence of *Sardinella lemuru* from the  
470 National Center for Biotechnology Information (NCBI) database, which was used  
471 as a reference for generating simulated clean and chimeric reads. These reads  
472 were subsequently indexed and mapped. The resulting collections then passed

473 through a feature extraction pipeline that extracted k-mer profiles, supplementary  
 474 alignment (SA) features, and microhomology information to prepare the data for  
 475 model construction. The machine learning model was trained using the processed  
 476 input, and its precision and accuracy were assessed. It underwent tuning until it  
 477 reached the desired performance threshold, after which it proceeded to validation  
 478 and will undergo testing.

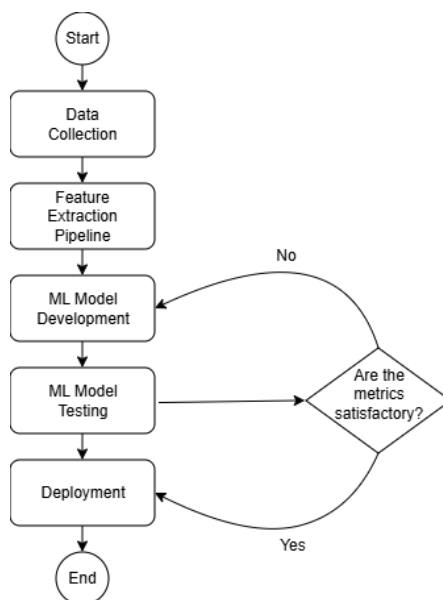


Figure 3.1: Process Diagram of Special Project

### 479 3.1.1 Data Collection

480 The mitochondrial genome reference sequence of *S. lemuru* was obtained from the  
 481 NCBI database (accession number NC\_039553.1) in FASTA format. This sequence  
 482 served as the basis for generating simulated reads for model development.

483 This step was scheduled to begin in the first week of November 2025 and  
 484 expected to be completed by the end of that week, with a total duration of ap-

485 proximately one (1) week.

## 486 Data Preprocessing

487 To reduce manual repetition, all steps in the simulation and preprocessing pipeline  
488 were executed using a custom script in Python (Version 3.11). The script runs  
489 each stage, including read simulation, reference indexing, mapping, and alignment  
490 processing, in a fixed sequence.

491 Sequencing data were simulated from the NCBI reference genome using `wgsim`  
492 (Version 1.13). First, a total of 10,000 paired-end fragments were simulated,  
493 producing 20,000 reads (10,000 forward and 10,000 reverse) from the the original  
494 reference (`original_reference.fasta`) and and designated as clean reads using  
495 the command:

```
496 wgsim -1 150 -2 150 -r 0 -R 0 -X 0 -e 0.001 -N 10000 \  
497     original_reference.fasta ref1.fastq ref2.fastq
```

498 The command parameters are as follows:

- 499 • `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.
- 500 • `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability,  
501 all set to a default value of 0.
- 502 • `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 503 • `-N`: number of read pairs, set to 10,000.

504 Chimeric sequences were then generated from the same NCBI reference using a  
505 separate Python script. Two non-adjacent segments were randomly selected such  
506 that their midpoint distances fell within specified minimum and maximum thresh-  
507 olds. The script attempts to retain microhomology, or short identical sequences  
508 at segment junctions, to mimic PCR-induced template switching. The resulting  
509 chimeras were written to `chimera_reference.fasta`, with headers recording seg-  
510 ment positions and microhomology length. The `chimera_reference.fasta` was  
511 processed with `wgsim` to simulate 10,000 paired-end fragments, generating 20,000  
512 chimeric reads (10,000 forward reads in `chimeric1.fastq` and 10,000 reverse reads  
513 in `chimeric2.fastq`) using the command format.

514 Next, a `minimap2` index of the reference genome was created using:

```
515 minimap2 -d ref.mmi original_reference.fasta
```

516 Minimap2 (Version 2.28) is a tool used to map reads to a reference genome.  
517 The index `ref.mmi` of the original reference sequence is required by `minimap2` for  
518 efficient read mapping. Mapping allows extraction of alignment features from each  
519 read, which were used as input for the machine learning model. The simulated  
520 clean and chimeric reads were then mapped to the reference index as follows:

```
521 minimap2 -ax sr -t 8 ref.mmi ref1.fastq ref2.fastq > clean.sam
```

```
522 minimap2 -ax sr -t 8 ref.mmi \  
523 chimeric1.fastq chimeric2.fastq > chimeric.sam
```

524 Here, `-ax sr` specifies short-read alignment mode, and `-t 8` uses 8 CPU

525 threads. The resulting clean and chimeric SAM files contain the alignment posi-  
526 tions of each read relative to the original reference genome.

527 The SAM files were then converted to BAM format, sorted, and indexed using  
528 **samtools** (Version 1.20):

```
529 samtools view -bS clean.sam -o clean.bam  
530 samtools view -bS chimeric.sam -o chimeric.bam  
531  
532 samtools sort clean.bam -o clean.sorted.bam  
533 samtools index clean.sorted.bam  
534  
535 samtools sort chimeric.bam -o chimeric.sorted.bam  
536 samtools index chimeric.sorted.bam
```

537 BAM files are the compressed binary version of SAM files, which enables faster  
538 processing and reduced storage. Sorting arranges reads by genomic coordinates,  
539 and indexing allows detection of SA as a feature for the machine learning model.

540 The total number of simulated reads was expected to be 40,000. The final col-  
541 lection of reads contained 19,984 clean reads and 20,000 chimeric reads (39,984 en-  
542 tries in total), providing a roughly balanced distribution between the two classes.  
543 After alignment with **minimap2**, only 19,984 clean reads remained because un-  
544 mapped reads were not included in the BAM file. Some sequences failed to align  
545 due to the 5% error rate defined during **wgsim** simulation, which produced mis-  
546 matches that caused certain reads to fall below the aligner's matching threshold.

547 This whole process was scheduled to start in the second week of November 2025

548 and was expected to be completed by the last week of November 2025, with a total  
549 duration of approximately three (3) weeks.

### 550 **3.1.2 Feature Extraction Pipeline**

551 This stage directly follows the previous alignment phase, utilizing the resulting  
552 BAM files (specifically `chimeric.sorted.bam` and `clean.sorted.bam`). A custom  
553 Python script was created to efficiently process each primary-mapped read to  
554 extract the necessary set of analytical features, which are then compiled into a  
555 structured feature matrix in TSV format. The pipeline's core functionality relies on  
556 libraries, namely `Pysam` (Version 0.22) for the robust parsing of BAM structures and  
557 `NumPy` (Version 1.26) for array operations and computations. To ensure correctness  
558 and adherence to best practices, bioinformatics experts at the PGC Visayas will  
559 be consulted to validate the pipeline design, feature extraction logic, and overall  
560 data integrity. This stage of the study was scheduled to begin in the last week  
561 of November 2025 and conclude by the first week of December 2025, with an  
562 estimated total duration of approximately two (2) weeks.

563 The pipeline focuses on three features that collectively capture biological sig-  
564 natures associated with PCR-induced chimeras: (1) Supplementary alignment flag  
565 (SA count), (2) k-mer composition difference, and (3) microhomology.

#### 566 **Supplementary Alignment Flag**

567 Split-alignment information was derived from the SA (Supplementary Alignment)  
568 tag embedded in each primary read of the BAM file. This tag is typically asso-



569 ciated with reads that map to multiple genomic locations, suggesting a chimeric  
570 structure. To extract this information, the script first checked whether the read  
571 carried an `SA:Z` tag. If present, the tag string was parsed using the function  
572 `parse_sa_tag`, yielding a structure for each alignment containing the reference  
573 name, mapped position, strand, mapping quality, and number of mismatches.

574 After parsing, the function `sa_feature_stats` was applied to establish the fun-  
575 damental split indicators, `has_sa` and `sa_count`. Along with these initial counts,  
576 the function synthesized a summarization by aggregating metrics related to the  
577 structure and reliability of the split alignments.

## 578 **K-mer Composition Difference**

579 Chimeric reads often comprise fragments from distinct genomic regions, resulting  
580 in a compositional discontinuity between segments. Comparing k-mer frequency  
581 profiles between the left and right halves of a read allows for the detection of such  
582 abrupt compositional shifts, independent of alignment information.

583 The script implemented this by inferring a likely junction breakpoint using  
584 the function `infer_breakpoints`, prioritizing the boundaries defined by soft-  
585 clipping operations in the `CIGAR` string. If no clipping was present, the midpoint  
586 of the alignment or the read length was utilized as a fallback. The read sequence  
587 was then divided into left and right segments at this inferred breakpoint, and  
588 k-mer frequency profiles ( $k = 5$ ) were generated for both halves, ignoring any  
589 k-mers containing ambiguous 'N' bases. The resulting k-mer frequency vectors  
590 will be normalized and compared using the functions `cosine_difference` and  
591 `js_divergence`.

## 592 Microhomology

593 The process of extracting the microhomology feature started by utilizing the func-  
594 tion `infer_breakpoints` similar to the k-mer workflow. Once a breakpoint was es-  
595 tablished, the script scanned a  $\pm 40$  base pair window surrounding the breakpoint  
596 and used the function `longest_suffix_prefix_overlap` to identify the longest  
597 exact suffix-prefix overlap between the left and right read segments. This overlap,  
598 which represents consecutive bases shared at the junction, was recorded as the  
599 `microhomology_length` in the dataset. The 40-base pair window was chosen to  
600 ensure that short shared sequences at or near the breakpoint were captured, with-  
601 out including distant sequences that are unrelated. Additionally, the GC content  
602 of the overlapping sequence was calculated using the function `gc_content`, which  
603 counts guanine (G) and cytosine (C) bases within the detected microhomology  
604 and divides by the total length, yielding a proportion between 0 and 1, and was  
605 stored under the `microhomology_gc` attribute. Short microhomologies, typically  
606 3-20 base pairs in length, are recognized signatures of PCR-induced template  
607 switching (Peccoud et al., 2018).

608 A k-mer length of 6 was used to capture patterns within the same 40-base pair  
609 window surrounding each breakpoint. These profiles complement microhomology  
610 measurements and help identify junctions that are potentially chimeric.

### 611 3.1.3 Machine Learning Model Development

612 After feature extraction, the per-read feature matrices for clean and chimeric  
613 reads were merged into a single dataset. Each row corresponded to one paired-

614 end read, and columns encoded alignment-structure features (e.g., supplementary  
615 alignment count and spacing between segments), CIGAR-derived soft-clipping  
616 statistics (e.g., left and right soft-clipped length, total clipped bases), k-mer com-  
617 position discontinuity between read segments, and microhomology descriptors  
618 near candidate junctions. The resulting feature set was restricted to quantities  
619 that can be computed from standard BAM/FASTQ files in typical mitochondrial  
620 sequencing workflows.

621 The labelled dataset was randomly partitioned into training (80%) and test  
622 (20%) subsets using stratified sampling to preserve the 1:1 ratio of clean to  
623 chimeric reads. Model development and evaluation were implemented in Python  
624 (Version 3.11) using the `scikit-learn`, `xgboost`, `lightgbm`, and `catboost` li-  
625 braries. A broad panel of classification algorithms was then benchmarked on the  
626 training data to obtain a fair comparison of different model families under identical  
627 feature conditions. The panel included: a trivial dummy classifier,  $L_2$ -regularized  
628 logistic regression, a calibrated linear support vector machine (SVM),  $k$ -nearest  
629 neighbours, Gaussian Naïve Bayes, decision-tree ensembles (Random Forest, Ex-  
630 tremely Randomized Trees, and Bagging with decision trees), gradient boosting  
631 methods (Gradient Boosting, XGBoost, LightGBM, and CatBoost), and a shallow  
632 multilayer perceptron (MLP).

633 For each model, five-fold stratified cross-validation was performed on the train-  
634 ing set. In every fold, four-fifths of the data were used for fitting and the remaining  
635 one-fifth for validation. Mean cross-validation accuracy, precision, recall, F1-score  
636 for the chimeric class, and area under the receiver operating characteristic curve  
637 (ROC-AUC) were computed to summarize performance and rank candidate meth-  
638 ods. This baseline screen allowed comparison of linear, probabilistic, neural, and

ensemble-based approaches and identified tree-based ensemble and boosting models as consistently strong performers relative to simpler baselines.

### 3.1.4 Model Benchmarking, Hyperparameter Optimization, and Evaluation

Model selection and refinement proceeded in two stages. First, the cross-validation results from the broad panel were used to identify a subset of competitive models for more detailed optimization. Specifically, ten model families were carried forward:  $L_2$ -regularized logistic regression, calibrated linear SVM, Random Forest, ExtraTrees, Gradient Boosting, XGBoost, LightGBM, CatBoost, Bagging with decision trees, and a shallow MLP. This subset spans both linear and non-linear decision boundaries, but emphasizes ensemble and boosting methods, which showed superior F1 and ROC-AUC in the initial benchmark.

Second, hyperparameter optimization was conducted for each of the ten selected models using randomized search with five-fold stratified cross-validation (`RandomizedSearchCV`). For tree-based ensembles, the search space included the number of trees, maximum depth, minimum samples per split and leaf, and the fraction of features considered at each split. For boosting methods, key hyperparameters such as the number of boosting iterations, learning rate, tree depth, subsampling rate, and column subsampling rate were tuned. For the MLP, the number and size of hidden layers, learning rate, and  $L_2$ -regularization strength were varied. In all cases, the primary optimisation criterion was the F1-score of the chimeric class, averaged across folds.

For each model family, the hyperparameter configuration with the highest mean cross-validation F1-score was selected as the best-tuned estimator. These tuned models were then refitted on the full training set and evaluated once on the held-out test set to obtain unbiased estimates of performance. Test-set metrics included accuracy, precision, recall, F1-score for the chimeric class, and ROC-AUC. Confusion matrices and ROC curves were generated for the top-performing models to characterise common error modes, such as false negatives (missed chimeric reads) and false positives (clean reads incorrectly labelled as chimeric). The final model or small set of models for downstream interpretation was chosen based on a combination of test-set F1-score, ROC-AUC, and practical considerations such as model complexity and ease of deployment within a feature extraction pipeline.

### 3.1.5 Feature Importance and Interpretation

To relate model decisions to biologically meaningful signals, feature-importance analyses were performed on the best-performing tree-based models. Two complementary approaches were used. First, built-in importance measures from ensemble methods (e.g., split-based importances in Random Forest and Gradient Boosting) were examined to obtain an initial ranking of features based on their contribution to reducing impurity. Second, model-agnostic permutation importance was computed on the test set by repeatedly permuting each feature column while keeping all others fixed and measuring the resulting decrease in F1-score. Features whose permutation led to a larger performance drop were interpreted as more influential for chimera detection.

For interpretability, individual features were grouped into four conceptual

684 families: (i) supplementary alignment and alignment-structure features (e.g., SA  
685 count, spacing between alignment segments, strand consistency), (ii) CIGAR-  
686 derived soft-clipping features (e.g., left and right soft-clipped length, total clipped  
687 bases), (iii) k-mer composition discontinuity features (e.g., cosine distance and  
688 Jensen–Shannon divergence between k-mer profiles of read segments), and (iv) mi-  
689 crohomology descriptors (e.g., microhomology length and local GC content around  
690 putative breakpoints). Aggregating permutation importance scores within each  
691 family allowed assessment of which biological signatures contributed most strongly  
692 to the classifier’s performance. This analysis provided a basis for interpreting the  
693 trained models in terms of known mechanisms of PCR-induced template switching  
694 and for identifying which alignment- and sequence-derived cues are most informa-  
695 tive for distinguishing chimeric from clean mitochondrial reads.

### 696 **3.1.6 Validation and Testing**

697 Validation will involve both internal and external evaluations. Internal valida-  
698 tion was achieved through five-fold cross-validation on the training data to verify  
699 model generalization and reduce variance due to random sampling. External vali-  
700 dation will be achieved through testing on the 20% hold-out dataset derived from  
701 the simulated reads, which will be an unbiased benchmark to evaluate how well  
702 the trained models generalized to unseen data. All feature extraction and prepro-  
703 cessing steps were performed using the same feature extraction pipeline to ensure  
704 consistency and comparability across validation stages.

705 Comparative evaluation was performed across all candidate algorithms, in-  
706 cluding a trivial dummy classifier,  $L_2$ -regularized logistic regression, a calibrated

707 linear SVM, k-nearest neighbours, Gaussian Naïve Bayes, decision-tree ensembles,  
708 gradient boosting methods, and a shallow MLP. This evaluation determined which  
709 models demonstrated the highest predictive performance and computational effi-  
710 ciency under identical data conditions. Their metrics were compared to identify  
711 which algorithms were most suitable for further refinement.

### 712 **3.1.7 Documentation**

713 Comprehensive documentation was maintained throughout the study to ensure  
714 transparency and reproducibility. All stages of the research, including data gath-  
715 ering, preprocessing, feature extraction, model training, and validation, were sys-  
716 tematically recorded in a `.README` file in the GitHub repository. For each ana-  
717 lytical step, the corresponding parameters, software versions, and command line  
718 scripts were documented to enable exact replication of results.

719 The repository structure followed standard research data management prac-  
720 tices, with clear directories for datasets and scripts. Computational environments  
721 were standardized using Conda, with an environment file (`environment.arm.yml`)  
722 specifying dependencies and package versions to maintain consistency across sys-  
723 tems.

724 For manuscript preparation and supplementary materials, Overleaf (L<sup>A</sup>T<sub>E</sub>X)  
725 was used to produce publication-quality formatting and consistent referencing.

## 3.2 Calendar of Activities

727 Table 3.1 presents the project timeline in the form of a Gantt chart, where each  
 728 bullet point corresponds to approximately one week of planned activity.

Table 3.1: Timetable of Activities

Activities (2025)	Nov	Dec	Jan	Feb	Mar	Apr	May
Data Collection and Simulation	• • • •						
Feature Extraction Pipeline	•	•					
Machine Learning Development		•	• •	• • • •	• • • •	• •	
Testing and Validation						• •	• • • •
Documentation	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •



## Chapter 4

# Results and Discussion

### 4.1 Descriptive Analysis of Features

This chapter presents the performance of the proposed feature set and machine learning models for detecting PCR-induced chimeric reads in simulated mitochondrial Illumina data. We first describe the behaviour of the main features, then compare baseline classifiers, assess the effect of hyperparameter tuning, and finally analyse feature importance in terms of individual variables and biologically motivated feature families.

The final dataset contained 31,986 reads for training and 7,997 reads for testing, with classes balanced (approximately 4,000 clean and 4,000 chimeric reads in the test split).

### 741 4.1.1 Exploratory Data Analysis

742 An exploratory data analysis (EDA) was conducted on the extracted feature ma-  
743 trix to characterize general patterns in the data and gain preliminary insight into  
744 which variables might meaningfully contribute to classification. Histograms of  
745 key features indicated that alignment-based variables showed clear class separa-  
746 tion as chimeric reads have higher frequencies of split alignments and and no-  
747 ticeably broader long-tailed distribution on soft-clipped regions (`softclip_left`  
748 and `softclip_right`). In contrast, sequence-based variables such a microhomol-  
749 ogy length and k-mer divergence displayed substantial overlap between classes,  
750 suggesting more limited discriminative value. The complete set of histograms is  
751 provided in in Appendix A.

752 As shown in Figure 4.1, the feature correlation heatmap shows that alignment-  
753 derived variables form a strongly correlated cluster, whereas sequence-derived  
754 measures show weak correlations with both the alignment-based features and with  
755 one another. This heterogeneity indicates that no single feature family captures  
756 all relevant signal sources.

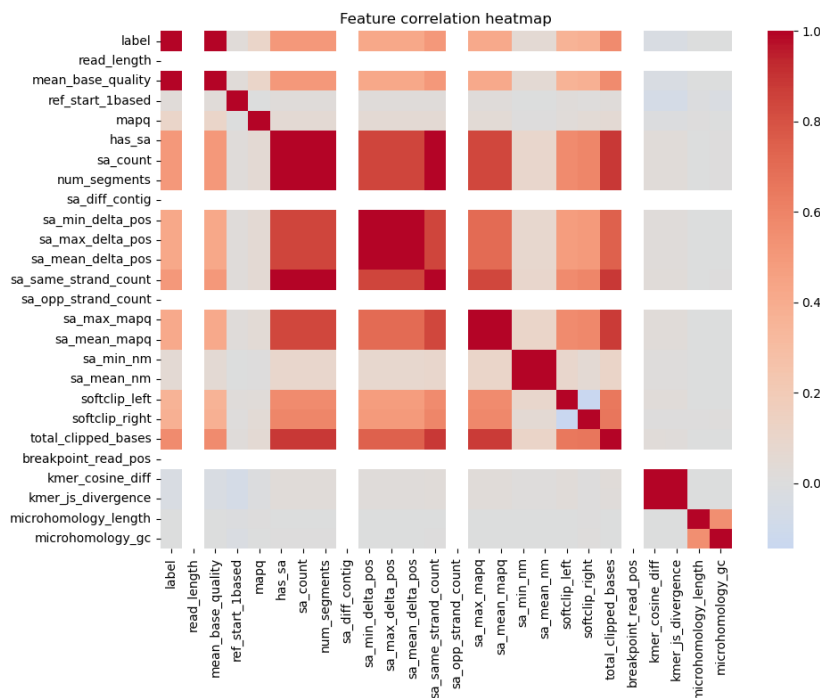


Figure 4.1: Feature correlation heatmap showing relationships among alignment-derived and sequence-derived variables.

## 4.2 Baseline Classification Performance

Table 4.1 summarises the performance of eleven classifiers trained on the engineered feature set using five-fold cross-validation and evaluated on the held-out test set. All models were optimised using default hyperparameters, without dedicated tuning.

The dummy baseline, which always predicts the same class regardless of the input features, achieved an accuracy of 0.50 and test F1-score of 0.67. This reflects the balanced class distribution and provides a lower bound for meaningful performance.

766 Across other models, test F1-scores clustered in a narrow band between ap-  
767 proximately 0.74 and 0.77 and ROC-AUC values between 0.82 and 0.84. Gradi-  
768 ent boosting, CatBoost, LightGBM, XGBoost, bagging trees, random forest, and  
769 multilayer perceptron (MLP) all produced very similar scores, with CatBoost and  
770 gradient boosting slightly ahead (test F1  $\approx 0.77$ , ROC-AUC  $\approx 0.84$ ). Linear  
771 models (logistic regression and calibrated linear SVM) performed only marginally  
772 worse (test F1  $\approx 0.74$ ), while Gaussian Naive Bayes lagged behind with substan-  
773 tially lower F1 ( $\approx 0.65$ ) despite very high precision for the chimeric class.

Table 4.1: Performance of baseline classifiers on the held-out test set.

model	test_accuracy	test_precision	test_recall	test_f1	test_roc_auc
dummy_baseline	0.500000	0.500000	1.000000	0.667000	0.500000
logreg_l2	0.789000	0.945000	0.614000	0.744000	0.821000
linear_svm_calibrated	0.789000	0.945000	0.614000	0.744000	0.820000
random_forest	0.788000	0.894000	0.654000	0.755000	0.834000
extra_trees	0.788000	0.901000	0.647000	0.753000	0.824000
gradient_boosting	0.802000	0.936000	0.648000	0.766000	0.840000
xgboost	0.800000	0.929000	0.650000	0.765000	0.839000
lightgbm	0.799000	0.926000	0.650000	0.764000	0.838000
catboost	0.803000	0.936000	0.650000	0.767000	0.839000
knn	0.782000	0.892000	0.642000	0.747000	0.815000
gaussian_nb	0.741000	0.996000	0.483000	0.651000	0.819000
bagging_trees	0.792000	0.900000	0.657000	0.760000	0.837000
mlp	0.789000	0.931000	0.625000	0.748000	0.819000

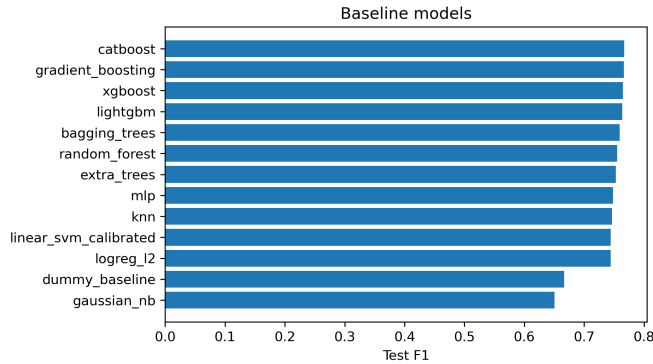


Figure 4.2: Test F1 of all baseline classifiers, showing that no single model clearly dominates and several achieve comparable performance.

### 4.3 Effect of Hyperparameter Tuning

To assess whether performance could be improved further, ten model families underwent randomised hyperparameter search (Chapter 3). The tuned metrics are summarised in Table 4.2. Overall, tuning yielded modest but consistent gains for tree-based ensembles and boosting methods, while leaving linear models essentially unchanged or slightly worse.

CatBoost, gradient boosting, LightGBM, XGBoost, random forest, bagging trees, and MLP all experienced small increases in test F1 (typically  $\Delta F1 \approx 0.002$ – $0.009$ ) and ROC–AUC (up to  $\Delta AUC \approx 0.008$ ). After tuning, CatBoost remained the best performer with test accuracy 0.802, precision 0.924, recall 0.658, F1-score 0.769, and ROC–AUC 0.844. Gradient boosting achieved almost identical performance (F1 0.767, AUC 0.843). Random forest and bagging trees also improved to F1 scores around 0.763 with  $AUC \approx 0.842$ .

Table 4.2: Performance of tuned classifiers on the held-out test set.

model	test_accuracy	test_precision	test_recall	test_f1	test_roc_auc
logreg_l2_tuned	0.788000	0.946000	0.612000	0.743000	0.818000
linear_svm_calibrated_tuned	0.788000	0.944000	0.612000	0.743000	0.818000
random_forest_tuned	0.797000	0.915000	0.655000	0.763000	0.842000
extra_trees_tuned	0.794000	0.910000	0.652000	0.760000	0.837000
gradient_boosting_tuned	0.802000	0.928000	0.654000	0.767000	0.843000
xgboost_tuned	0.799000	0.922000	0.653000	0.765000	0.839000
lightgbm_tuned	0.801000	0.930000	0.651000	0.766000	0.842000
catboost_tuned	0.802000	0.924000	0.658000	0.769000	0.844000
bagging_trees_tuned	0.798000	0.922000	0.650000	0.763000	0.842000
mlp_tuned	0.790000	0.934000	0.625000	0.749000	0.821000

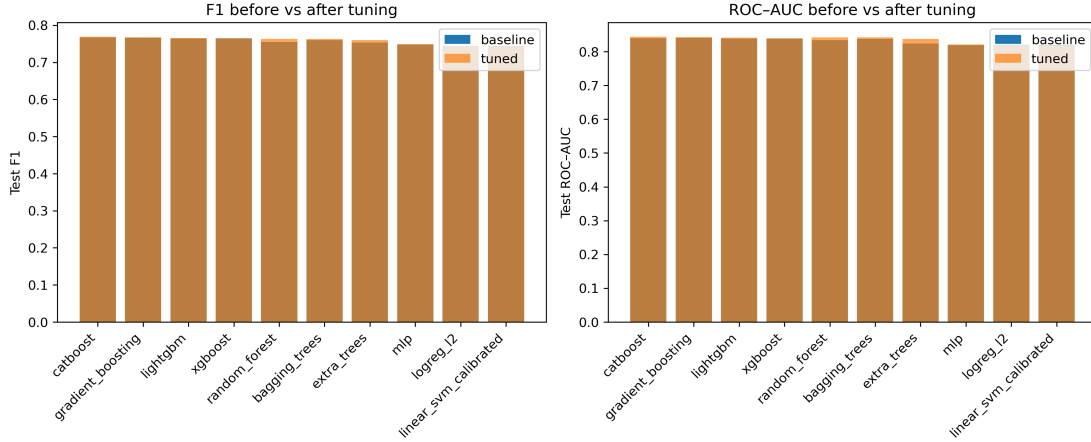


Figure 4.3: Comparison of test F1 (left) and ROC-AUC (right) for baseline and tuned models. Hyperparameter tuning yields small but consistent gains, particularly for tree-based ensembles.

Because improvements are small and within cross-validation variability, we interpret tuning as stabilising and slightly refining the models rather than fundamentally altering their behaviour or their relative ranking.

## 4.4 Detailed Evaluation of Representative Models

For interpretability and diversity, four tuned models were selected for deeper analysis: CatBoost (best-performing boosted tree), scikit-learn gradient boosting (canonical gradient-boosting implementation), random forest (non-boosted ensemble baseline), and  $L_2$ -regularised logistic regression (linear baseline). All models were trained on the engineered feature set and evaluated on the same held-out test data.

#### 798 4.4.1 Confusion Matrices and Error Patterns

799 Classification reports and confusion matrices for the four models reveal consistent  
800 patterns. CatBoost and gradient boosting both reached overall accuracy of ap-  
801 proximately 0.80 with similar macro-averaged F1 scores ( $\sim 0.80$ ). For CatBoost,  
802 precision and recall for clean reads were 0.73 and 0.95, respectively, while for  
803 chimeric reads they were 0.92 and 0.66 ( $F1 = 0.77$ ). Gradient boosting showed  
804 nearly identical trade-offs.

805 Random forest attained slightly lower accuracy (0.80) and chimeric F1 (0.76),  
806 whereas logistic regression achieved the lowest accuracy among the four (0.79)  
807 and chimeric F1 (0.74), although it provided the highest chimeric precision (0.95)  
808 at the cost of lower recall (0.61).

809 Across all models, errors were asymmetric. False negatives (chimeric reads  
810 predicted as clean) were more frequent than false positives. For example, CatBoost  
811 misclassified 1 369 chimeric reads as clean but only 215 clean reads as chimeric.  
812 This pattern indicates that the models are conservative: they prioritise avoiding  
813 spurious chimera calls at the expense of missing some true chimeras. Depending on  
814 downstream application, alternative decision thresholds or cost-sensitive training  
815 could be explored to adjust this balance.

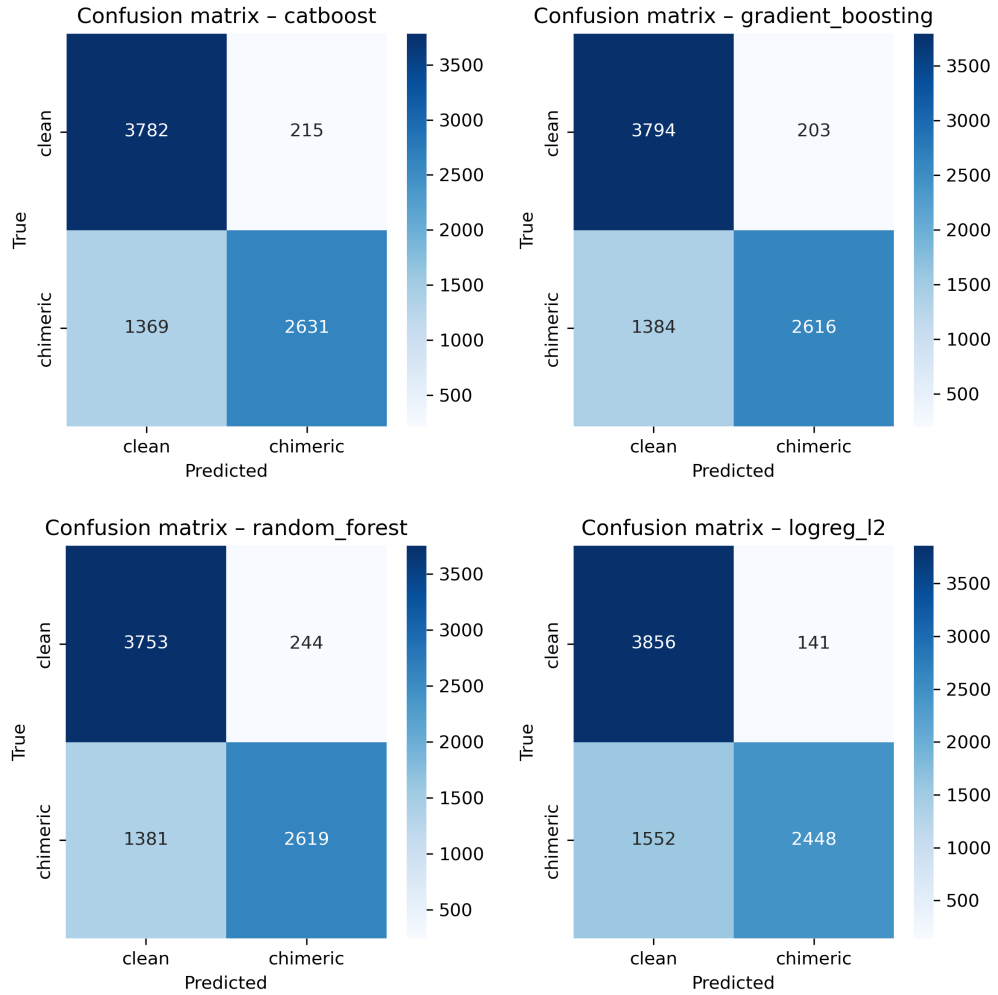


Figure 4.4: Confusion matrices for the four representative models on the held-out test set. All models show more false negatives (chimeric reads called clean) than false positives.

#### 4.4.2 ROC and Precision–Recall Curves

Receiver operating characteristic (ROC) and precision–recall (PR) curves (Figure 4.5) further support the similarity among the top models. The three tree-based ensembles (CatBoost, gradient boosting, random forest) achieved ROC–AUC values of approximately 0.84 and average precision (AP) around 0.88. Logistic re-



gression performed slightly worse ( $AUC \approx 0.82$ ,  $AP \approx 0.87$ ) but still substantially better than random guessing.

The PR curves show that precision remains above 0.9 across a broad range of recall values (up to roughly 0.5–0.6), after which precision gradually declines. This behaviour indicates that the models can assign very high confidence to a subset of chimeric reads, while more ambiguous reads can only be recovered by accepting lower precision.

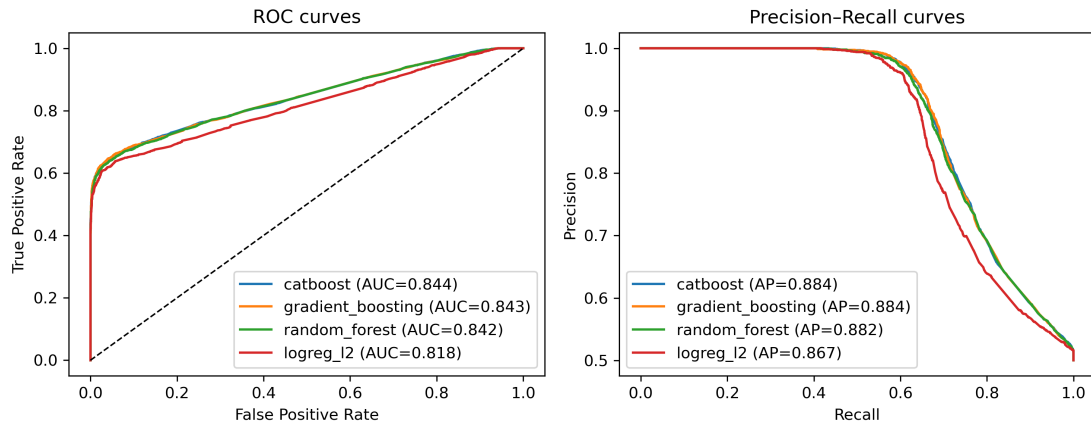


Figure 4.5: ROC (left) and precision–recall (right) curves for the four representative models on the held-out test set. Tree-based ensembles cluster closely, with logistic regression performing slightly but consistently worse.

## 828 4.5 Feature Importance and Biological Interpre- 829 tation

### 830 4.5.1 Permutation Importance of Individual Features

831 To understand how each classifier made predictions, feature importance was quan-  
832 tified using permutation importance. In this approach, the values of a single fea-  
833 ture are randomly shuffled, and the resulting drop in  $F_1$  score ( $\Delta F_1$ ) reflects how  
834 strongly the model depends on that feature. Greater decreases in  $F_1$  indicate  
835 stronger reliance on that feature. This analysis was applied to four representa-  
836 tive models: CatBoost, Gradient Boosting, Random Forest, and  $L_2$ -regularized  
837 Logistic Regression.

838 As shown in Figure 4.6, the total number of clipped bases consistently pro-  
839 vides a strong predictive signal, particularly in Random Forest, Gradient Boosting,  
840 and  $L_2$ -regularized Logistic Regression. CatBoost differs by assigning the highest  
841 importance to k-mer divergence metrics such as `kmer_js_divergence`, which cap-  
842 ture subtle sequence changes resulting from structural variants or PCR-induced  
843 chimeras. Soft-clipping features (`softclip_left` and `softclip_right`) provide  
844 additional context around breakpoints, complementing these primary signals in  
845 all models except Gradient Boosting.  $L_2$ -regularized Logistic Regression relies  
846 more on alignment-based split-read metrics when breakpoints are simple, but it is  
847 less effective at detecting complex rearrangements that introduce novel sequences.

848 Overall, these results indicate that accurate detection of chimeric reads relies  
849 on both alignment-based signals and k-mer compositional information. Explicit

850 microhomology features contribute minimally in this analysis, and combining both  
851 alignment-based and sequence-level features enhances model sensitivity and speci-  
852 ficity.

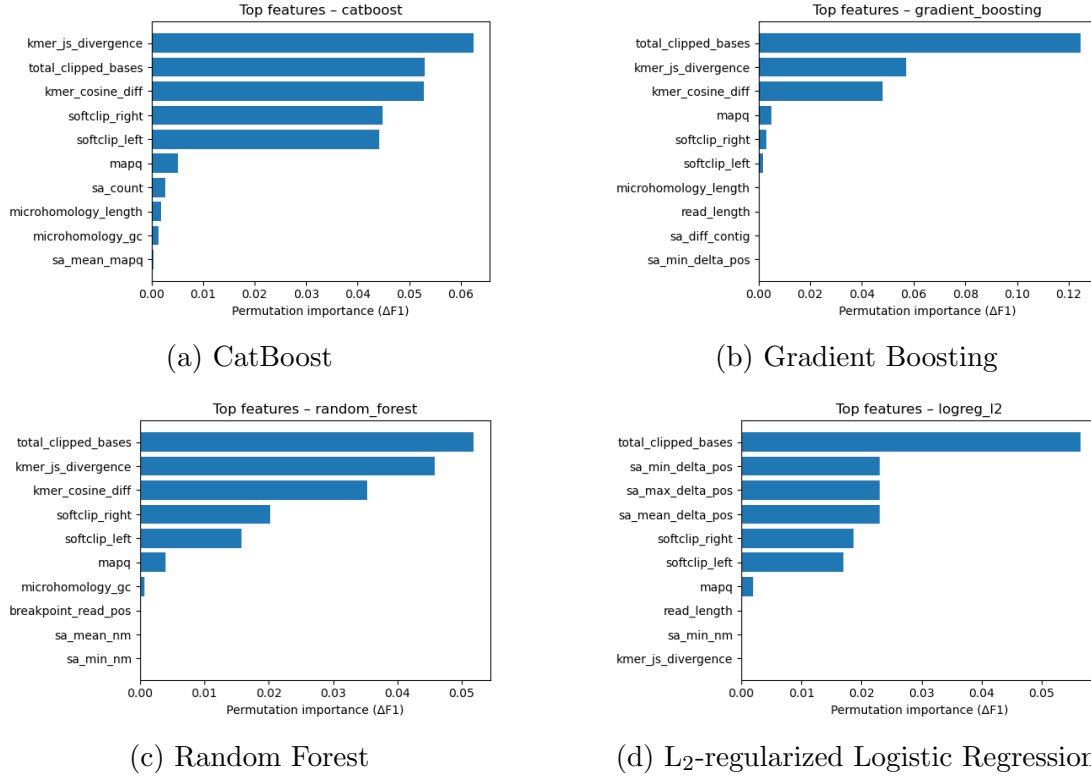


Figure 4.6: Permutation-based feature importance for four representative classifiers. Clipping and k-mer composition features are generally the strongest predictors, whereas microhomology and other alignment metrics contribute minimally.

## 853 4.5.2 Feature Family Importance

854 To evaluate the contribution of broader biological signals, features were  
855 grouped into five families: SA\_structure (supplementary alignment and seg-  
856 ment metrics, e.g., `has_sa`, `sa_count`, `sa_min_delta_pos`, `sa_mean_nm`), Clipping  
857 (`softclip_left`, `softclip_right`, `total_clipped_bases`, `breakpoint_read_pos`),

858 Kmer\_jump (`kmer_cosine_diff`, `kmer_js_divergence`), `Micro_homology`, and  
859 Other (e.g., `mapq`).

860     Aggregated analyses reveal consistent patterns across models. In CatBoost,  
861 the Clipping family has the largest cumulative contribution (0.14), followed  
862 by Kmer\_jump (0.12), with Other features contributing modestly (0.005) and  
863 SA\_structure (0.003) and Micro\_homology (0.003) providing minimal predictive  
864 power. Gradient Boosting shows a similar trend, with Clipping (0.13) domi-  
865 nating, Kmer\_jump (0.11) secondary, and the remaining families contributing  
866 negligibly. Random Forest integrates both Clipping (0.088) and Kmer\_jump  
867 (0.08) effectively, while SA\_structure, Micro\_homology, and Other remain minor  
868 contributors. L<sub>2</sub>-regularized Logistic Regression emphasizes Clipping (0.09)  
869 and SA\_structure (0.07), with Kmer\_jump and Micro\_homology having minimal  
870 impact.

871     Both feature-level and aggregated analyses indicate that detection of chimeric  
872 reads in this dataset relies primarily on alignment disruptions (Clipping) and  
873 k-mer compositional shifts (Kmer\_jump), which often arise from PCR-induced  
874 recombination events, while explicit microhomology features contribute minimally.

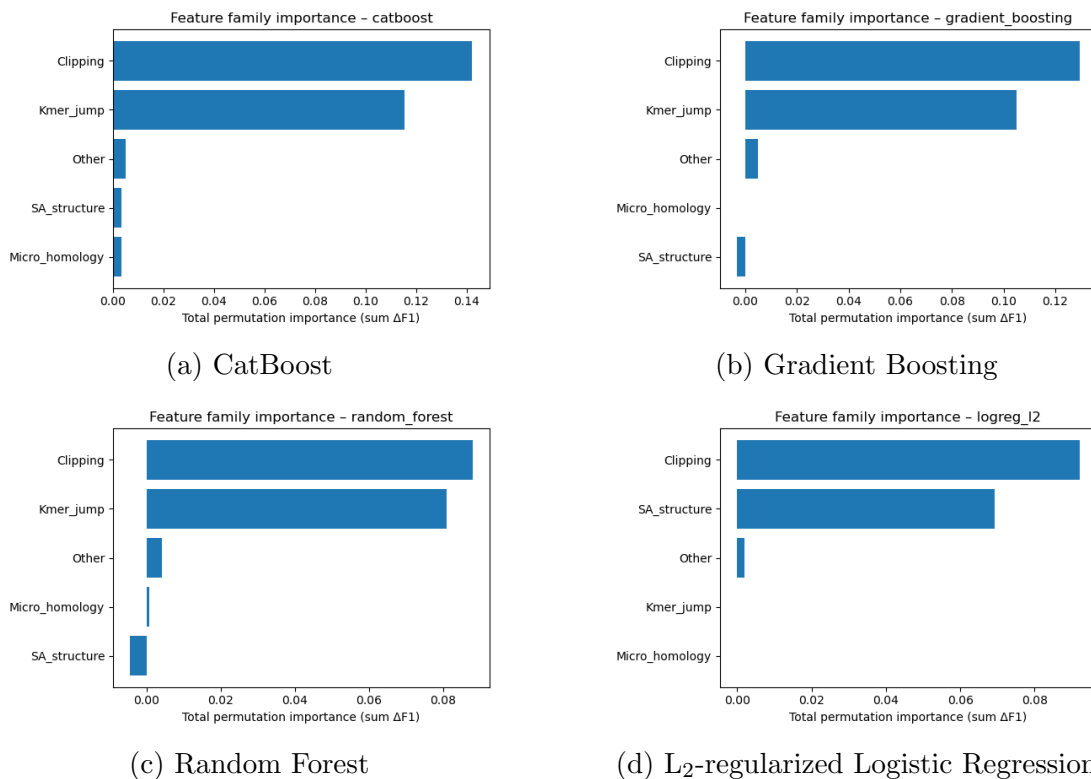


Figure 4.7: Aggregated feature family importance across four models. Clipping and k-mer compositional shifts are consistently the dominant contributors, while SA\_structure, Micro\_homology, and other features contribute minimally.

## 4.6 Summary of Findings

After removing trivially discriminative metadata, all models performed substantially better than the dummy baseline, with test F1-scores around 0.76 and ROC-AUC values near 0.84. Hyperparameter tuning yielded modest improvements, with boosting methods, particularly CatBoost and gradient boosting, achieving the highest performance. Confusion matrices and precision-recall curves indicate that these models prioritise precision for chimeric reads while accepting lower recall, which is a conservative strategy appropriate for scenarios where false positives

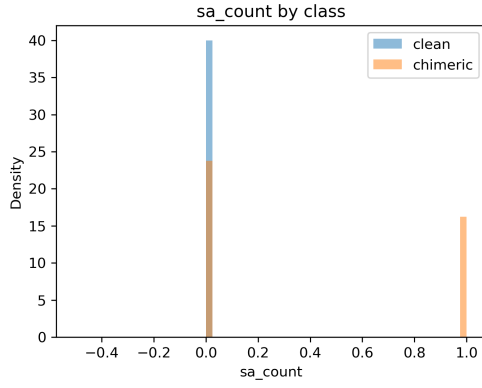
883 are costly.

884 Feature importance analyses revealed that alignment disruptions, such as clip-  
885 ping, and abrupt k-mer composition changes accounted for most predictive power.  
886 In contrast, microhomology metrics and supplementary alignment descriptors con-  
887 tributed minimally. These results indicate that features based on read alignment  
888 and k-mer composition are sufficient to train classifiers for detecting mitochon-  
889 drial PCR-induced chimera reads, without needing additional quality-score or  
890 positional information in the conditions tested.

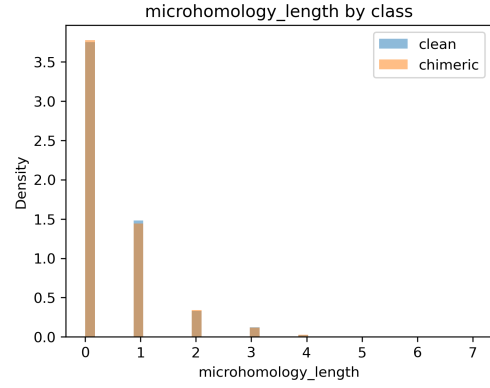
# 891 **Appendix A**

## 892 **Exploratory Data Analysis**

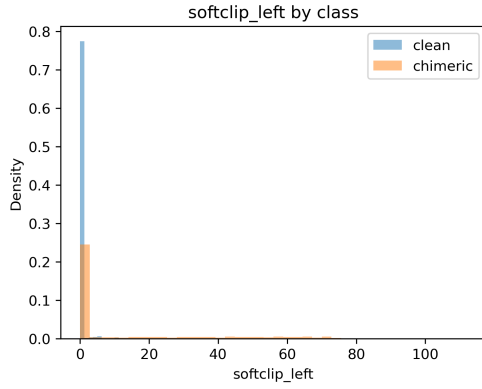
### 893 **A.1 Histograms of Key Features**



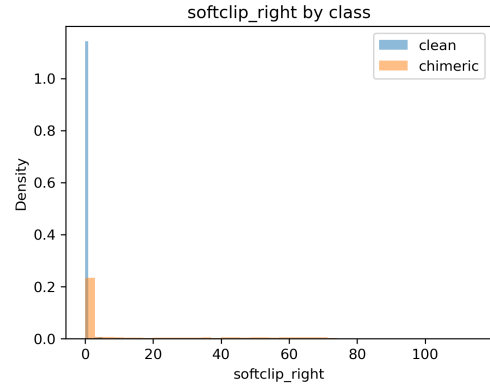
(a) sa\_count



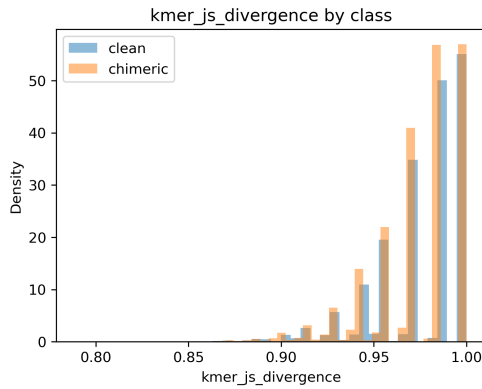
(b) Microhomology length



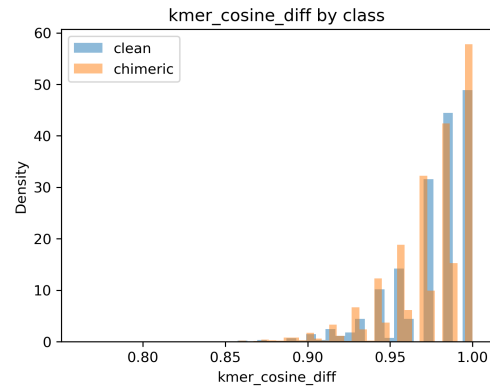
(c) softclip\_left



(d) softclip\_right



(e) k-mer Jensen-Shannon divergence



(f) k-mer cosine difference

Figure A.1: Histogram plots of six key features comparing clean and chimeric reads.



## 894 References

- 895 Anderson, S., Bankier, A., Barrell, B., Bruijn, M., Coulson, A., Drouin, J., ...  
896 Young, I. (1981, 04). Sequence and organization of the human mitochondrial  
897 genome. *Nature*, 290, 457-465. doi: 10.1038/290457a0
- 898 Arango, G., Garner, E., Pruden, A., Heath, L., Vikesland, P., & Zhang, L. (2018,  
899 02). Deeparg: A deep learning approach for predicting antibiotic resistance  
900 genes from metagenomic data. *Microbiome*, 6. doi: 10.1186/s40168-018  
901 -0401-z
- 902 Bentley, D. R., Balasubramanian, S., Swerdlow, H. P., Smith, G. P., Milton, J.,  
903 Brown, C. G., ... Smith, A. J. (2008). Accurate whole human genome  
904 sequencing using reversible terminator chemistry. *Nature*, 456(7218), 53–  
905 59. doi: 10.1038/nature07517
- 906 Boore, J. L. (1999). Animal mitochondrial genomes. *Nucleic Acids Research*,  
907 27(8), 1767–1780. doi: 10.1093/nar/27.8.1767
- 908 Cameron, S. L. (2014). Insect mitochondrial genomics: Implications for evolution  
909 and phylogeny. *Annual Review of Entomology*, 59, 95–117. doi: 10.1146/  
910 annurev-ento-011613-162007
- 911 Dierckxsens, N., Mardulyn, P., & Smits, G. (2017). Novoplasty: de novo assembly  
912 of organelle genomes from whole genome data. *Nucleic Acids Research*,

913 45(4), e18. doi: 10.1093/nar/gkw955

914 Edgar, R. C. (2016). Uchime2: improved chimera prediction for amplicon se-

915 quencing. *bioRxiv*. Retrieved from [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:88955007)

916 CorpusID:88955007

917 Edgar, R. C. (n.d). *Uchime in practice*. Retrieved from [https://www.drive5](https://www.drive5.com/usearch/manual7/uchime_practical.html)

918 [.com/usearch/manual7/uchime\\_practical.html](https://www.drive5.com/usearch/manual7/uchime_practical.html)

919 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011).

920 Uchime improves sensitivity and speed of chimera detection. *Bioinformatics*,

921 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

922 Glenn, T. C. (2011). Field guide to next-generation dna sequencers. *Molecular*

923 *Ecology Resources*, 11(5), 759–769. doi: 10.1111/j.1755-0998.2011.03024.x

924 Gonzalez, J. M., Zimmermann, J., & Saiz-Jimenez, C. (2004, 09). Evalu-

925 ating putative chimeric sequences from pcr-amplified products. *Bioin-*

926 *formatics*, 21(3), 333-337. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bti008)

927 [bioinformatics/bti008](https://doi.org/10.1093/bioinformatics/bti008) doi: 10.1093/bioinformatics/bti008

928 Gray, M. W. (2012). Mitochondrial evolution. *Cold Spring Harbor perspectives*

929 *in biology*, 4. Retrieved from [https://doi.org/10.1101/cshperspect](https://doi.org/10.1101/cshperspect.a011403)

930 [.a011403](https://doi.org/10.1101/cshperspect.a011403) doi: 10.1101/cshperspect.a011403

931 Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial

932 genomes directly from genomic next-generation sequencing reads—a baiting

933 and iterative mapping approach. *Nucleic Acids Research*, 41(13), e129. doi:

934 10.1093/nar/gkt371

935 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,

936 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de

937 novo assembly of organelle genomes. *Genome Biology*, 21(1), 241. doi:

938 10.1186/s13059-020-02154-5

- 939 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-  
940 pression of pcr-mediated recombination. *Nucleic Acids Research*, 26(7),  
941 1819–1825. doi: 10.1093/nar/26.7.1819
- 942 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.  
943 (2021). Mitochondrial dna reveals genetically structured haplogroups of  
944 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*  
945 *Marine Science*, 41, 101588. doi: 10.1016/j.rsma.2020.101588
- 946 Li, H. (2018, 05). Minimap2: pairwise alignment for nucleotide sequences. *Bioin-*  
947 *formatics*, 34(18), 3094–3100. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bty191)  
948 [bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191) doi: 10.1093/bioinformatics/bty191
- 949 Liang, Q., Bible, P. W., Liu, Y., Zou, B., & Wei, L. (2020, 02). Deepmi-  
950 crobes: taxonomic classification for metagenomics with deep learning. *NAR*  
951 *Genomics and Bioinformatics*, 2(1), lqaa009. Retrieved from [https://](https://doi.org/10.1093/nargab/lqaa009)  
952 [doi.org/10.1093/nargab/lqaa009](https://doi.org/10.1093/nargab/lqaa009) doi: 10.1093/nargab/lqaa009
- 953 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*  
954 *Reviews Genetics*, 11(1), 31–46. doi: 10.1038/nrg2626
- 955 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,  
956 an ensemble classifier for chimera detection in 16s rna sequencing stud-  
957 ies. *Applied and Environmental Microbiology*, 81(5), 1573–1584. Retrieved  
958 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:  
959 10.1128/AEM.02896-14
- 960 Peccoud, J., Lequime, S., Moltini-Conclois, I., Giraud, I., Lambrechts, L., &  
961 Gilbert, C. (2018, 04). A survey of virus recombination uncovers canon-  
962 ical features of artificial chimeras generated during deep sequencing li-  
963 brary preparation. *G3 Genes—Genomes—Genetics*, 8(4), 1129–1138. Re-  
964 trieved from <https://doi.org/10.1534/g3.117.300468> doi: 10.1534/

g3.117.300468

Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J. (2023). Effects of error, chimera, bias, and gc content on the accuracy of amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from <https://journals.asm.org/doi/abs/10.1128/msystems.01025-23> doi: 10.1128/msystems.01025-23

Qiu, X., Wu, L., Huang, H., McDonel, P. E., Palumbo, A. V., Tiedje, J. M., & Zhou, J. (2001). Evaluation of pcr-generated chimeras, mutations, and heteroduplexes with 16s rna gene-based cloning. *Applied and Environmental Microbiology*, 67(2), 880–887. doi: 10.1128/AEM.67.2.880-887.2001

Ren, J., Song, K., Deng, C., Ahlgren, N., Fuhrman, J., Li, Y., ... Sun, F. (2020, 01). Identifying viruses from metagenomic data using deep learning. *Quantitative Biology*, 8. doi: 10.1007/s40484-019-0187-4

Rodriguez-Martin, B., Palumbo, E., Marco-Sola, S., Griebel, T., Ribeca, P., Alonso, G., ... Djebali, S. (2017, 01). Chimpipes: Accurate detection of fusion genes and transcription-induced chimeras from rna-seq data. *BMC Genomics*, 18. doi: 10.1186/s12864-016-3404-9

Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). Vsearch: a versatile open source tool for metagenomics. *PeerJ*, 4, e2584. doi: 10.7717/peerj.2584

Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler, A., & Schatz, M. (2018, 06). Accurate detection of complex structural variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10.1038/s41592-018-0001-7

Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A back-up survival mechanism or dedicated pathway? *Trends in Biochemical*

991 *Sciences*, 40(11), 701-714. Retrieved from <https://www.sciencedirect>  
 992 [.com/science/article/pii/S0968000415001589](https://www.sciencedirect.com/science/article/pii/S0968000415001589) doi: [https://doi.org/](https://doi.org/10.1016/j.tibs.2015.08.006)  
 993 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)  
 994 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,  
 995 11). Large-scale machine learning for metagenomics sequence classifica-  
 996 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)  
 997 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683  
 998 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*  
 999 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI  
 1000 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)  
 1001 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)