

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                             February 7, 2026

## 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 lemur* Illumina paired-end data without relying on external reference  
 31 databases.

32 Using simulated datasets containing clean and chimeric reads, a feature  
 33 set was extracted, combining alignment-based metrics (e.g., supplementary  
 34 alignments, soft-clipping) with sequence-derived statistics (e.g., k-mer com-  
 35 position, microhomology). A comparative evaluation of supervised learning  
 36 models identified tree-based ensembles CatBoost and Gradient Boosting as top  
 37 performers, achieving an F1-score of 0.77 and an ROC-AUC of 0.84 on held-  
 38 out 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

43

# 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 . . . . .	3
48	1.3.1 General Objective . . . . .	3
49	1.3.2 Specific Objectives . . . . .	4
50	1.4 Scope and Limitations of the Research . . . . .	4
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 . . . . .	23
69	3.1.3	Machine Learning Model Development . . . . .	26
70	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
71		Evaluation . . . . .	27
72	3.1.5	Feature Importance, Feature Selection, and Interpretation	29

73	3.1.6	Validation and Testing . . . . .	31
74	3.1.7	Documentation . . . . .	32
75	3.2	Calendar of Activities . . . . .	32
76	<b>4</b>	<b>Results and Discussion</b>	<b>34</b>
77	4.1	Descriptive Analysis of Features . . . . .	35
78	4.1.1	Summary Statistics Per Class . . . . .	35
79	4.1.2	Correlation Analysis of Extracted Features . . . . .	40
80	4.2	Baseline Classification Performance . . . . .	41
81	4.3	Effect of Hyperparameter Tuning . . . . .	43
82	4.4	Detailed Evaluation of Representative Models . . . . .	45
83	4.4.1	Confusion Matrices and Error Patterns . . . . .	45
84	4.4.2	ROC and Precision–Recall Curves . . . . .	47
85	4.5	Feature Importance . . . . .	48
86	4.5.1	Permutation Importance of Individual Features . . . . .	48
87	4.5.2	Feature Family Importance . . . . .	50
88	4.6	Feature Selection . . . . .	52
89	4.6.1	Cumulative Importance Curve . . . . .	53

90	4.6.2	Performance Comparison Across Feature Sets . . . . .	53
91	4.6.3	Interpretation and Final Feature Set Choice . . . . .	55
92	4.7	Summary of Findings . . . . .	55
93	<b>A</b>	<b>Complete Per-Class Summary Statistics</b>	<b>57</b>
94	<b>B</b>	<b>Boxplots for All Numeric Features by Feature Family</b>	<b>61</b>
95	B.0.1	SA Structure (Supplementary Alignment and Segment Met-	
96		rics) . . . . .	61
97	B.0.2	Clipping-Based Features . . . . .	63
98	B.0.3	K-mer Features . . . . .	63
99	B.0.4	Microhomology Features . . . . .	64
100	B.0.5	Others . . . . .	64

# 101 List of Figures

102	3.1	Process diagram of the study workflow. . . . .	20
103	4.1	Boxplots of selected features for clean and chimeric reads. . . . .	39
104	4.2	Feature correlation heatmap showing relationships among alignment-	
105		derived and sequence-derived features. . . . .	41
106	4.3	Test F1 of all baseline classifiers. . . . .	43
107	4.4	Comparison of test F1 (left) and ROC-AUC (right) for baseline	
108		and tuned models. . . . .	44
109	4.5	Confusion matrices for the four representative models on the held-	
110		out test set. . . . .	47
111	4.6	ROC (left) and precision-recall (right) curves for the four represen-	
112		tative models on the held-out test set. . . . .	48
113	4.7	Permutation-based feature importance for four representative clas-	
114		sifiers. . . . .	50

115	4.8	Aggregated feature family importance across four models. . . . .	52
116	4.9	Cumulative importance curve of features sorted by importance. . .	53
117	4.10	Comparison of F1 and ROC–AUC for the full, top-4 selected, and	
118		no-microhomology feature set variants. . . . .	54
119	B.1	Boxplots of SA Structure features by class (1/2). . . . .	61
120	B.2	Boxplots of SA Structure features by class (2/2). . . . .	62
121	B.3	Boxplots of clipping-based features by class. . . . .	63
122	B.4	Boxplots of k-mer features by class. . . . .	63
123	B.5	Boxplots of microhomology features by class. . . . .	64
124	B.6	Boxplots of other numeric features by class. . . . .	64



# 125 List of Tables

126	2.1	Comparison of Chimera Detection Approaches and Tools . . . . .	17
127	3.1	Timetable of activities. . . . .	33
128	4.1	Summary statistics of selected key features by class. . . . .	38
129	4.2	Performance of baseline classifiers on the held-out test set. . . . .	42
130	4.3	Performance of tuned classifiers on the held-out test set. . . . .	44
131	4.4	Test set performance of three feature set variants using tuned Gra-	
132		dient Boosting. . . . .	54
133	A.1	Complete per-class summary statistics for all extracted features. .	58

# 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 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 PCR-induced chimeric reads using both alignment-based and sequence-derived statistical features. The tool aims to provide bioinformatics laboratories, particularly the Philippine Genome Center Visayas (PGC Visayas), with an efficient

173 solution for improving mitochondrial genome reconstruction.

## 174 1.2 Problem Statement

175 Chimeric reads can distort assembly graphs and cause misassemblies, with par-  
176 ticularly severe effects in mitochondrial genomes (Boore, 1999; Cameron, 2014).  
177 Existing assembly pipelines such as GetOrganelle, MITObim, and NOVOPlasty  
178 assume that sequencing reads are free of such artifacts (Dierckxsens et al., 2017;  
179 Hahn et al., 2013; Jin et al., 2020). At PGC Visayas, several mitochondrial as-  
180 semblies have failed or yielded incomplete contigs despite sufficient coverage, sug-  
181 gesting that undetected chimeric reads compromise assembly reliability. Mean-  
182 while, existing chimera detection tools such as UCHIME and VSEARCH were  
183 developed primarily for amplicon-based community analysis and rely heavily on  
184 reference or taxonomic comparisons (Edgar, Haas, Clemente, Quince, & Knight,  
185 2011; Rognes, Flouri, Nichols, Quince, & Mahé, 2016). These approaches are un-  
186 suitable for single-species organellar data, where complete reference genomes are  
187 often unavailable.

## 188 1.3 Research Objectives

### 189 1.3.1 General Objective

190 This study aims to develop and evaluate a machine learning-based pipeline (Mi-  
191 toChime) that detects PCR-induced chimeric reads in *Sardinella lemuru* mito-

192 chondrial sequencing data in order to improve the quality and reliability of down-  
193 stream mitochondrial genome assemblies.

### 194 **1.3.2 Specific Objectives**

195 Specifically, the study aims to:

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

## 207 **1.4 Scope and Limitations of the Research**

208 This study focuses solely on PCR-induced chimeric reads in *Sardinella lemuru*  
209 mitochondrial sequencing data, with the species choice guided by four consid-  
210 erations: (1) to limit interspecific variation in mitochondrial genome size, GC

211 content, and repetitive regions so that differences in read patterns can be at-  
212 tributed more directly to PCR-induced chimerism, (2) to align the analysis with  
213 relevant *S. lemuru* sequencing projects at PGC Visayas, (3) to take advantage of  
214 the availability of *S. lemuru* mitochondrial assemblies and raw datasets in public  
215 repositories such as the National Center for Biotechnology Information (NCBI),  
216 which facilitates reference selection and benchmarking, and (4) to develop a tool  
217 that directly supports local studies on *S. lemuru* population structure and fisheries  
218 management.

219 The study emphasizes **wgsim**-based simulations and selected empirical mito-  
220 chondrial datasets from *S. lemuru*. It excludes naturally occurring chimeras, nu-  
221 clear mitochondrial pseudogenes (NUMTs), and large-scale assembly rearrange-  
222 ments in nuclear genomes. Feature extraction is restricted to low-dimensional  
223 alignment and sequence statistics, such as k-mer frequency profiles, GC con-  
224 tent, soft and hard clipping metrics, and split-alignment counts rather than high-  
225 dimensional deep learning embeddings. This design keeps model behaviour inter-  
226 pretable and ensures that the pipeline can be run on standard workstations at  
227 PGC Visayas. Testing on long-read platforms (e.g., Nanopore, PacBio) and other  
228 taxa is outside the scope of this project.

229 Other limitations in this study include the following: simulations with vary-  
230 ing error rates were not performed, so the effect of different sequencing errors on  
231 model performance remains unexplored; alternative parameter settings, including  
232 k-mer lengths and microhomology window sizes, were not systematically tested,  
233 which could affect the sensitivity of both k-mer and microhomology feature de-  
234 tection; and the machine learning models rely on supervised training with labeled  
235 examples, which may limit their ability to detect novel or unexpected chimeric

236 patterns.

## 237 1.5 Significance of the Research

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

## 248 Chapter 2

## 249 Review of Related Literature

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

### 256 2.1 The Mitochondrial Genome

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



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

### 267 **2.1.1 Mitochondrial Genome Assembly**

268 Mitochondrial genome assembly refers to the reconstruction of the complete mito-  
269 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is  
270 conducted to obtain high-quality, continuous representations of the mitochondrial  
271 genome that can be used for a wide range of analyses, including species identi-  
272 fication, phylogenetic reconstruction, evolutionary studies, and investigations of  
273 mitochondrial diseases. Because mtDNA evolves rapidly, its assembled sequence  
274 provides valuable insights into population structure, lineage divergence, and adap-  
275 tive evolution across taxa (Boore, 1999). Compared to nuclear genome assembly,  
276 assembling the mitochondrial genome is often considered more straightforward but  
277 still encounters technical challenges such as the formation of chimeric reads. Com-  
278 monly used tools for mitogenome assembly such as GetOrganelle and MITObim  
279 operate under the assumption of organelle genome circularity, and are vulnerable  
280 when chimeric reads disrupt this circular structure, resulting in assembly errors  
281 (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

326 sequence correspond to distinct high-abundance sequences.

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

### 332 **2.3.1 UCHIME**

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

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

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

### 354 **2.3.2 UCHIME2**

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

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

### 371 2.3.3 CATCh

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

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

388 Benchmarking in both reference-based and de novo modes demonstrated signif-  
389 icant performance improvements. CATCh achieved sensitivities of approximately  
390 85 percent in reference-based mode and 92 percent in de novo mode, with corre-  
391 sponding specificities of approximately 96 percent and 95 percent. These results  
392 indicate that CATCh detected 7 to 12 percent more chimeras than any individual  
393 algorithm while maintaining high precision.

### 394 2.3.4 ChimPipe

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

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

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

## 412 **2.4 Machine Learning Approaches for Chimera** 413 **and Sequence Quality Detection**

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

### 425 **2.4.1 Feature-Based Representations of Genomic Se-** 426 **quences**

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



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

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

443 Microhomology, or short sequences shared between adjacent segments, is an-  
444 other biologically meaningful feature. Short microhomologies, typically 3–20 bp,  
445 are involved in template switching both in cellular repair pathways and during  
446 PCR, where they act as signatures of chimera formation (Peccoud et al., 2018;  
447 Sfeir & Symington, 2015). In PCR-induced chimeras, short identical sequences  
448 at junctions provide a clear signature of chimerism. Measuring the longest exact  
449 overlap at each breakpoint complements k-mer and alignment features and helps  
450 identify reads that are potentially chimeric.

## 451 2.5 Synthesis of Chimera Detection Approaches

452 To provide an integrated overview of the literature discussed in this chapter, Ta-  
453 ble 2.1 summarizes the major chimera detection studies, their methodological  
454 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.

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

## 467 Chapter 3

# 468 Research Methodology

469 This chapter outlines the steps involved in completing the study, including data  
470 gathering, generating simulated mitochondrial Illumina reads, preprocessing and  
471 indexing the data, developing a feature extraction pipeline to obtain read-level fea-  
472 tures, applying machine learning algorithms for chimera detection, implementing  
473 feature selection methods, and validating and comparing model performance.

### 474 3.1 Research Activities

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

481 through a feature extraction pipeline that computed k-mer profiles, supplementary  
 482 alignment (SA) features, and microhomology information to prepare the data  
 483 for model construction. The machine learning models were trained using the  
 484 processed input, evaluated using cross-validation and held-out testing, tuned for  
 485 improved performance, and then subjected to feature importance and feature  
 486 selection analyses before final validation.

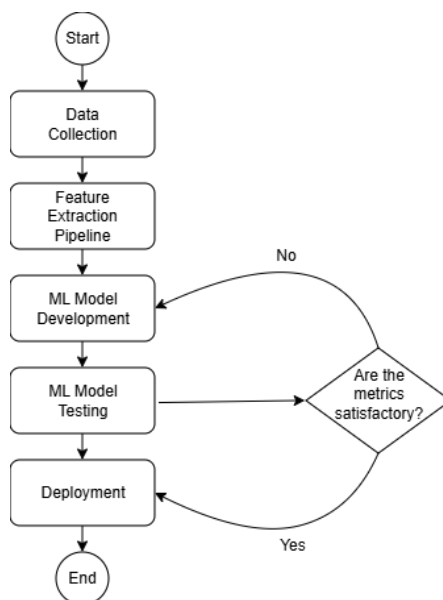


Figure 3.1: Process diagram of the study workflow.

### 487 3.1.1 Data Collection

488 The mitochondrial genome reference sequence of *S. lemur* was obtained from the  
 489 NCBI database (accession number NC\_039553.1) in FASTA format and was used  
 490 to generate simulated reads.

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

493 proximately one (1) week.

## 494 Data Preprocessing

495 All steps in the simulation and preprocessing pipeline were executed using a cus-  
496 tom script in Python (Version 3.11). The script runs each stage, including read  
497 simulation, reference indexing, mapping, and alignment processing, in a fixed se-  
498 quence.

499 `wgsim` (Version 1.13) was used to simulate 10,000 paired-end fragments, pro-  
500 ducing 20,000 reads (10,000 forward and 10,000 reverse) from the original refer-  
501 ence (`original_reference.fasta`) and designated as clean reads. The tool was  
502 selected because it provides fast generation of Illumina-like reads with controllable  
503 error rates, using the following command:

```
504 wgsim -1 150 -2 150 -r 0 -R 0 -X 0 -e 0.05 -N 10000 \  
505     original_reference.fasta ref1.fastq ref2.fastq
```

506 Chimeric sequences were then generated from the same reference FASTA  
507 file using a separate Python script. Two non-adjacent segments were ran-  
508 domly selected such that their midpoint distances fell within specified minimum  
509 and maximum thresholds. The script attempted to retain microhomology to  
510 mimic PCR-induced template switching. The resulting chimeras were written  
511 to `chimera_reference.fasta` and processed with `wgsim` to simulate 10,000  
512 paired-end fragments, generating 20,000 chimeric reads (10,000 forward reads in  
513 `chimeric1.fastq` and 10,000 reverse reads in `chimeric2.fastq`) using the same  
514 command format as above.

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

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

517       Minimap2 (Version 2.28) was used to map simulated clean and chimeric reads  
518 to the original reference. An index (`ref.mmi`) was first generated to enable efficient  
519 alignment, and mapping produced the alignment features used as input for the  
520 machine learning model. The reads were mapped using the following commands:

```
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       The resulting clean and chimeric SAM files contain the alignment positions of  
525 each read relative to the original reference genome. These files were then converted  
526 to BAM format, sorted, and indexed using `samtools` (Version 1.20):

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

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

542 This whole process was scheduled to start in the second week of November 2025  
543 and was expected to be completed by the last week of November 2025, with a total  
544 duration of approximately three (3) weeks.

### 545 **3.1.2 Feature Extraction Pipeline**

546 This stage directly followed the alignment phase, utilizing the resulting BAM files  
547 (specifically `chimeric.sorted.bam` and `clean.sorted.bam`). A custom Python  
548 script was created to efficiently process each primary-mapped read to extract  
549 the necessary set of features, which were then compiled into a structured feature  
550 matrix in TSV format. The pipeline’s core functionality relied on the `Pysam`  
551 (Version 0.22) library for parsing BAM structures and `NumPy` (Version 1.26) for  
552 array operations and computations. To ensure correctness and adherence to best  
553 practices, bioinformatics experts at PGC Visayas were consulted to validate the  
554 pipeline design, feature extraction logic, and overall data integrity.

555 This stage of the study was scheduled to begin in the last week of Novem-  
556 ber 2025 and conclude by the first week of December 2025, with an estimated



557 total duration of approximately two (2) weeks.

558 The pipeline focused on three feature families that collectively capture bi-  
559 ological signatures associated with PCR-induced chimeras: (1) supplementary  
560 alignment (SA) and alignment-structure metrics, (2) k-mer composition differ-  
561 ence, and (3) microhomology around putative junctions. Additional alignment  
562 quality indicators such as mapping quality were also included.

### 563 **Supplementary Alignment and Alignment-Structure Features**

564 Split-alignment information was derived from the SA tag embedded in each pri-  
565 mary read of the BAM file. This tag is typically associated with reads that map to  
566 multiple genomic locations, suggesting a chimeric structure. To extract this infor-  
567 mation, the script first checked whether the read carried an **SA:Z** tag. If present,  
568 the tag string was parsed using the function `parse_sa_tag`, yielding metadata for  
569 each alignment containing the reference name, mapped position, strand, mapping  
570 quality, and number of mismatches.

571 After parsing, the function `sa_feature_stats` was applied to establish the fun-  
572 damental split indicators, `has_sa` and `sa_count`. Along with these initial counts,  
573 the function aggregated metrics related to the structure and reliability of the  
574 split alignments, including the number of alignment segments, strand consistency,  
575 minimum, maximum, and mean distance between split segments, and summary  
576 statistics of mapping quality and mismatch counts across segments.

## 577 **K-mer Composition Difference**

578 Comparing k-mer frequency profiles between the left and right halves of a read  
579 allows for the detection of abrupt compositional shifts, independent of alignment  
580 information.

581 The script implemented this by inferring a likely junction breakpoint using the  
582 function `infer_breakpoints`, prioritizing the boundaries defined by soft-clipping  
583 operations. If no clipping was present, the midpoint of the alignment or the read  
584 length was used as a fallback. The read sequence was then divided into left and  
585 right segments at this inferred breakpoint, and k-mer frequency profiles ( $k =$   
586 6) were generated for both halves, ignoring any k-mers containing ambiguous N  
587 bases. The resulting k-mer frequency vectors were normalised and compared using  
588 the functions `cosine_difference` and `js_divergence` to quantify compositional  
589 discontinuity across the inferred breakpoint.

## 590 **Microhomology**

591 The process of extracting the microhomology feature also started by using  
592 `infer_breakpoints` to identify a candidate junction. Once a breakpoint was  
593 established, the script scanned a  $\pm 40$  base-pair window surrounding the break-  
594 point and applied the function `longest_suffix_prefix_overlap` to identify the  
595 longest exact suffix–prefix overlap between the left and right read segments. This  
596 overlap, representing consecutive bases shared at the junction, was recorded as  
597 `microhomology_length` in the dataset. The 40 base-pair window was chosen  
598 to ensure that short shared sequences at or near the breakpoint were captured

599 without including distant sequences that are unlikely to be mechanistically  
600 related.

601 Additionally, the GC content of the overlapping sequence was calculated using  
602 the function `gc_content`, which counts guanine (G) and cytosine (C) bases within  
603 the detected microhomology and divides by the total length, yielding a proportion  
604 between 0 and 1 that was stored under the `microhomology_gc` attribute. Micro-  
605 homology was quantified using a 3–20 bp window, consistent with values reported  
606 in prior research on PCR-induced chimeras. A k-mer length of 6 was used to cap-  
607 ture patterns within the 40 bp window surrounding each breakpoint, providing  
608 sufficient resolution to detect informative sequence shifts.

### 609 **3.1.3 Machine Learning Model Development**

610 After feature extraction, the per-read feature matrices for clean and chimeric  
611 reads were merged into a single dataset. Each row corresponded to one paired-  
612 end read, and columns encoded alignment-structure features (e.g., supplementary  
613 alignment count and spacing between segments), CIGAR-derived soft-clipping  
614 statistics (e.g., left and right soft-clipped length, total clipped bases), k-mer com-  
615 position discontinuity between read segments, microhomology descriptors near  
616 candidate junctions, and alignment quality (e.g., mapping quality). The result-  
617 ing feature set comprised 23 numeric features and was restricted to quantities  
618 that can be computed from standard BAM/FASTQ files in typical mitochondrial  
619 sequencing workflows.

620 The labelled dataset was randomly partitioned into training (80%) and test

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

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

### 640 **3.1.4 Model Benchmarking, Hyperparameter Optimiza-** 641 **tion, and Evaluation**

642 Model selection and refinement proceeded in two stages. First, the cross-validation  
 643 results from the broad panel were used to identify a subset of competitive mod-

644 els for more detailed optimization. Specifically, ten model families were carried  
645 forward:  $L_2$ -regularized logistic regression, calibrated linear SVM, Random For-  
646 est, ExtraTrees, Gradient Boosting, XGBoost, LightGBM, CatBoost, Bagging  
647 with decision trees, and a shallow MLP. This subset spans both linear and non-  
648 linear decision boundaries, but emphasizes ensemble and boosting methods, which  
649 showed superior F1 and ROC-AUC in the initial benchmark.

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

660 For each model family, the hyperparameter configuration with the highest  
661 mean cross-validation F1-score was selected as the best-tuned estimator. These  
662 tuned models were then refitted on the full training set and evaluated once on the  
663 held-out test set to obtain unbiased estimates of performance. Test-set metrics in-  
664 cluded accuracy, precision, recall, F1-score for the chimeric class, and ROC-AUC.  
665 Confusion matrices and ROC curves were generated for the top-performing mod-  
666 els to characterise common error modes, such as false negatives (missed chimeric  
667 reads) and false positives (clean reads incorrectly labelled as chimeric). The final  
668 model or small set of models for downstream interpretation was chosen based on

669 a combination of test-set F1-score and ROC-AUC.

### 670 **3.1.5 Feature Importance, Feature Selection, and Inter-** 671 **pretation**

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

682 For interpretability, individual features were grouped into conceptual families:  
683 (i) supplementary alignment and alignment-structure features (e.g., SA count,  
684 spacing between alignment segments, strand consistency), (ii) soft-clipping fea-  
685 tures (e.g., left and right soft-clipped length, total clipped bases, inferred break-  
686 point position), (iii) k-mer composition discontinuity features (e.g., cosine dis-  
687 tance and Jensen-Shannon divergence between k-mer profiles of read segments),  
688 (iv) microhomology descriptors (e.g., microhomology length and local GC content  
689 around putative breakpoints), and (v) other alignment quality features (e.g., map-  
690 ping quality). This analysis provided a basis for interpreting the trained models  
691 in terms of known mechanisms of PCR-induced template switching and for iden-

692 tifying which alignment-based and sequence-derived cues are most informative for  
693 distinguishing chimeric from clean mitochondrial reads.

694 Building on these importance results, an explicit feature selection step was  
695 implemented using CatBoost as the reference model, since it was among the top-  
696 performing classifiers. Permutation importance scores were re-estimated for Cat-  
697 Boost on the held-out test set using the F1-score of the chimeric class as the  
698 scoring function. Negative importance scores, which indicate that permuting a  
699 feature did not reliably harm performance, were set to zero and interpreted as  
700 noise. The remaining non-negative importances were sorted in descending order  
701 and converted into a cumulative importance curve by expressing each feature’s  
702 importance as a fraction of the total positive importance.

703 A compact feature subset was then defined by selecting the smallest number of  
704 features whose cumulative importance reached at least 95% of the total positive  
705 importance. This procedure yielded a reduced set of four strongly predictive  
706 variables dominated by soft-clipping and k-mer divergence metrics (for example,  
707 total clipped bases and k-mer divergence between read halves).

708 To quantify the impact of this reduction, CatBoost was retrained using only  
709 the selected feature subset, with the same tuned hyperparameters as the full 23-  
710 feature model, and evaluated on the held-out test set. Performance of the reduced  
711 model was then compared to that of the full model in terms of F1-score and ROC-  
712 AUC to assess whether dimensionality could be reduced without appreciable loss  
713 in predictive accuracy.

714 In addition, an ablation experiment was performed to specifically evaluate  
715 the contribution of explicit microhomology features. The microhomology vari-

ables (`microhomology_length` and `microhomology_gc`) were removed from the full feature set to obtain a 21-feature configuration. CatBoost was refitted on this microhomology-ablated feature set, using the same tuned hyperparameters, and evaluated on the held-out test set. Comparing the full, reduced-subset, and microhomology-ablated variants allowed the study to quantify both the degree of redundancy among features and the practical contribution of microhomology to classification accuracy.

Taken together, the feature importance and feature selection analyses provided a more parsimonious model variant and a clearer interpretation of which alignment-based and sequence-derived signals are most informative for detecting PCR-induced chimeras.

### 3.1.6 Validation and Testing

Validation involved both internal and external evaluations. Internal validation was achieved through five-fold stratified cross-validation on the training data to verify model generalization and reduce variance due to random sampling. External testing was performed on the 20% hold-out dataset from the simulated reads, providing an unbiased assessment of model generalization. Feature extraction and preprocessing were applied consistently across all splits.

Comparative evaluation was performed across all candidate algorithms and CatBoost feature-set variants to determine which models demonstrated the highest predictive performance and computational efficiency under identical data conditions. Their metrics were compared to identify which algorithms and feature



738 configurations were most suitable for further refinement and potential integration  
739 into downstream mitochondrial assembly workflows.

### 740 **3.1.7 Documentation**

741 Comprehensive documentation was maintained throughout the study to ensure  
742 transparency and reproducibility. All stages of the research, including data gath-  
743 ering, preprocessing, feature extraction, model training, feature selection, and  
744 validation, were systematically recorded in a **README** file in the GitHub reposi-  
745 tory. For each analytical step, the corresponding parameters, software versions,  
746 and command line scripts were documented to enable exact replication of results.

747 The repository structure followed standard research data management prac-  
748 tices, with clear directories for datasets and scripts. Computational environments  
749 were standardised using Conda, with an environment file (**environment.yml**)  
750 specifying dependencies and package versions to maintain consistency across sys-  
751 tems.

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

## 754 **3.2 Calendar of Activities**

755 Table 3.1 presents the project timeline in the form of a Gantt chart, where each  
756 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

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. The behaviour of the extracted features is first examined through descriptive and correlation analyses, followed by a comparison of baseline and tuned classifiers. The chapter then examines model performance in detail and investigates the contribution of individual features and feature families, including the impact of feature selection on classification performance.

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).

## 769 4.1 Descriptive Analysis of Features

### 770 4.1.1 Summary Statistics Per Class

771 Summary statistics were computed separately for clean reads (class 0) and  
772 chimeric reads (class 1) to characterize the distributional behavior of the features.  
773 For each feature, the mean, standard deviation, median, first and third quartiles  
774 (Q1, Q3), interquartile range (IQR), minimum, maximum, and sample size ( $n$ )  
775 were calculated.

776 Only a subset of the features is summarized in the main text to highlight key  
777 trends, and not all summary statistics columns are shown for brevity. The com-  
778 plete set of per-class summary statistics for all features is provided in Appendix A  
779 (Table A.1).

### 780 Alignment and Supplementary Alignment Features

781 Features related to supplementary alignments show strong separation between  
782 classes. Chimeric reads frequently exhibit supplementary alignments, reflected  
783 by higher values of `has_sa`, `sa_count`, and `num_segments`, whereas clean reads  
784 consistently show a single alignment segment with no supplementary mappings.  
785 Table 4.1 shows that `has_sa` is present in chimeric reads (mean = 0.406) but absent  
786 in clean reads (mean = 0.000), while `num_segments` increases from a constant value  
787 of 1.000 in clean reads to a mean of 1.406 in chimeric reads. These patterns align  
788 with the expected structure of chimeric reads and indicate that alignment-based  
789 features are highly informative.

## 790 Clipping-Based Features

791 Clipping-related features, including `softclip_left`, `softclip_right`, and  
792 `total_clipped_bases`, display higher values and broader distributions in chimeric  
793 reads. In chimeric reads, `total_clipped_bases` reaches 25.44 on average, with a  
794 median of 19.0 and an IQR of 48.0, while `softclip_left` and `softclip_right`  
795 have averages of 12.55 and 12.90, medians of 0.0, and IQRs of 19.0. Clean  
796 reads maintain values near zero across all these metrics. These patterns indi-  
797 cate substantial clipping and increased variability in chimeric reads, reflecting  
798 junction-like alignment fragmentation, whereas clean reads remain unaltered.

## 799 K-mer Distribution Features

800 K-mer-based features, including `kmer_js_divergence` and `kmer_cosine_diff`,  
801 show only minor differences between clean and chimeric reads. In chimeric  
802 reads, `kmer_js_divergence` has a mean of 0.974 with a median of 0.986, and  
803 `kmer_cosine_diff` has a mean of 0.974 with a median of 0.986. Clean reads show  
804 similar values, with `kmer_js_divergence` at 0.976 with a median of 0.986, and  
805 `kmer_cosine_diff` at 0.976 with a median of 0.986. The close similarity of the  
806 means, medians, and overall ranges of values indicates that these features alone  
807 provide limited ability to distinguish clean from chimeric reads.

## 808 Microhomology Features

809 Microhomology-related features, including `microhomology_length` and  
810 `microhomology_gc`, exhibit nearly identical summary statistics between clean

811 and chimeric reads. Most reads in both classes have short or zero-length micro-  
812 homologies. Table 4.1 shows that `microhomology_gc` has a mean of 0.172 and  
813 a median of 0.0 in both clean and chimeric reads, while `microhomology_length`  
814 averages 0.458 with a median of 0.0 in chimeric reads and 0.462 with a median  
815 of 0.0 in clean reads. These values indicate that microhomology features alone  
816 provide limited discriminatory power and are more appropriately considered as  
817 supporting evidence.

818 Overall, the summary statistics indicate that alignment-based and clipping-  
819 based features provide the strongest class separation, k-mer features contribute  
820 limited but complementary signal, and microhomology features exhibit minimal  
821 discriminative power on their own. These observations motivate the combined  
822 multi-feature approach used in subsequent modeling and evaluation.

Table 4.1: Summary statistics of selected key features by class.

Feature	Class	Mean	Std	Median	IQR
has_sa	chimeric	0.406	0.491	0.0	1.0
has_sa	clean	0.000	0.000	0.0	0.0
num_segments	chimeric	1.406	0.491	1.0	1.0
num_segments	clean	1.000	0.000	1.0	0.0
softclip_left	chimeric	12.55	21.90	0.0	19.0
softclip_left	clean	0.23	1.54	0.0	0.0
softclip_right	chimeric	12.90	22.12	0.0	19.0
softclip_right	clean	0.21	1.51	0.0	0.0
total_clipped_bases	chimeric	25.44	25.48	19.0	48.0
total_clipped_bases	clean	0.44	2.16	0.0	0.0
kmer_js_divergence	chimeric	0.974	0.025	0.986	0.043
kmer_js_divergence	clean	0.976	0.025	0.986	0.040
kmer_cosine_diff	chimeric	0.974	0.026	0.986	0.042
kmer_cosine_diff	clean	0.976	0.025	0.986	0.041
microhomology_length	chimeric	0.458	0.755	0.0	1.0
microhomology_length	clean	0.462	0.758	0.0	1.0
microhomology_gc	chimeric	0.172	0.361	0.0	0.0
microhomology_gc	clean	0.172	0.361	0.0	0.0

Boxplots were generated for each feature, with the x-axis representing the class (clean reads and chimeric reads) and the y-axis representing the feature value. Figure 4.1 presents a panel of selected key features, while boxplots for all numeric features are provided in Appendix B.

For clipping-related features (`softclip_left`, `softclip_right`, and `total_clipped_bases`), chimeric reads exhibit higher medians and longer upper whiskers than clean reads, indicating increased variability and the presence of split alignments.

Supplementary alignment features (`has_sa` and `sa_count`), show that clean reads are largely zero, whereas chimeric reads display a wider distribution, re-

833 flecting frequent supplementary alignments.

834 K-mer metrics (`kmer_js.divergence` and `kmer_cosine.diff`) show a slight  
 835 upward shift for chimeric reads, but substantial overlap with clean reads indicates  
 836 low discriminative power.

837 Microhomology features (`microhomology_length` and `microhomology_gc`)  
 838 have nearly overlapping distributions for both classes, consistent with their low  
 839 standalone predictive importance.

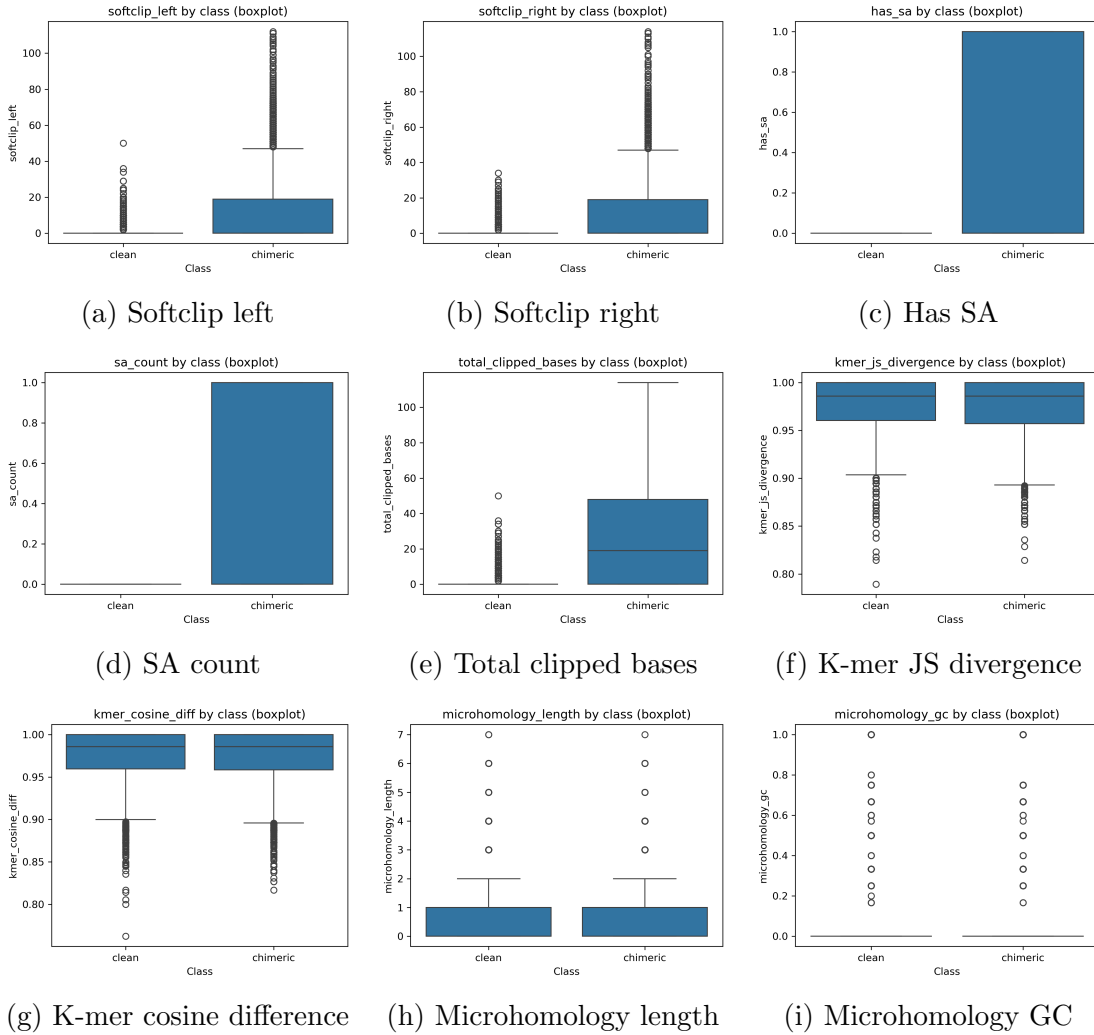


Figure 4.1: Boxplots of selected features for clean and chimeric reads.



## 840 4.1.2 Correlation Analysis of Extracted Features

841 A feature correlation heatmap (Figure 4.2) was generated to examine relationships  
842 among the extracted variables and to identify patterns of redundancy and inde-  
843 pendence within the feature set. The analysis shows that alignment-related and  
844 clipping-related features form a strongly correlated cluster, including indicators  
845 of supplementary alignments, alignment segment counts, positional differences,  
846 and soft-clipping measures. These features capture related aspects of alignment  
847 fragmentation, which is a known characteristic of chimeric reads, and several  
848 show moderate correlations with the class label, supporting their relevance for  
849 distinguishing chimeric from clean reads. In contrast, general read-quality and  
850 alignment-quality metrics, such as read length, base quality, and mapping qual-  
851 ity, exhibit weak correlations with most split-alignment features, indicating that  
852 they provide distinct information rather than overlapping with alignment-derived  
853 signals. Sequence-based features display a similar pattern of independence, as  
854 k-mer divergence metrics show weak correlations with other feature groups, while  
855 microhomology features exhibit generally low correlations with both alignment-  
856 based and k-mer-based features. Overall, the correlation structure highlights in-  
857 tentional redundancy within alignment-derived features and clear separation be-  
858 tween feature families, supporting the use of features that capture different aspects  
859 of chimeric read characteristics to improve chimera classification.

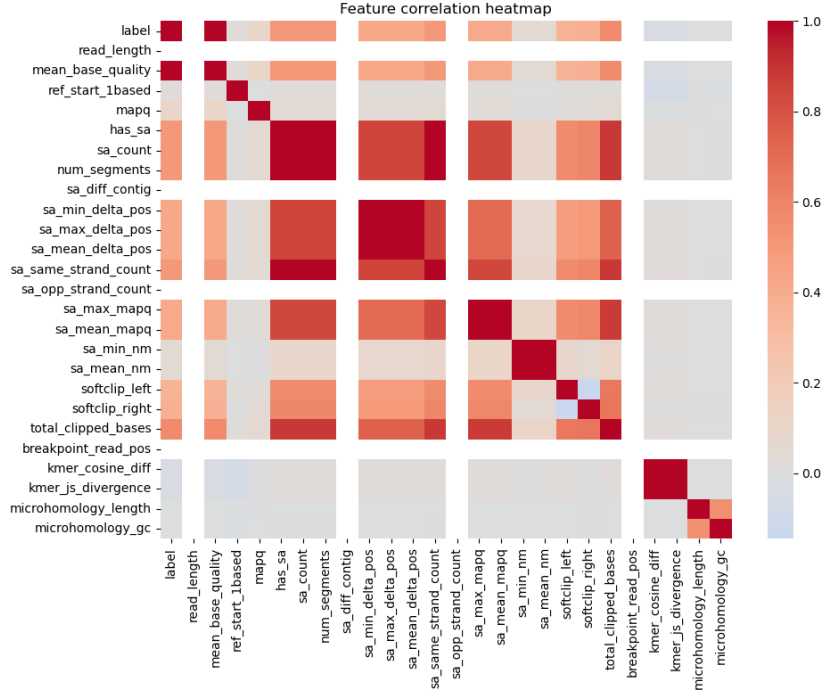


Figure 4.2: Feature correlation heatmap showing relationships among alignment-derived and sequence-derived features.

## 4.2 Baseline Classification Performance

Table ?? summarises the performance of thirteen baseline classifiers trained on the engineered feature set and evaluated on a held-out test set. All models were trained using default hyperparameters, without dedicated tuning.

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

Across the remaining models, test F1-scores clustered within a narrow range,

869 from approximately 0.75 to 0.78, with ROC–AUC values between about 0.82  
 870 and 0.85. Ensemble methods, including gradient boosting, CatBoost, LightGBM,  
 871 XGBoost, bagging trees, and random forest, exhibited very similar performance.  
 872 Among these, CatBoost and gradient boosting achieved the highest scores, with  
 873 test F1-scores of approximately 0.775 and ROC–AUC values of approximately  
 874 0.84. Linear models, namely logistic regression and the calibrated linear SVM,  
 875 performed slightly worse, with test F1-scores around 0.75. In contrast, Gaussian  
 876 Naive Bayes lagged behind with a substantially lower F1-score of approximately  
 877 0.66, despite exhibiting extremely high precision for the chimeric class.

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

model	test_accuracy	test_precision	test_recall	test_f1	test_roc_auc
dummy_baseline	0.500188	0.500188	1.000000	0.666833	0.500000
logreg_l2	0.790797	0.945956	0.617000	0.746860	0.829807
linear_svm_calibrated	0.791422	0.947773	0.617000	0.747426	0.829602
random_forest	0.800050	0.910427	0.665750	0.769097	0.832766
extra_trees	0.797924	0.918833	0.653750	0.763950	0.826517
gradient_boosting	0.809053	0.947521	0.654500	0.774213	0.844844
xgboost	0.807303	0.942107	0.655000	0.772747	0.841042
lightgbm	0.806052	0.936231	0.657000	0.772146	0.841671
catboost	0.808803	0.941408	0.658750	0.775114	0.843362
knn	0.789671	0.902990	0.649250	0.755381	0.820898
gaussian_nb	0.745780	0.997975	0.492750	0.659749	0.826918
bagging_trees	0.800550	0.910830	0.666500	0.769742	0.837357
mlp	0.793047	0.949062	0.619500	0.749660	0.829611

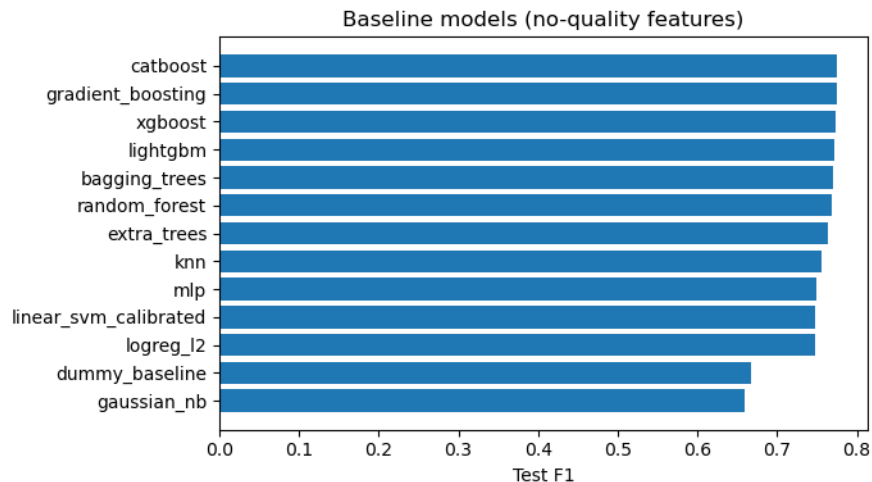


Figure 4.3: Test F1 of all baseline classifiers.

## 4.3 Effect of Hyperparameter Tuning

To assess whether performance could be improved further, ten model families underwent randomised hyperparameter search. The tuned metrics are summarised in Table 4.3. 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, random forest, bagging trees, and extra trees experienced small increases in test F1 after tuning, typically on the order of  $\Delta F1 \approx 0.002$ – $0.006$ , with corresponding improvements in ROC–AUC of up to approximately  $\Delta AUC \approx 0.009$ . In contrast, XGBoost and the multilayer perceptron showed negligible change or slight decreases in F1, while linear models did not benefit from tuning.

After tuning, gradient boosting achieved the best overall test performance,

with a test F1-score of 0.776 and a ROC-AUC of 0.846. LightGBM and bagging trees followed closely, attaining test F1-scores of 0.774 and 0.772 and ROC-AUC values of 0.843 and 0.842, respectively. Random forest also improved modestly to a test F1-score of 0.772 with a ROC-AUC of 0.839. CatBoost, with a test F1-score of 0.775 and ROC-AUC of 0.843, achieved marginal changes relative to its baseline performance.

Table 4.3: Performance of tuned classifiers on the held-out test set.

model_name	test_f1_base	test_roc_auc_base	test_f1_tuned	test_roc_auc_tuned
gradient_boosting	0.774213	0.844844	0.776460	0.845858
catboost	0.775114	0.843362	0.775289	0.842918
lightgbm	0.772146	0.841671	0.773802	0.843451
bagging_trees	0.769742	0.837357	0.772422	0.841870
random_forest	0.769097	0.832766	0.772376	0.838799
xgboost	0.772747	0.841042	0.770118	0.843225
extra_trees	0.763950	0.826517	0.769878	0.834912
mlp	0.749660	0.829611	0.749167	0.828506
logreg_l2	0.746860	0.829807	0.745187	0.825632
linear_svm_calibrated	0.747426	0.829602	0.744848	0.825147

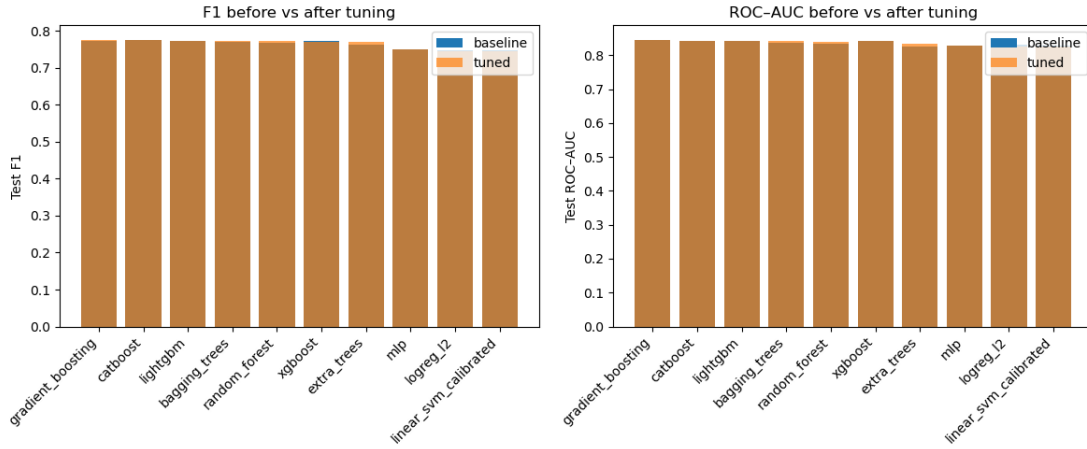


Figure 4.4: Comparison of test F1 (left) and ROC-AUC (right) for baseline and tuned models.

Because improvements are small and within cross-validation variability, tun-

ing was interpreted as stabilising and slightly refining the models rather than completely 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.

### 4.4.1 Confusion Matrices and Error Patterns

Classification reports and confusion matrices for the four models reveal consistent patterns. CatBoost and gradient boosting both achieved overall accuracy around 0.81, with similar macro-averaged F1 scores (0.80–0.805). For CatBoost, precision and recall for clean reads were 0.74 and 0.95, respectively, while for chimeric reads they were 0.94 and 0.66 ( $F1 = 0.775$ ). Gradient boosting showed nearly identical trade-offs, with clean read precision/recall of 0.74/0.96 and chimeric read precision/recall of 0.94/0.66 ( $F1 = 0.777$ ).

Bagging trees achieved slightly lower accuracy (0.805) and chimeric F1 (0.772), whereas the multilayer perceptron (MLP) attained the lowest accuracy (0.793) and

918 chimeric F1 (0.749), despite achieving high chimeric precision (0.95) at the cost  
919 of lower recall (0.62).

920     Across all models, errors were asymmetric: false negatives (chimeric reads  
921 predicted as clean) were more frequent than false positives. For instance, CatBoost  
922 misclassified 1,352 chimeric reads as clean but only 215 clean reads as chimeric,  
923 while gradient boosting misclassified 1,352 chimeric reads as clean and 181 clean  
924 reads as chimeric. This pattern indicates that both models are conservative,  
925 prioritizing the avoidance of false chimera calls even if some true chimeras are  
926 missed. Consultation with PGC Visayas suggested that this conservative behavior  
927 is generally acceptable, although further evaluation is needed to assess its impact  
928 on downstream analyses.

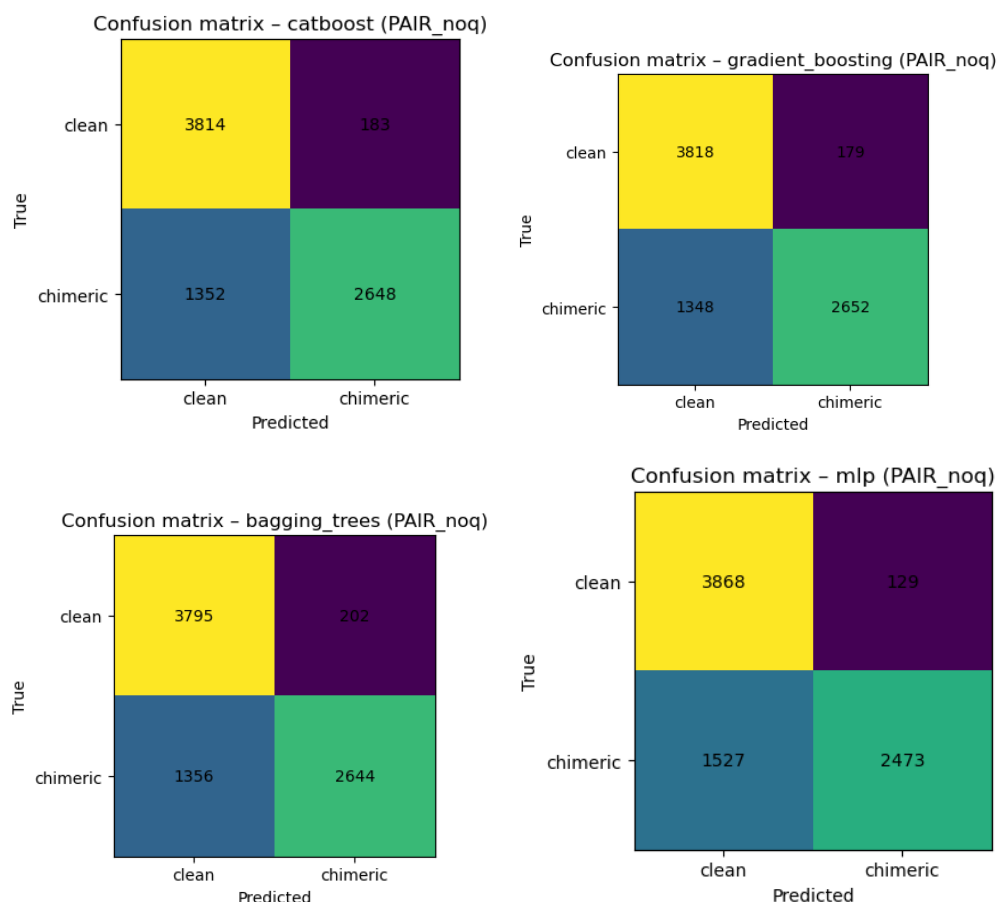


Figure 4.5: Confusion matrices for the four representative models on the held-out test set.

## 929 4.4.2 ROC and Precision–Recall Curves

930 Receiver operating characteristic (ROC) and precision–recall (PR) curves as  
 931 shown in Figure 4.6 further support the similarity among the top models. The  
 932 three tree-based ensembles (CatBoost, gradient boosting, bagging trees) achieved  
 933 ROC–AUC values of approximately 0.84 and average precision (AP) around 0.88.  
 934 MLP performed slightly worse ( $AUC \approx 0.82$ ,  $AP \approx 0.87$ ) but still substantially  
 935 better than the dummy baseline.



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

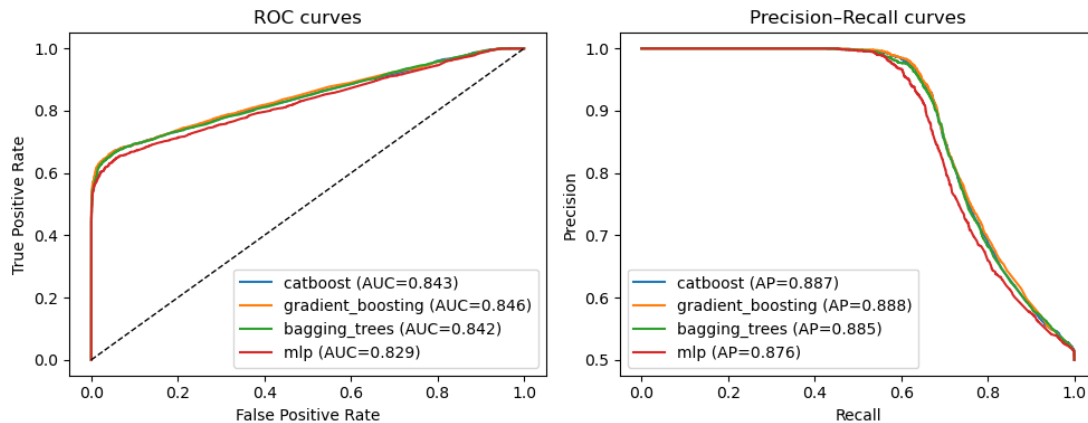


Figure 4.6: ROC (left) and precision–recall (right) curves for the four representative models on the held-out test set.

## 941 4.5 Feature Importance

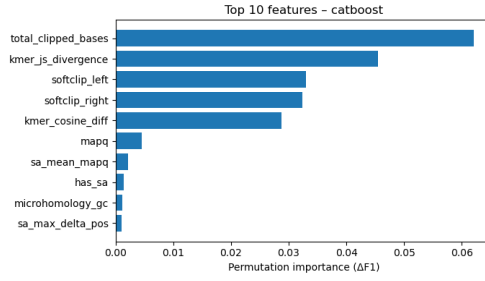
### 942 4.5.1 Permutation Importance of Individual Features

943 To understand how each classifier made predictions, feature importance was quan-  
 944 tified using permutation importance. This analysis was applied to four represen-  
 945 tative models: CatBoost, Gradient Boosting, Bagging Trees, and an MLP.

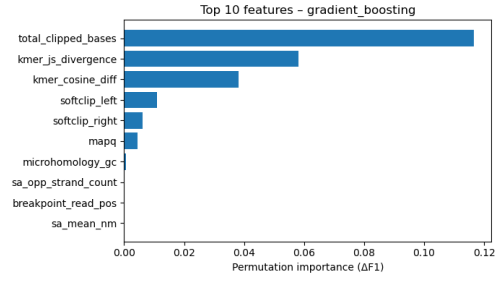
946 As shown in Figure 4.7, `total_clipped_bases` consistently provides a  
 947 strong predictive signal across all models, particularly in Gradient Boost-

948 ing (importance = 0.117) and Bagging Trees (importance = 0.274). Cat-  
949 Boost assigns high importance to both `total_clipped_bases` (0.062) and  
950 `kmer_js_divergence` (0.045), while MLP relies on `total_clipped_bases` and  
951 soft-clipping features (`softclip_left`, `softclip_right`) as primary signals. Gra-  
952 dient Boosting emphasizes `kmer_js_divergence` and `kmer_cosine_diff` alongside  
953 `total_clipped_bases`, but soft-clipping features contribute less.

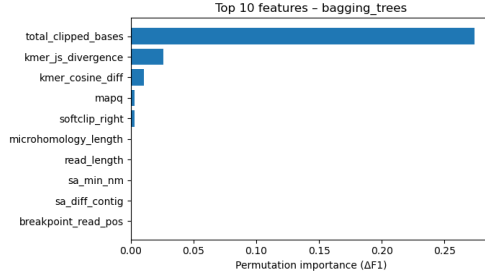
954     Microhomology features (`microhomology_length` and `microhomology_gc`)  
955 provide minimal predictive value in all models, and some alignment-based split-  
956 read metrics (e.g., `sa_min_delta_pos`, `sa_max_delta_pos`) are leveraged primarily  
957 by the MLP. Overall, these results indicate that accurate detection of chimeric  
958 reads relies on both alignment-based signals and k-mer compositional information,  
959 with explicit microhomology features contributing little. Combining multiple  
960 feature types enhances model sensitivity and specificity.



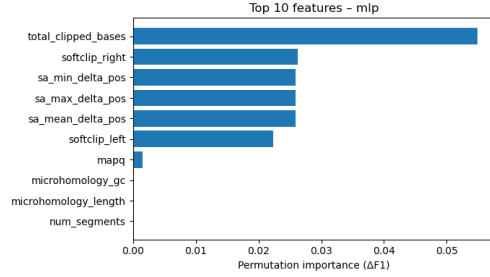
(a) CatBoost



(b) Gradient Boosting



(c) Bagging Trees



(d) Multi-Layer Perceptron (MLP)

Figure 4.7: Permutation-based feature importance for four representative classifiers.

## 961 4.5.2 Feature Family Importance

962 To evaluate broader predictive signals, features were grouped into five fam-  
 963 ilies: SA\_structure (supplementary alignment and segment metrics, e.g.,  
 964 has\_sa, sa\_count, sa\_min\_delta\_pos, sa\_mean\_nm), Clipping (softclip\_left,  
 965 softclip\_right, total\_clipped\_bases, breakpoint\_read\_pos), Kmer\_jump  
 966 (kmer\_cosine\_diff, kmer\_js\_divergence), Micro\_homology (microhomology\_length,  
 967 microhomology\_gc), and Other (e.g., mapq).

968 Aggregated analyses reveal consistent patterns across models. In CatBoost,  
 969 the Clipping family dominates with cumulative importance 0.127, followed  
 970 by Kmer\_jump (0.074), while Other (0.0045), SA\_structure (0.0033), and Mi-

cro\_homology (0.0013) contribute minimally. Gradient Boosting shows a similar trend, with Clipping (0.134) and Kmer\_jump (0.096) providing most predictive power, and the remaining families contributing negligibly. Bagging Trees emphasizes Clipping even more strongly (0.277), with Kmer\_jump secondary (0.037), and SA\_structure, Micro\_homology, and Other remaining minor contributors. Interestingly, the MLP exhibits a different pattern, prioritizing Clipping (0.104) and SA\_structure (0.078), while Kmer\_jump (0.000034) and Micro\_homology (0.000091) have almost no effect.

Both feature-level and aggregated analyses indicate that accurate detection of chimeric reads in this dataset relies primarily on alignment irregularities captured by Clipping features and, in most tree-based models, on k-mer compositional shifts (Kmer\_jump), which often arise from PCR-induced template switching events. Explicit microhomology features contribute minimally, and some reliance on SA\_structure signals is observed only in the MLP.

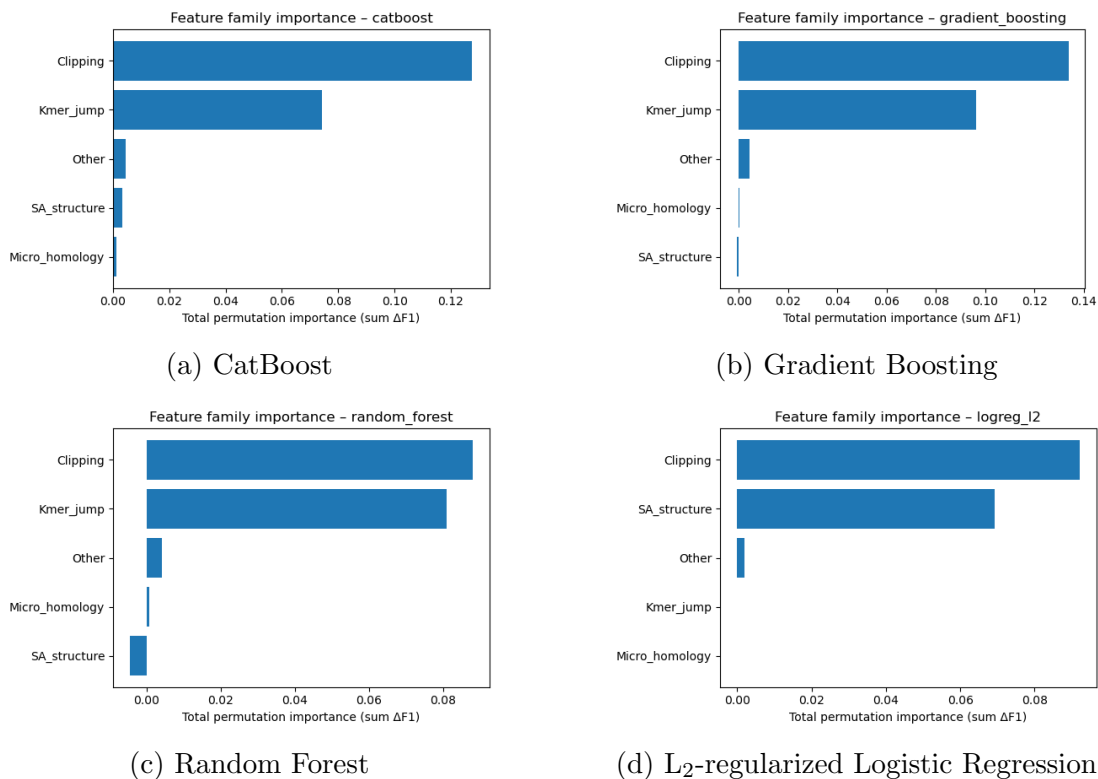


Figure 4.8: Aggregated feature family importance across four models.

## 985 4.6 Feature Selection

986 Feature selection was performed to identify the smallest subset reaching 95% cu-  
 987 mulative importance. Three models were evaluated as references: the full model  
 988 with all 23 features, a reduced model with the top- $k$  features, and an ablation  
 989 model excluding microhomology features, using a tuned CatBoost classifier to  
 990 assess feature contributions and overall classification performance.

### 991 4.6.1 Cumulative Importance Curve

992 The cumulative importance curve was computed using the tuned Gradient Boost-  
993 ing classifier. Figure 4.9 illustrates the contribution of features sorted by impor-  
994 tance. The curve rises steeply for the top features and then gradually plateaus,  
995 indicating that a small number of features capture most of the model’s pre-  
996 dictive power. A cumulative importance of 95% is reached at  $k = 4$  features,  
997 which are `total_clipped_bases`, `kmer_js_divergence`, `kmer_cosine_diff`, and  
998 `softclip_left`.

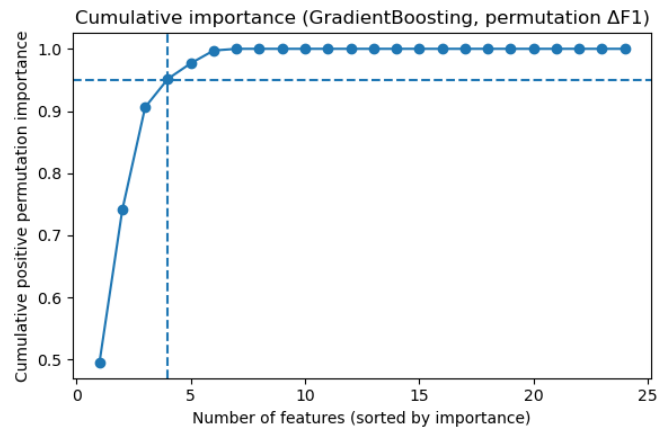


Figure 4.9: Cumulative importance curve of features sorted by importance.

### 999 4.6.2 Performance Comparison Across Feature Sets

1000 Classification performance was compared across three feature sets using a tuned  
1001 Gradient Boosting classifier. The full model, incorporating all 24 engineered fea-  
1002 tures, achieved an F1 score of 0.7765 and a ROC-AUC of 0.8459. A reduced model  
1003 using only the top four features (`total_clipped_bases`, `kmer_js_divergence`,  
1004 `kmer_cosine_diff`, and `softclip_left`) achieved nearly equivalent performance

1005 with an F1 of 0.7768 and a ROC-AUC of 0.8369. An ablation model excluding mi-  
 1006 crohomology features (`microhomology_length` and `microhomology_gc`) also per-  
 1007 formed comparably, with an F1 of 0.7761 and ROC-AUC of 0.8444. These results  
 1008 indicate that clipping and k-mer features capture almost all predictive signal,  
 1009 while microhomology features are largely redundant in this dataset.

Table 4.4: Test set performance of three feature set variants using tuned Gradient Boosting.

Variant	No. of Features	Test F1	ROC-AUC
Full Gradient Boost	24	0.7765	0.8459
Selected (top-4)	4	0.7768	0.8369
No microhomology	22	0.7761	0.8444

1010 Figure 4.10 presents a bar chart comparing F1 and ROC-AUC across the  
 1011 three variants, with the x-axis showing the model variants and two bars per group  
 1012 representing the F1 and ROC-AUC values.

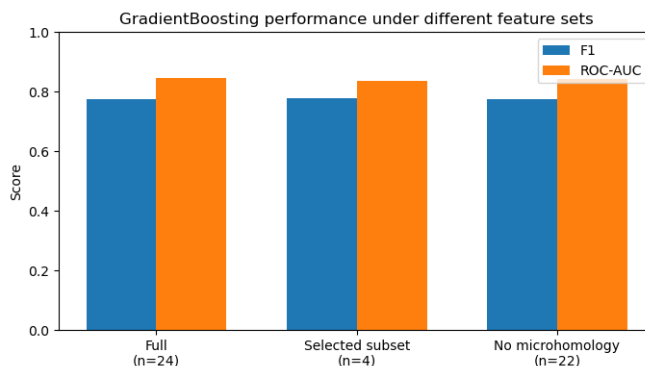


Figure 4.10: Comparison of F1 and ROC-AUC for the full, top-4 selected, and no-microhomology feature set variants.

### 1013 4.6.3 Interpretation and Final Feature Set Choice

1014 The full 23-feature model is retained as the primary configuration for the re-  
1015 mainder of the study, while the four-feature subset serves as a lightweight al-  
1016 ternative. Clipping features reflect alignment junctions and mapping disruptions  
1017 typical of chimeric reads, and k-mer divergence captures changes in sequence com-  
1018 position across breakpoints. Microhomology features appear largely redundant,  
1019 as their signal is either indirectly represented by clipping and k-mer features or  
1020 not strongly expressed in the simulation dataset.

## 1021 4.7 Summary of Findings

1022 All evaluated machine learning models substantially outperformed the dummy  
1023 baseline, demonstrating that the engineered feature set contains meaningful  
1024 signals for detecting PCR-induced chimeric reads. Across classifiers, the best-  
1025 performing models achieved test F1-scores of approximately 0.77 and ROC-AUC  
1026 values around 0.84 on held-out simulated mitochondrial reads, indicating reli-  
1027 able discrimination between clean and chimeric sequences. Among the tested  
1028 approaches, tree-based ensemble and boosting methods consistently showed the  
1029 strongest and most stable performance. In particular, CatBoost and Gradient  
1030 Boosting ranked among the top models across multiple evaluation metrics,  
1031 both before and after hyperparameter tuning. These results suggest that non-  
1032 linear ensemble methods are well suited to capturing the interaction between  
1033 alignment-derived and sequence-derived features in this setting.

1034 Analysis of feature behaviour revealed clear differences in how effectively fea-



1035 ture groups distinguished clean and chimeric reads. Alignment- and clipping-  
1036 based features, such as soft-clipping measures and total clipped bases, showed  
1037 strong separation between clean and chimeric reads and emerged as the most  
1038 informative signals. K-mer divergence features provided additional but weaker  
1039 separation, contributing complementary information beyond alignment irregular-  
1040 ities. In contrast, microhomology features and several supplementary alignment  
1041 (SA) structure metrics exhibited minimal class separation and contributed little  
1042 to overall predictive performance.

1043 Feature selection results further supported these observations. A reduced sub-  
1044 set of four features, dominated by clipping-based and k-mer divergence metrics,  
1045 achieved nearly identical performance to the full 23-feature model. Moreover,  
1046 removing explicit microhomology features did not degrade performance and in  
1047 some cases resulted in slightly improved metrics, suggesting that these features  
1048 are largely redundant under the simulated conditions tested.

1049 Overall, these findings suggest that alignment-based and k-mer-based fea-  
1050 tures provide sufficient signal to detect PCR-induced chimeric reads in simulated  
1051 mitochondrial data, supporting the use of a compact and interpretable machine  
1052 learning approach as a pre-assembly chimera detection step.

# 1053 **Appendix A**

## 1054 **Complete Per-Class Summary**

### 1055 **Statistics**

Table A.1: Complete per-class summary statistics for all extracted features.

Feature	Class	Mean	Std	Median	Q1	Q3	IQR	Min	Max	n
breakpoint_read_pos	chimeric	75.000	0.000	75.000	75.000	75.000	0.000	75.000	75.000	20000
breakpoint_read_pos	clean	75.000	0.000	75.000	75.000	75.000	0.000	75.000	75.000	19983
has_sa	chimeric	0.406	0.491	0.000	0.000	1.000	1.000	0.000	1.000	20000
has_sa	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
kmer_cosine_diff	chimeric	0.974	0.026	0.986	0.958	1.000	0.042	0.817	1.000	20000
kmer_cosine_diff	clean	0.976	0.025	0.986	0.959	1.000	0.041	0.814	1.000	19983
kmer_js_divergence	chimeric	0.974	0.025	0.986	0.957	1.000	0.043	0.811	1.000	20000
kmer_js_divergence	clean	0.976	0.025	0.986	0.959	1.000	0.040	0.817	1.000	19983
mapq	chimeric	59.987	0.355	60.000	60.000	60.000	0.000	43.000	60.000	20000
mapq	clean	59.663	2.036	60.000	60.000	60.000	0.000	0.000	60.000	19983
mean_base_quality	chimeric	40.000	0.000	40.000	40.000	40.000	0.000	40.000	40.000	20000
mean_base_quality	clean	13.000	0.000	13.000	13.000	13.000	0.000	13.000	13.000	19983
microhomology_gc	chimeric	0.172	0.361	0.000	0.000	0.000	0.000	0.000	1.000	20000
microhomology_gc	clean	0.172	0.361	0.000	0.000	0.000	0.000	0.000	1.000	19983
microhomology_length	chimeric	0.458	0.755	0.000	0.000	1.000	1.000	0.000	5.000	20000
microhomology_length	clean	0.462	0.758	0.000	0.000	1.000	1.000	0.000	5.000	19983

*Continued on next page*

Feature	Class	Mean	Std	Median	Q1	Q3	IQR	Min	Max	n
num_segments	chimeric	1.406	0.491	1.000	1.000	2.000	1.000	1.000	2.000	20000
num_segments	clean	1.000	0.000	1.000	1.000	1.000	0.000	1.000	1.000	19983
read_length	chimeric	150.000	0.000	150.000	150.000	150.000	0.000	150.000	150.000	20000
read_length	clean	150.000	0.000	150.000	150.000	150.000	0.000	150.000	150.000	19983
ref_start_1based	chimeric	8428.635	4248.348	8433.000	5013.000	11786.250	6773.250	1.000	16521.000	20000
ref_start_1based	clean	8200.121	4626.918	8240.000	3639.000	11565.000	7926.000	1.000	16521.000	19983
sa_count	chimeric	0.406	0.491	0.000	0.000	1.000	1.000	0.000	1.000	20000
sa_count	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_diff_contig	chimeric	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	20000
sa_diff_contig	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_max_delta_pos	chimeric	1573.531	2364.996	0.000	0.000	2826.250	2826.250	0.000	16519.000	20000
sa_max_delta_pos	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_max_mapq	chimeric	14.104	21.424	0.000	0.000	27.000	27.000	0.000	60.000	20000
sa_max_mapq	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_mean_delta_pos	chimeric	1573.531	2364.996	0.000	0.000	2826.250	2826.250	0.000	16519.000	20000
sa_mean_delta_pos	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_mean_mapq	chimeric	14.104	21.424	0.000	0.000	27.000	27.000	0.000	60.000	20000
sa_mean_mapq	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983

*Continued on next page*

Feature	Class	Mean	Std	Median	Q1	Q3	IQR	Min	Max	n
sa_mean_nm	chimeric	0.022	0.319	0.000	0.000	0.000	0.000	0.000	6.000	20000
sa_mean_nm	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_min_delta_pos	chimeric	1573.531	2364.996	0.000	0.000	2826.250	2826.250	0.000	16519.000	20000
sa_min_delta_pos	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_min_nm	chimeric	0.022	0.319	0.000	0.000	0.000	0.000	0.000	6.000	20000
sa_min_nm	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_opp_strand_count	chimeric	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	20000
sa_opp_strand_count	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
sa_same_strand_count	chimeric	0.406	0.491	0.000	0.000	1.000	1.000	0.000	1.000	20000
sa_same_strand_count	clean	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	19983
softclip_left	chimeric	12.546	21.898	0.000	0.000	19.000	19.000	0.000	150.000	20000
softclip_left	clean	0.225	1.543	0.000	0.000	0.000	0.000	0.000	56.000	19983
softclip_right	chimeric	12.896	22.123	0.000	0.000	19.000	19.000	0.000	150.000	20000
softclip_right	clean	0.212	1.513	0.000	0.000	0.000	0.000	0.000	55.000	19983
total_clipped_bases	chimeric	25.442	25.481	19.000	0.000	48.000	48.000	0.000	150.000	20000
total_clipped_bases	clean	0.437	2.157	0.000	0.000	0.000	0.000	0.000	110.000	19983

## Appendix B

### Boxplots for All Numeric Features by Feature Family

#### B.0.1 SA Structure (Supplementary Alignment and Segment Metrics)

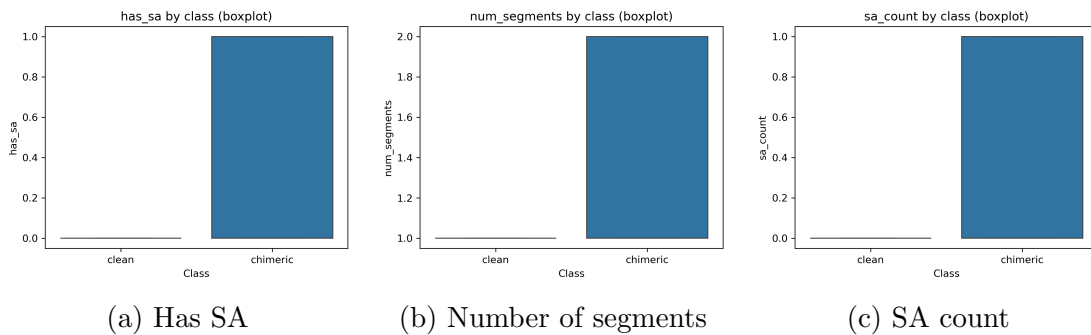
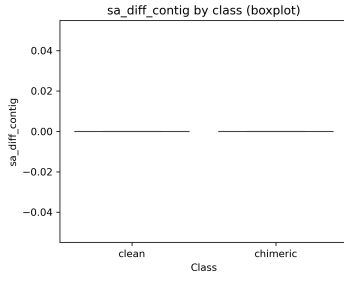
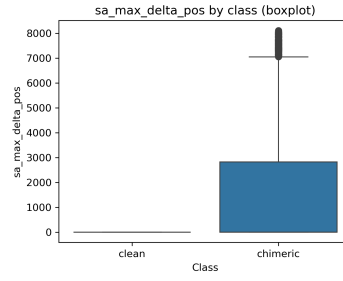


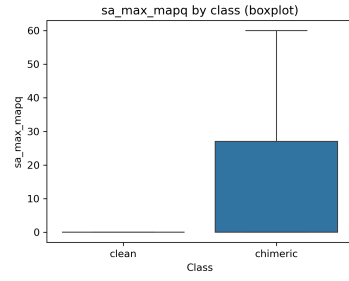
Figure B.1: Boxplots of SA Structure features by class (1/2).



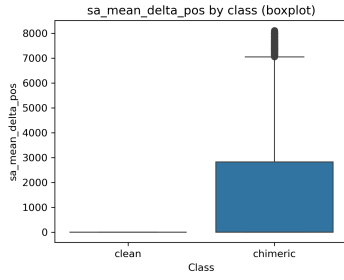
(a) SA different contig



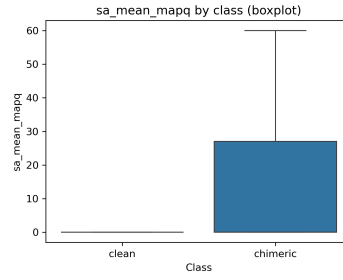
(b) SA max  $\Delta$  position



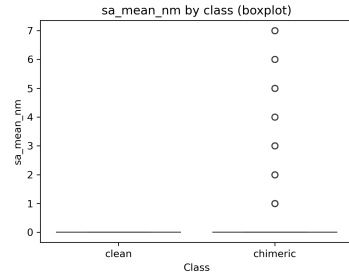
(c) SA max MAPQ



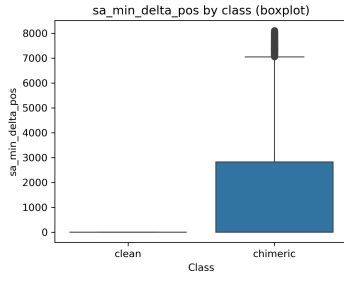
(d) SA mean  $\Delta$  position



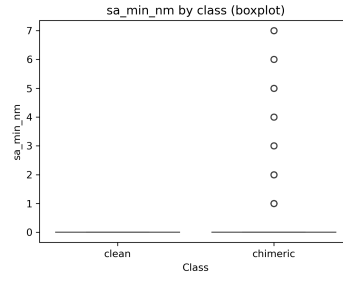
(e) SA mean MAPQ



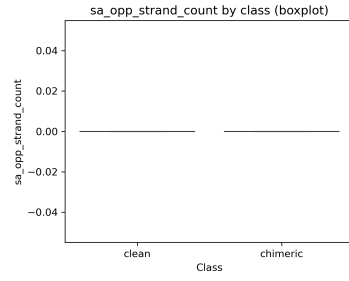
(f) SA mean NM



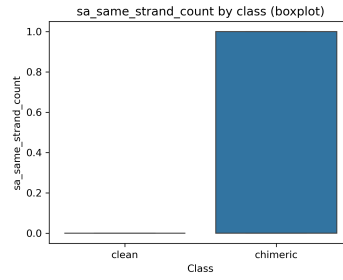
(g) SA min  $\Delta$  position



(h) SA min NM



(i) SA opposite strand count



(j) SA same strand count

Figure B.2: Boxplots of SA Structure features by class (2/2).

## 1062 B.0.2 Clipping-Based Features

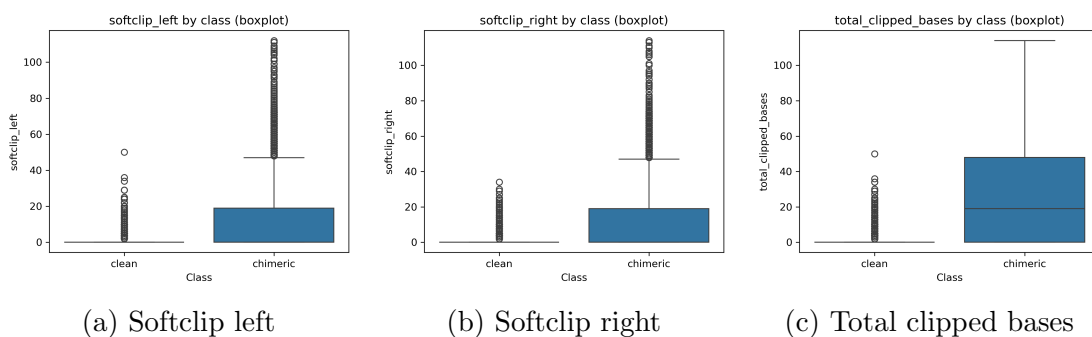


Figure B.3: Boxplots of clipping-based features by class.

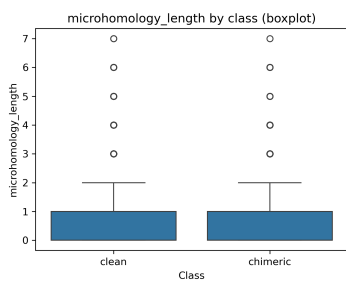
## 1063 B.0.3 K-mer Features



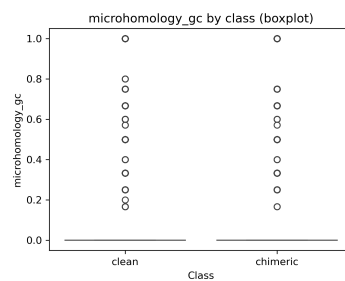
Figure B.4: Boxplots of k-mer features by class.



## 1064 B.0.4 Microhomology Features



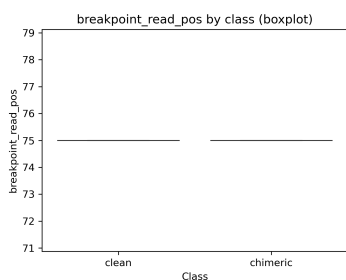
(a) Microhomology length



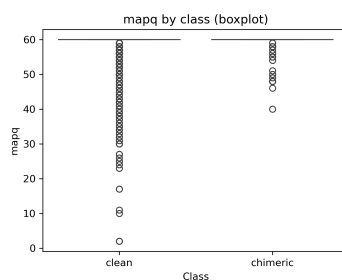
(b) Microhomology GC

Figure B.5: Boxplots of microhomology features by class.

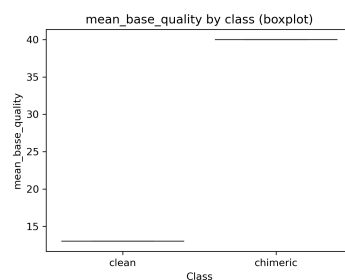
## 1065 B.0.5 Others



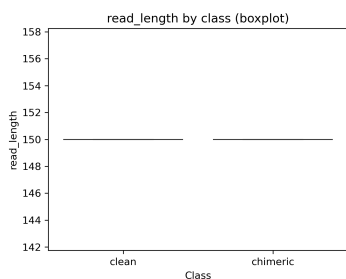
(a) Breakpoint read position



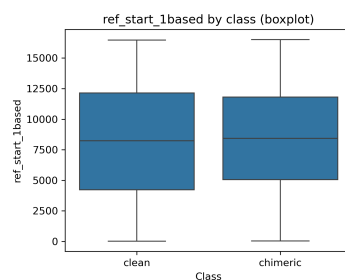
(b) MAPQ



(c) Mean base quality



(d) Read length



(e) Reference start (1-based)

Figure B.6: Boxplots of other numeric features by class.

# References

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

1085 45(4), e18. doi: 10.1093/nar/gkw955

1086 Edgar, R. C. (n.d.). *Uchime in practice*. Retrieved from [https://www.drive5](https://www.drive5.com/usearch/manual7/uchime_practical.html)

1087 [.com/usearch/manual7/uchime\\_practical.html](https://www.drive5.com/usearch/manual7/uchime_practical.html)

1088 Edgar, R. C. (2016). Uchime2: improved chimera prediction for amplicon se-

1089 quencing. *bioRxiv*. Retrieved from [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:88955007)

1090 [CorpusID:88955007](https://api.semanticscholar.org/CorpusID:88955007)

1091 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011).

1092 Uchime improves sensitivity and speed of chimera detection. *Bioinformatics*,

1093 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

1094 Glenn, T. C. (2011). Field guide to next-generation dna sequencers. *Molecular*

1095 *Ecology Resources*, 11(5), 759–769. doi: 10.1111/j.1755-0998.2011.03024.x

1096 Gonzalez, J. M., Zimmermann, J., & Saiz-Jimenez, C. (2004, 09). Evalu-

1097 ating putative chimeric sequences from pcr-amplified products. *Bioin-*

1098 *formatics*, 21(3), 333–337. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bti008)

1099 [bioinformatics/bti008](https://doi.org/10.1093/bioinformatics/bti008) doi: 10.1093/bioinformatics/bti008

1100 Gray, M. W. (2012). Mitochondrial evolution. *Cold Spring Harbor perspectives*

1101 *in biology*, 4. Retrieved from [https://doi.org/10.1101/cshperspect](https://doi.org/10.1101/cshperspect.a011403)

1102 [.a011403](https://doi.org/10.1101/cshperspect.a011403) doi: 10.1101/cshperspect.a011403

1103 Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial

1104 genomes directly from genomic next-generation sequencing reads—a baiting

1105 and iterative mapping approach. *Nucleic Acids Research*, 41(13), e129. doi:

1106 10.1093/nar/gkt371

1107 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,

1108 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de

1109 novo assembly of organelle genomes. *Genome Biology*, 21(1), 241. doi:

1110 10.1186/s13059-020-02154-5

- 1111 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-  
 1112 pression of pcr-mediated recombination. *Nucleic Acids Research*, 26(7),  
 1113 1819–1825. doi: 10.1093/nar/26.7.1819
- 1114 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.  
 1115 (2021). Mitochondrial dna reveals genetically structured haplogroups of  
 1116 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*  
 1117 *Marine Science*, 41, 101588. doi: 10.1016/j.rsma.2020.101588
- 1118 Li, H. (2018, 05). Minimap2: pairwise alignment for nucleotide sequences. *Bioin-*  
 1119 *formatics*, 34(18), 3094–3100. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bty191)  
 1120 [bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191) doi: 10.1093/bioinformatics/bty191
- 1121 Liang, Q., Bible, P. W., Liu, Y., Zou, B., & Wei, L. (2020, 02). Deepmi-  
 1122 crobes: taxonomic classification for metagenomics with deep learning. *NAR*  
 1123 *Genomics and Bioinformatics*, 2(1), lqaa009. Retrieved from [https://](https://doi.org/10.1093/nargab/lqaa009)  
 1124 [doi.org/10.1093/nargab/lqaa009](https://doi.org/10.1093/nargab/lqaa009) doi: 10.1093/nargab/lqaa009
- 1125 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*  
 1126 *Reviews Genetics*, 11(1), 31–46. doi: 10.1038/nrg2626
- 1127 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,  
 1128 an ensemble classifier for chimera detection in 16s rna sequencing stud-  
 1129 ies. *Applied and Environmental Microbiology*, 81(5), 1573–1584. Retrieved  
 1130 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:  
 1131 10.1128/AEM.02896-14
- 1132 Peccoud, J., Lequime, S., Moltini-Conclois, I., Giraud, I., Lambrechts, L., &  
 1133 Gilbert, C. (2018, 04). A survey of virus recombination uncovers canon-  
 1134 ical features of artificial chimeras generated during deep sequencing li-  
 1135 brary preparation. *G3 Genes—Genomes—Genetics*, 8(4), 1129–1138. Re-  
 1136 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*

1163 *Sciences*, 40(11), 701-714. Retrieved from <https://www.sciencedirect>  
 1164 [.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)  
 1165 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)  
 1166 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,  
 1167 11). Large-scale machine learning for metagenomics sequence classifica-  
 1168 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)  
 1169 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683  
 1170 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*  
 1171 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI  
 1172 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)  
 1173 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)