

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 January 3, 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 lemuru* 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

Contents

44	1 Introduction	1
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	2 Review of Related Literature	7
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	3	Research Methodology	19
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	4	Results and Discussion	34
77	4.1	Descriptive Analysis of Features	35
78	4.1.1	Summary Statistics Per Class	35
79	4.1.2	Boxplots By Class	37
80	4.1.3	Correlation Analysis of Extracted Features	39
81	4.2	Baseline Classification Performance	40
82	4.3	Effect of Hyperparameter Tuning	42
83	4.4	Detailed Evaluation of Representative Models	43
84	4.4.1	Confusion Matrices and Error Patterns	44
85	4.4.2	ROC and Precision–Recall Curves	45
86	4.5	Feature Importance	46
87	4.5.1	Permutation Importance of Individual Features	46
88	4.5.2	Feature Family Importance	48
89	4.6	Feature Selection	50

90	4.6.1	Cumulative Importance Curve	51
91	4.6.2	Performance Comparison Across Feature Sets	51
92	4.6.3	Interpretation and Final Feature Set Choice	53
93	4.7	Summary of Findings	53
94	A	Complete Per-Class Summary Statistics	55
95	B	Boxplots for All Numeric Features by Feature Family	59
96	B.0.1	SA Structure (Supplementary Alignment and Segment Met-	
97		rics)	59
98	B.0.2	Clipping-Based Features	61
99	B.0.3	K-mer Features	61
100	B.0.4	Microhomology Features	62
101	B.0.5	Others	62

List of Figures

103	3.1	Process diagram of the study workflow.	20
104	4.1	Boxplots of key features by class	38
105	4.2	Feature correlation heatmap showing relationships among alignment-	
106		derived and sequence-derived features.	40
107	4.3	Test F1 of all baseline classifiers, showing that no single model	
108		clearly dominates and several achieve comparable performance. . .	41
109	4.4	Comparison of test F1 (left) and ROC-AUC (right) for baseline	
110		and tuned models.	43
111	4.5	Confusion matrices for the four representative models on the held-	
112		out test set.	45
113	4.6	ROC (left) and precision-recall (right) curves for the four represen-	
114		tative models on the held-out test set.	46
115	4.7	Permutation-based feature importance for four representative clas-	
116		sifiers.	48

117	4.8	Aggregated feature family importance across four models.	50
118	4.9	Cumulative importance curve of features sorted by importance. . .	51
119	4.10	Comparison of F1 and ROC–AUC for the full, top-4 selected, and	
120		no-microhomology feature set variants.	52
121	B.1	Boxplots of SA Structure features by class (1/2).	59
122	B.2	Boxplots of SA Structure features by class (2/2).	60
123	B.3	Boxplots of clipping-based features by class.	61
124	B.4	Boxplots of k-mer features by class.	61
125	B.5	Boxplots of microhomology features by class.	62
126	B.6	Boxplots of other numeric features by class.	62

127 List of Tables

128	2.1	Comparison of Chimera Detection Approaches and Tools	17
129	3.1	Timetable of activities.	33
130	4.1	Summary statistics of selected key features by class.	37
131	4.2	Performance of baseline classifiers on the held-out test set.	41
132	4.3	Performance of tuned classifiers on the held-out test set.	42
133	4.4	Test set performance of three feature set variants using tuned Cat-	
134		Boost.	52
135	A.1	Complete per-class summary statistics for all extracted features. .	56

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

175 solution for improving mitochondrial genome reconstruction.

176 1.2 Problem Statement

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

190 1.3 Research Objectives

191 1.3.1 General Objective

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

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

196 **1.3.2 Specific Objectives**

197 Specifically, the study aims to:

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

209 **1.4 Scope and Limitations of the Research**

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

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

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

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

238 patterns.

239 1.5 Significance of the Research

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

Chapter 2

Review of Related Literature

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

2.1 The Mitochondrial Genome

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

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

269 **2.1.1 Mitochondrial Genome Assembly**

270 Mitochondrial genome assembly refers to the reconstruction of the complete mito-
271 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is
272 conducted to obtain high-quality, continuous representations of the mitochondrial
273 genome that can be used for a wide range of analyses, including species identi-
274 fication, phylogenetic reconstruction, evolutionary studies, and investigations of
275 mitochondrial diseases. Because mtDNA evolves rapidly, its assembled sequence
276 provides valuable insights into population structure, lineage divergence, and adap-
277 tive evolution across taxa (Boore, 1999). Compared to nuclear genome assembly,
278 assembling the mitochondrial genome is often considered more straightforward but
279 still encounters technical challenges such as the formation of chimeric reads. Com-
280 monly used tools for mitogenome assembly such as GetOrganelle and MITObim
281 operate under the assumption of organelle genome circularity, and are vulnerable
282 when chimeric reads disrupt this circular structure, resulting in assembly errors
283 (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

sequence correspond to distinct high-abundance sequences.

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

2.3.1 UCHIME

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

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

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

356 2.3.2 UCHIME2

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

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

373 2.3.3 CATCh

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

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

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

396 2.3.4 ChimPipe

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

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

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

414 **2.4 Machine Learning Approaches for Chimera** 415 **and Sequence Quality Detection**

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

427 **2.4.1 Feature-Based Representations of Genomic Se-** 428 **quences**

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

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

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

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

453 **2.5 Synthesis of Chimera Detection Approaches**

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

Table 2.1: Comparison of Chimera Detection Approaches and Tools

Method / Tool	Core Approach	Key Limitations
Reference-based Detection	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.
De novo Detection	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.
UCHIME	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.
UCHIME2	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.
CATCh	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.
ChimPipe	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.

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

469 Chapter 3

470 Research Methodology

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

476 3.1 Research Activities

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

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

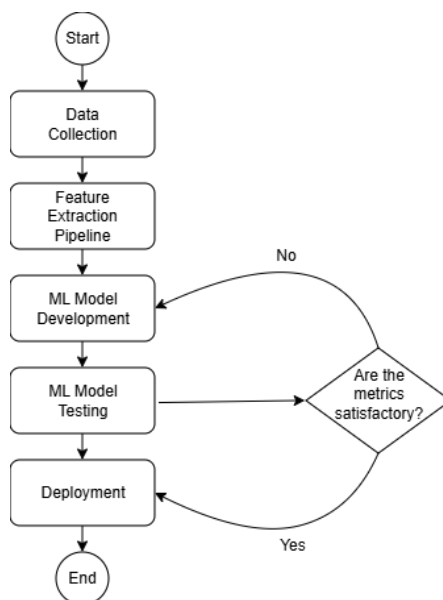


Figure 3.1: Process diagram of the study workflow.

489 3.1.1 Data Collection

490 The mitochondrial genome reference sequence of *S. lemur* was obtained from the
491 NCBI database (accession number NC_039553.1) in FASTA format and was used
492 to generate simulated reads.

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

495 proximately one (1) week.

496 Data Preprocessing

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

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

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

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

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

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

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

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

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

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

```
529 samtools view -bS clean.sam -o clean.bam
```

```
530 samtools view -bS chimeric.sam -o chimeric.bam
```

531

```
532 samtools sort clean.bam -o clean.sorted.bam
```

```
533 samtools index clean.sorted.bam
```

534

```
535 samtools sort chimeric.bam -o chimeric.sorted.bam
```

```
536 samtools index chimeric.sorted.bam
```

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

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

547 **3.1.2 Feature Extraction Pipeline**

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

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

559 total duration of approximately two (2) weeks.

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

565 **Supplementary Alignment and Alignment-Structure Features**

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

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

579 **K-mer Composition Difference**

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

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

592 **Microhomology**

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

601 without including distant sequences that are unlikely to be mechanistically
602 related.

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

611 **3.1.3 Machine Learning Model Development**

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

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

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

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

3.1.4 Model Benchmarking, Hyperparameter Optimization, and Evaluation

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

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

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

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

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

672 **3.1.5 Feature Importance, Feature Selection, and Inter-** 673 **pretation**

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

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

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

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

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

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

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

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

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

729 3.1.6 Validation and Testing

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

736 Comparative evaluation was performed across all candidate algorithms and
737 CatBoost feature-set variants to determine which models demonstrated the high-
738 est predictive performance and computational efficiency under identical data con-
739 ditions. Their metrics were compared to identify which algorithms and feature

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

742 **3.1.7 Documentation**

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

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

754 For manuscript preparation and supplementary materials, Overleaf (L^AT_EX)
755 was used to produce publication-quality formatting and consistent referencing.

756 **3.2 Calendar of Activities**

757 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
758 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	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •

759 Chapter 4

760 Results and Discussion

761 This chapter presents the performance of the proposed feature set and machine
762 learning models for detecting PCR-induced chimeric reads in simulated mito-
763 chondrial Illumina data. The behaviour of the extracted features is first examined
764 through descriptive and correlation analyses, followed by a comparison of baseline
765 and tuned classifiers. The chapter then examines model performance in detail and
766 investigates the contribution of individual features and feature families, including
767 the impact of feature selection on classification performance.

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

771 4.1 Descriptive Analysis of Features

772 4.1.1 Summary Statistics Per Class

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

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

782 Alignment and Supplementary Alignment Features

783 Features related to supplementary alignments show strong separation between
784 classes. Chimeric reads exhibit supplementary alignments, as reflected by higher
785 values of `has_sa`, `sa_count`, and `num_segments`, whereas clean reads consistently
786 show a single alignment segment with no supplementary mappings. This behavior
787 is consistent with the expected structure of chimeric reads and indicates that
788 alignment-based features are highly informative.

789 Clipping-Based Features

790 Clipping-related features, including `softclip_left`, `softclip_right`, and
791 `total_clipped_bases`, display higher means and broader distributions in chimeric
792 reads. Clean reads are dominated by zero or near-zero clipping, while chimeric
793 reads exhibit increased clipping and greater variability, which reflects the presence
794 of split alignments.

795 K-mer Distribution Features

796 K-mer-based features, such as `kmer_js_divergence` and `kmer_cosine_diff`, show
797 only modest differences between clean and chimeric reads. Chimeric reads show
798 slightly higher average divergence, but substantial overlap with clean reads means
799 this feature alone cannot reliably distinguish the classes.

800 Microhomology Features

801 Microhomology-related features (`microhomology_length` and `microhomology_gc`)
802 exhibit nearly identical summary statistics across both classes. The majority of
803 reads in both classes contain short or zero-length microhomologies, resulting in
804 minimal separation. This means that microhomology serves as a weak standalone
805 indicator and is more appropriately treated as supporting evidence.

806 Overall, the summary statistics indicate that alignment-based and clipping-
807 based features provide the strongest class separation, k-mer features contribute
808 limited but complementary signal, and microhomology features exhibit minimal

discriminative power on their own. These observations motivate the combined 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

4.1.2 Boxplots By Class

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, chimeric reads exhibit higher medians and longer upper whiskers than clean reads, indicating increased variability and the presence

818 of split alignments.

819 Supplementary alignment features show that clean reads are largely zero,
820 whereas chimeric reads display a wider distribution, reflecting frequent supple-
821 mentary alignments.

822 K-mer metrics show a slight upward shift for chimeric reads, but substantial
823 overlap with clean reads indicates modest discriminative power.

824 Microhomology features have nearly overlapping distributions for both classes,
825 consistent with their low standalone predictive importance.

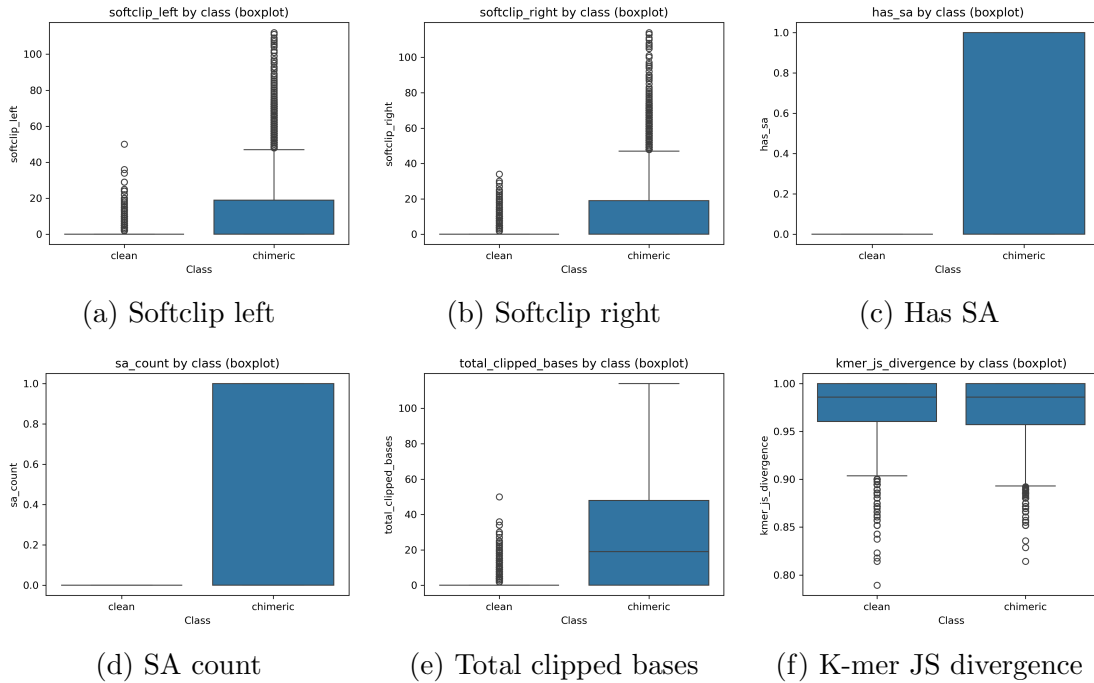


Figure 4.1: Boxplots of key features by class

826 4.1.3 Correlation Analysis of Extracted Features

827 A feature correlation heatmap (Figure 4.2) was generated to examine relation-
828 ships among the extracted variables and to identify patterns of redundancy and
829 independence within the feature set. The analysis shows that alignment- and
830 clipping-related features form a strongly correlated cluster, including indicators
831 of supplementary alignments, alignment segment counts, positional differences,
832 and soft-clipping measures. These features capture related aspects of alignment
833 fragmentation, which is a known characteristic of chimeric reads, and several show
834 moderate correlations with the class label, supporting their relevance for distin-
835 guishing chimeric from clean reads. In contrast, general read- and alignment-
836 quality metrics, such as read length, base quality, and mapping quality, exhibit
837 weak correlations with most split-alignment features, indicating that they pro-
838 vide distinct information rather than overlapping with alignment-derived signals.
839 Sequence-based features display a similar pattern of independence, as k-mer di-
840 vergence metrics show weak correlations with other feature groups, while micro-
841 homology features exhibit generally low correlations with both alignment- and
842 k-mer-based features. Overall, the correlation structure highlights intentional re-
843 dundancy within alignment-derived features and clear separation between feature
844 families, supporting the use of features that capture different aspects of chimeric
845 read characteristics to improve chimera classification.

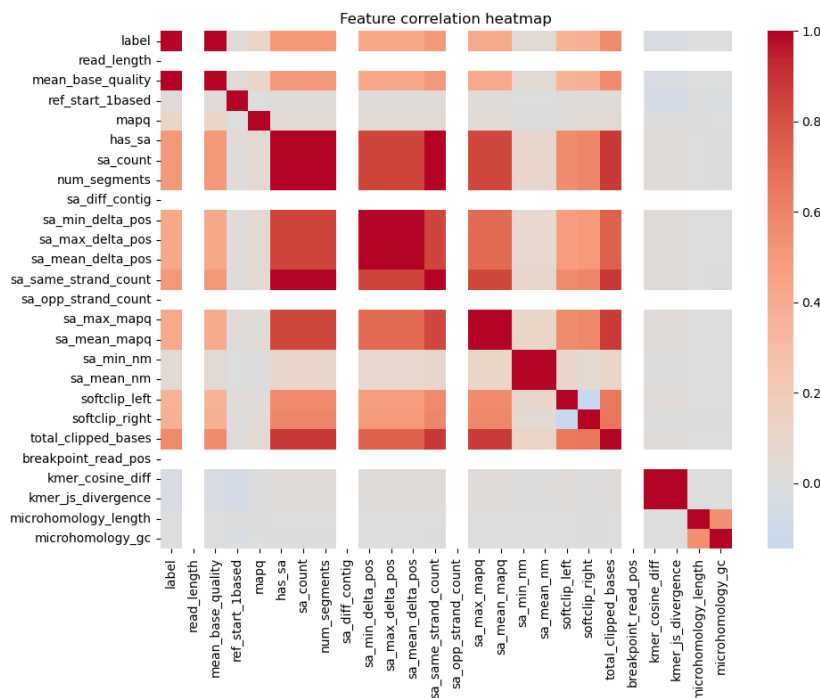


Figure 4.2: Feature correlation heatmap showing relationships among alignment-derived and sequence-derived features.

4.2 Baseline Classification Performance

Table 4.2 summarises the performance of eleven classifiers trained on the engineered feature set using five-fold cross-validation and evaluated on the held-out test set. All models were optimised using default hyperparameters, without dedicated tuning.

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

855 Across other models, test F1-scores clustered in a narrow band between ap-
856 proximately 0.74 and 0.77 and ROC-AUC values between 0.82 and 0.84. Gradi-
857 ent boosting, CatBoost, LightGBM, XGBoost, bagging trees, random forest, and
858 multilayer perceptron (MLP) all produced very similar scores, with CatBoost and
859 gradient boosting slightly ahead (test F1 ≈ 0.77 , ROC-AUC ≈ 0.84). Linear
860 models (logistic regression and calibrated linear SVM) performed only marginally
861 worse (test F1 ≈ 0.74), while Gaussian Naive Bayes lagged behind with substan-
862 tially lower F1 (≈ 0.65) despite very 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	0.500000	0.500000	1.000000	0.667000	0.500000
logreg_l2	0.789000	0.945000	0.614000	0.744000	0.821000
linear_svm_calibrated	0.789000	0.945000	0.614000	0.744000	0.820000
random_forest	0.788000	0.894000	0.654000	0.755000	0.834000
extra_trees	0.788000	0.901000	0.647000	0.753000	0.824000
gradient_boosting	0.802000	0.936000	0.648000	0.766000	0.840000
xgboost	0.800000	0.929000	0.650000	0.765000	0.839000
lightgbm	0.799000	0.926000	0.650000	0.764000	0.838000
catboost	0.803000	0.936000	0.650000	0.767000	0.839000
knn	0.782000	0.892000	0.642000	0.747000	0.815000
gaussian_nb	0.741000	0.996000	0.483000	0.651000	0.819000
bagging_trees	0.792000	0.900000	0.657000	0.760000	0.837000
mlp	0.789000	0.931000	0.625000	0.748000	0.819000

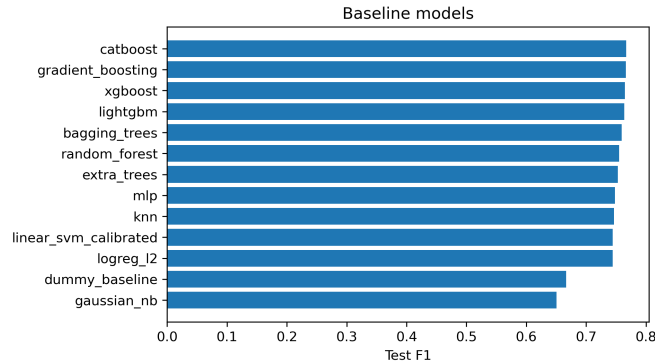


Figure 4.3: Test F1 of all baseline classifiers, showing that no single model clearly dominates and several achieve comparable performance.

4.3 Effect of Hyperparameter Tuning

To assess whether performance could be improved further, ten model families underwent randomised hyperparameter search. 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, XGBoost, random forest, bagging trees, and MLP all experienced small increases in test F1 (typically $\Delta\text{F1} \approx 0.002$ – 0.009) and ROC–AUC (up to $\Delta\text{AUC} \approx 0.008$). After tuning, CatBoost remained the best performer with test accuracy 0.80, precision 0.92, recall 0.66, F1-score 0.77, and ROC–AUC 0.84. Gradient boosting achieved almost identical performance (F1 0.77, AUC 0.84). Random forest and bagging trees also improved to F1 scores around 0.76 with $\text{AUC} \approx 0.84$.

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

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

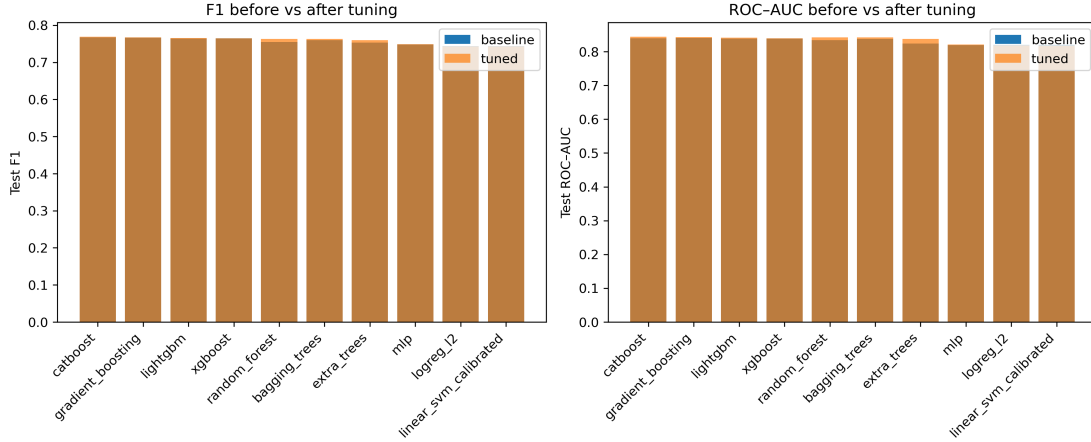


Figure 4.4: Comparison of test F1 (left) and ROC-AUC (right) for baseline and tuned models.

876 Because improvements are small and within cross-validation variability, tun-
877 ing was interpreted as stabilising and slightly refining the models rather than
878 completely altering their behaviour or their relative ranking.

879 4.4 Detailed Evaluation of Representative Mod- 880 els

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

887 4.4.1 Confusion Matrices and Error Patterns

888 Classification reports and confusion matrices for the four models reveal consistent
889 patterns. CatBoost and gradient boosting both reached overall accuracy of ap-
890 proximately 0.80 with similar macro-averaged F1 scores (~ 0.80). For CatBoost,
891 precision and recall for clean reads were 0.73 and 0.95, respectively, while for
892 chimeric reads they were 0.92 and 0.66 ($F1 = 0.77$). Gradient boosting showed
893 nearly identical trade-offs.

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

898 Across all models, errors were asymmetric. False negatives (chimeric reads pre-
899 dicted as clean) were more frequent than false positives. For example, CatBoost
900 misclassified 1,369 chimeric reads as clean but only 215 clean reads as chimeric.
901 This pattern indicates that the models are conservative and prioritise avoiding
902 false chimera calls at the expense of missing some true chimeras. Consultation
903 with PGC Visayas indicated that this conservative behavior is generally accept-
904 able, though further evaluation and testing will be required to assess its impact
905 on downstream analyses.

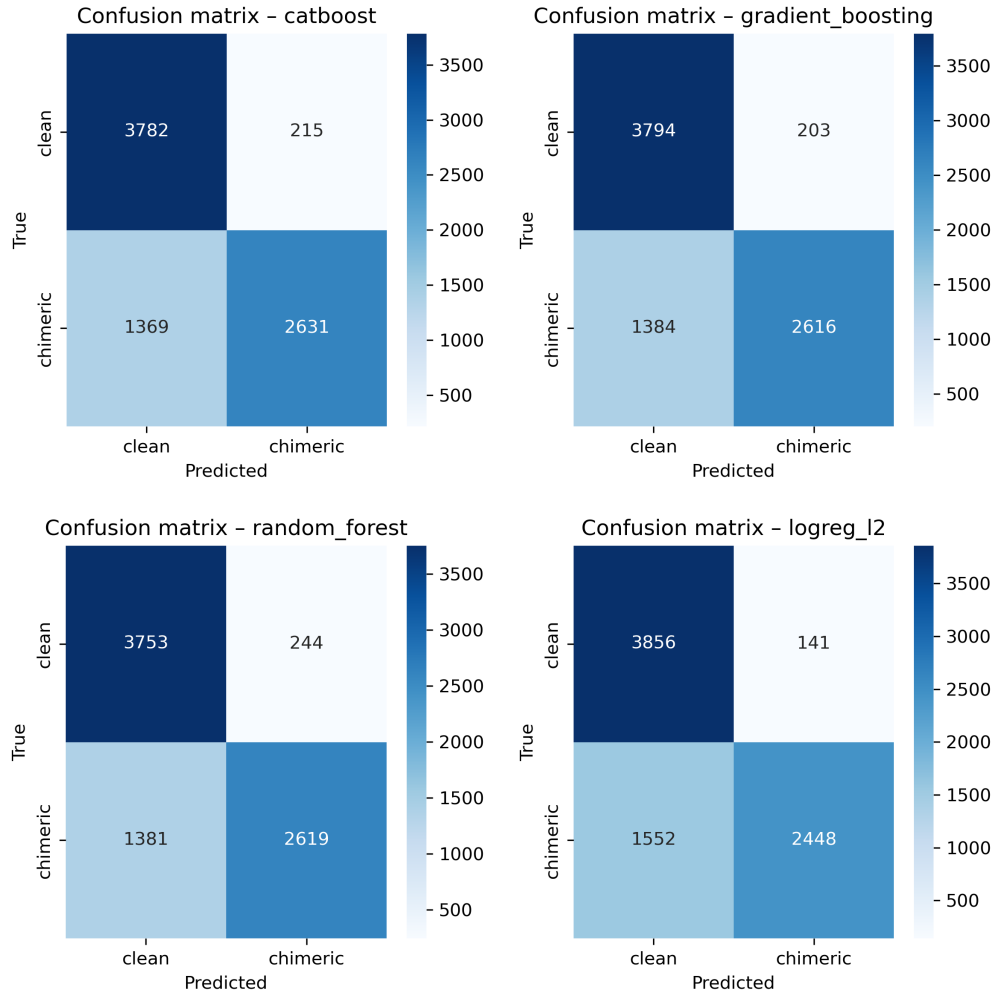


Figure 4.5: Confusion matrices for the four representative models on the held-out test set.

4.4.2 ROC and Precision–Recall Curves

Receiver operating characteristic (ROC) and precision–recall (PR) curves as shown in Figure 4.6 further support the similarity among the top models. The three tree-based ensembles (CatBoost, gradient boosting, random forest) achieved ROC–AUC values of approximately 0.84 and average precision (AP) around 0.88.

911 Logistic regression performed slightly worse ($AUC \approx 0.82$, $AP \approx 0.87$) but still
 912 substantially better than the dummy baseline.

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

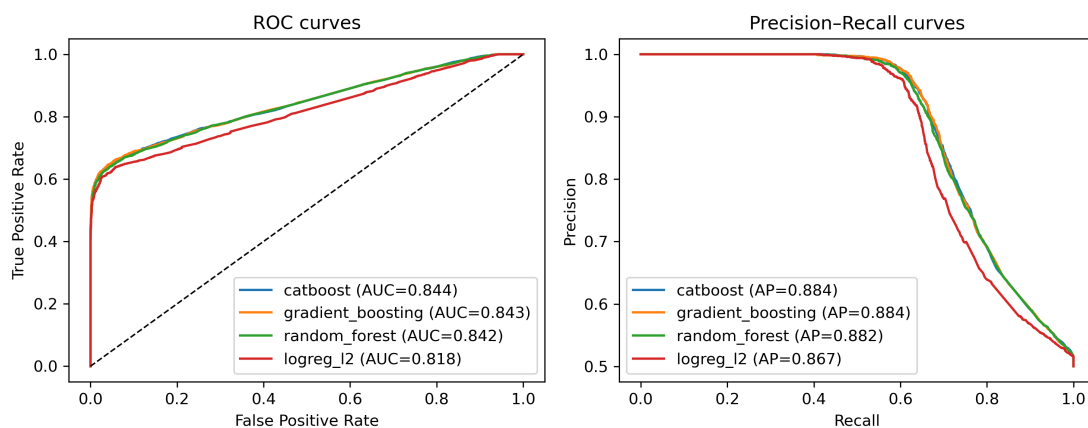


Figure 4.6: ROC (left) and precision–recall (right) curves for the four representative models on the held-out test set.

918 4.5 Feature Importance

919 4.5.1 Permutation Importance of Individual Features

920 To understand how each classifier made predictions, feature importance was quan-
 921 tified using permutation importance. This analysis was applied to four represen-
 922 tative models: CatBoost, Gradient Boosting, Random Forest, and L_2 -regularized

923 Logistic Regression.

924 As shown in Figure 4.7, the total number of clipped bases consistently pro-
925 vides a strong predictive signal, particularly in Random Forest, Gradient Boosting,
926 and L₂-regularized Logistic Regression. CatBoost differs by assigning the highest
927 importance to k-mer divergence metrics such as `kmer_js_divergence`, which cap-
928 ture subtle sequence changes resulting from structural variants or PCR-induced
929 chimeras. Soft-clipping features (`softclip_left` and `softclip_right`) provide
930 more information around breakpoints, complementing these primary signals in all
931 models except Gradient Boosting. L₂-regularized Logistic Regression relies more
932 on alignment-based split-read metrics.

933 Overall, these results indicate that accurate detection of chimeric reads relies
934 on both alignment-based signals and k-mer compositional information. Explicit
935 microhomology features contribute minimally in this analysis, and combining both
936 alignment-based and sequence-level features enhances model sensitivity and speci-
937 ficity.

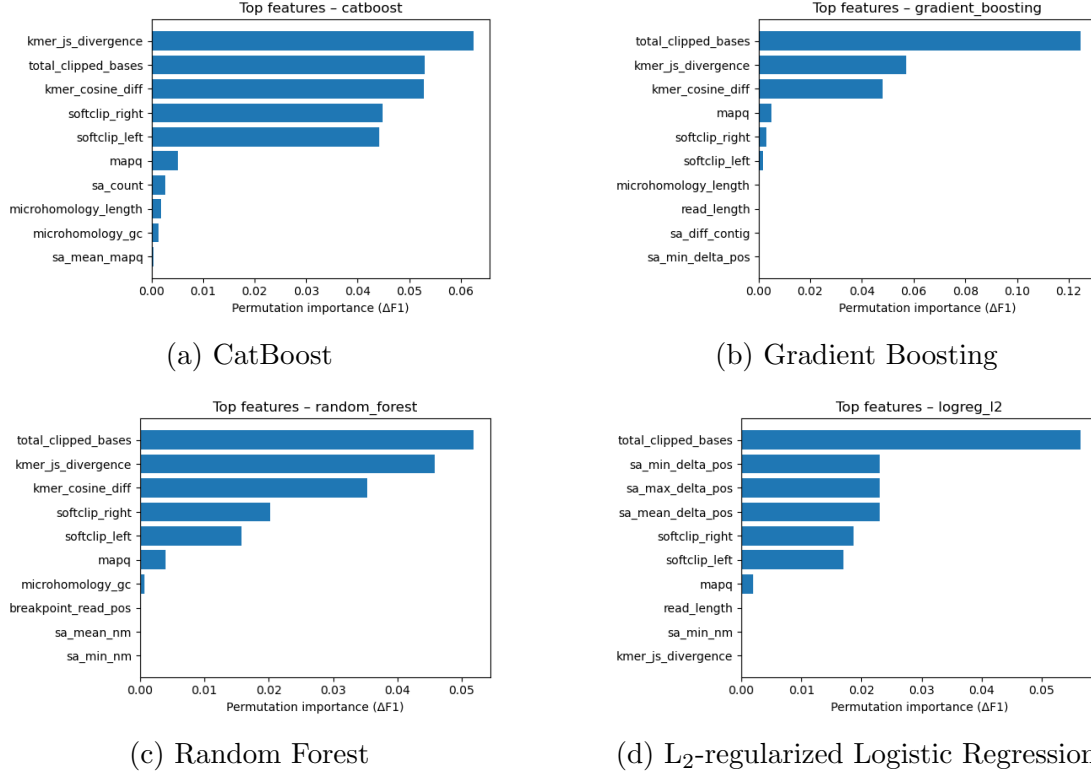


Figure 4.7: Permutation-based feature importance for four representative classifiers.

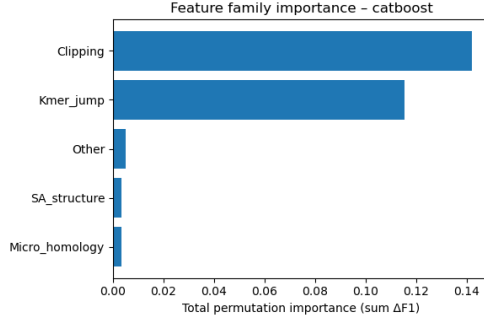
4.5.2 Feature Family Importance

To evaluate the contribution of broader signals, features were grouped into five families: SA-structure (supplementary alignment and segment metrics, e.g., `has_sa`, `sa_count`, `sa_min_delta_pos`, `sa_mean_nm`, etc.), Clipping (`softclip_left`, `softclip_right`, `total_clipped_bases`, `breakpoint_read_pos`), Kmer-jump (`kmer_cosine_diff`, `kmer_js_divergence`), Micro-homology (`microhomology_length`, `microhomology_gc`), and Other (e.g., `mapq`).

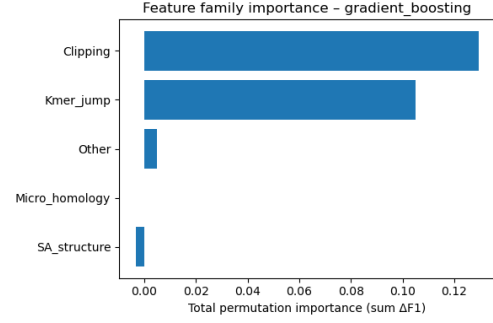
Aggregated analyses reveal consistent patterns across models. In CatBoost, the Clipping family has the largest cumulative contribution (0.14), followed

947 by Kmer_jump (0.12), with Other features contributing minimally (0.005) and
948 SA_structure (0.003) and Micro_homology (0.003) providing minimal predictive
949 power. Gradient Boosting shows a similar trend, with Clipping (0.13) domi-
950 nating, Kmer_jump (0.11) secondary, and the remaining families contributing
951 negligibly. Random Forest integrates both Clipping (0.088) and Kmer_jump
952 (0.08) effectively, while SA_structure, Micro_homology, and Other remain minor
953 contributors. L₂-regularized Logistic Regression emphasizes Clipping (0.09)
954 and SA_structure (0.07), with Kmer_jump and Micro_homology having minimal
955 impact.

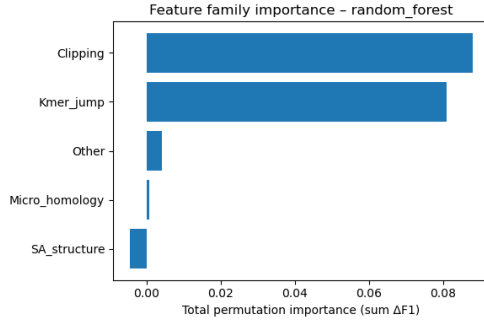
956 Both feature-level and aggregated analyses indicate that detection of chimeric
957 reads in this dataset relies primarily on alignment irregularities (Clipping) and
958 k-mer compositional shifts (Kmer_jump), which often arise from PCR-induced
959 template switching events, while explicit microhomology features contribute min-
960 imally.



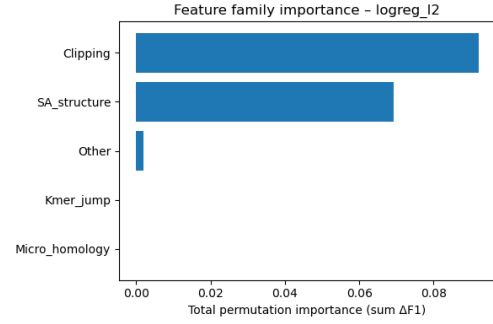
(a) CatBoost



(b) Gradient Boosting



(c) Random Forest



(d) L_2 -regularized Logistic Regression

Figure 4.8: Aggregated feature family importance across four models.

961 4.6 Feature Selection

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

967 4.6.1 Cumulative Importance Curve

968 The cumulative importance curve was computed using the tuned CatBoost clas-
969 sifier. Figure 4.9 illustrates the contribution of features sorted by importance.
970 The curve rises steeply for the first few features and then gradually plateaus,
971 indicating that a small number of features capture most of the model’s pre-
972 dictive power. A cumulative importance of 95% is reached at $k = 4$ features,
973 which are `total_clipped_bases`, `kmer_js_divergence`, `kmer_cosine_diff`, and
974 `softclip_left`.

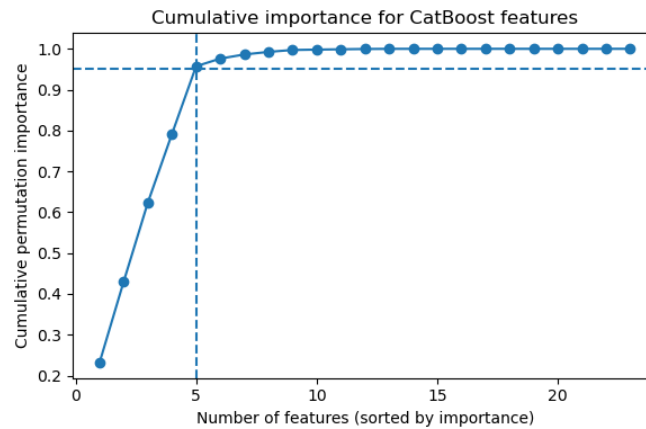


Figure 4.9: Cumulative importance curve of features sorted by importance.

975 4.6.2 Performance Comparison Across Feature Sets

976 Classification performance was compared across three feature sets using a tuned
977 CatBoost classifier. The full model, incorporating all 23 engineered features,
978 achieved an F1 score of approximately 0.7686 and a ROC-AUC of 0.8436.
979 A reduced model using only the top four features (`total_clipped_bases`,
980 `kmer_js_divergence`, `kmer_cosine_diff`, and `softclip_left`) achieved nearly

981 equivalent performance with an F1 of 0.7670 and a ROC-AUC of 0.8353. An
 982 ablation model excluding microhomology features (`microhomology_length` and
 983 `microhomology_gc`) also performed comparably, with an F1 of 0.7679 and ROC-
 984 AUC of 0.8447. These results indicate that clipping and k-mer features capture
 985 almost all predictive signal, while microhomology features are largely redundant
 986 in this dataset.

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

Variant	No. of Features	Test F1	ROC-AUC
Full CatBoost	23	0.7686	0.8436
Selected (top-4)	4	0.7670	0.8353
No microhomology	21	0.7679	0.8447

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

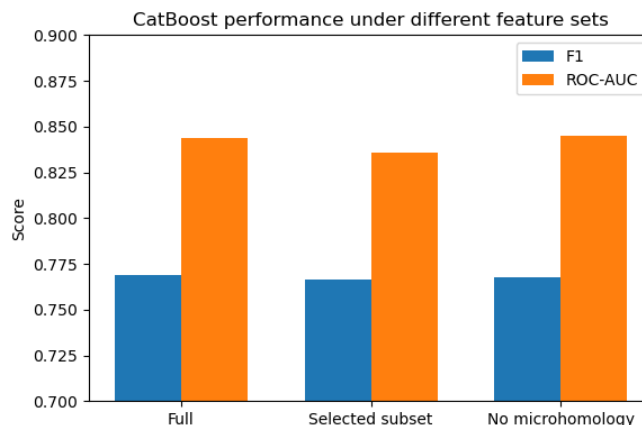


Figure 4.10: Comparison of F1 and ROC-AUC for the full, top-4 selected, and no-microhomology feature set variants.

990 4.6.3 Interpretation and Final Feature Set Choice

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

998 4.7 Summary of Findings

999 All evaluated machine learning models substantially outperformed the dummy
1000 baseline, demonstrating that the engineered feature set contains meaningful
1001 signals for detecting PCR-induced chimeric reads. Across classifiers, the best-
1002 performing models achieved test F1-scores of approximately 0.77 and ROC-AUC
1003 values around 0.84 on held-out simulated mitochondrial reads, indicating reli-
1004 able discrimination between clean and chimeric sequences. Among the tested
1005 approaches, tree-based ensemble and boosting methods consistently showed the
1006 strongest and most stable performance. In particular, CatBoost and Gradient
1007 Boosting ranked among the top models across multiple evaluation metrics,
1008 both before and after hyperparameter tuning. These results suggest that non-
1009 linear ensemble methods are well suited to capturing the interaction between
1010 alignment-derived and sequence-derived features in this setting.

1011 Analysis of feature behaviour revealed clear differences in how effectively fea-

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

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

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

1030 **Appendix A**

1031 **Complete Per-Class Summary**

1032 **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)

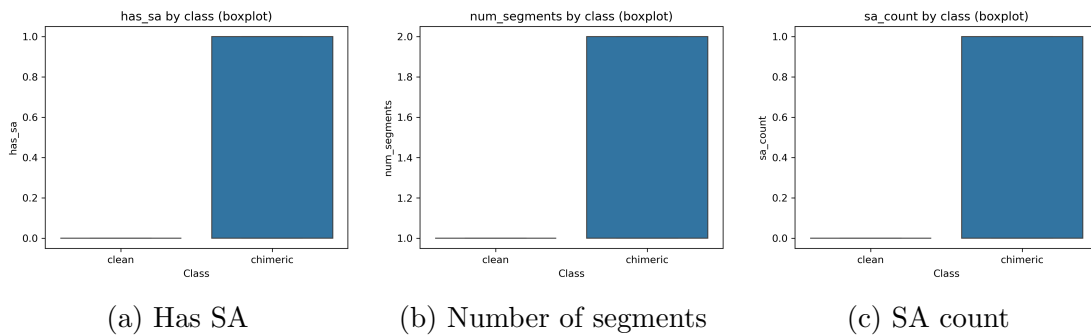
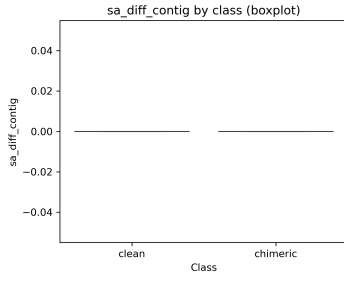
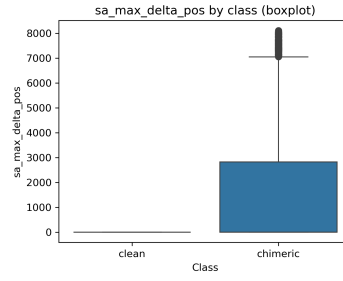


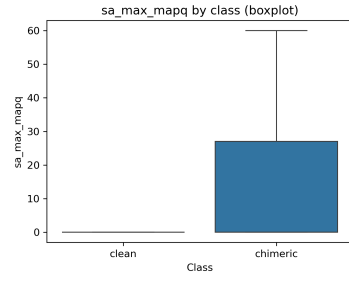
Figure B.1: Boxplots of SA Structure features by class (1/2).



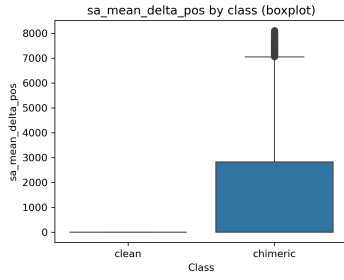
(a) SA different contig



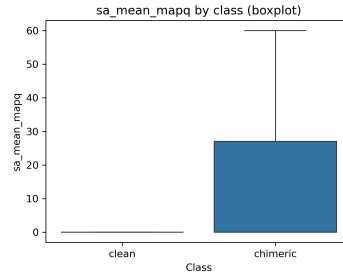
(b) SA max Δ position



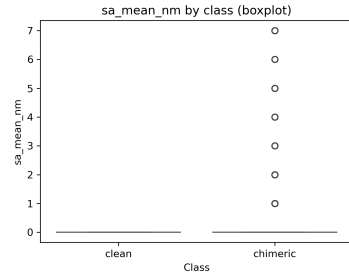
(c) SA max MAPQ



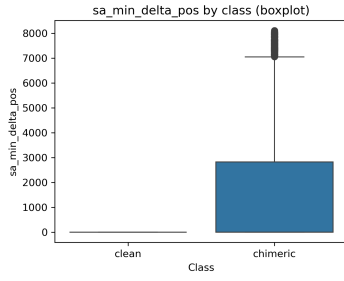
(d) SA mean Δ position



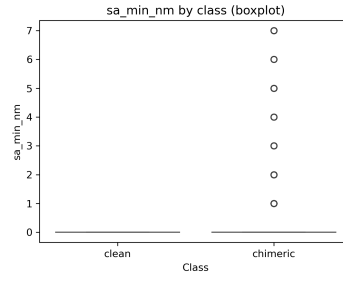
(e) SA mean MAPQ



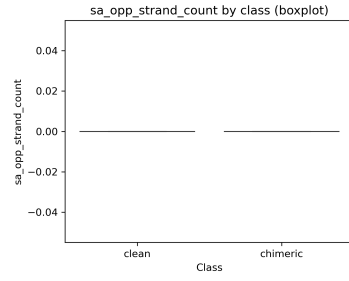
(f) SA mean NM



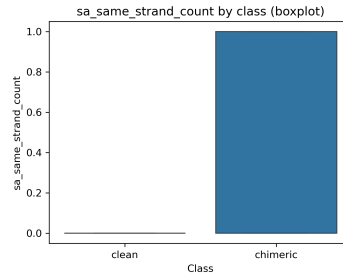
(g) SA min Δ position



(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).

1039 B.0.2 Clipping-Based Features

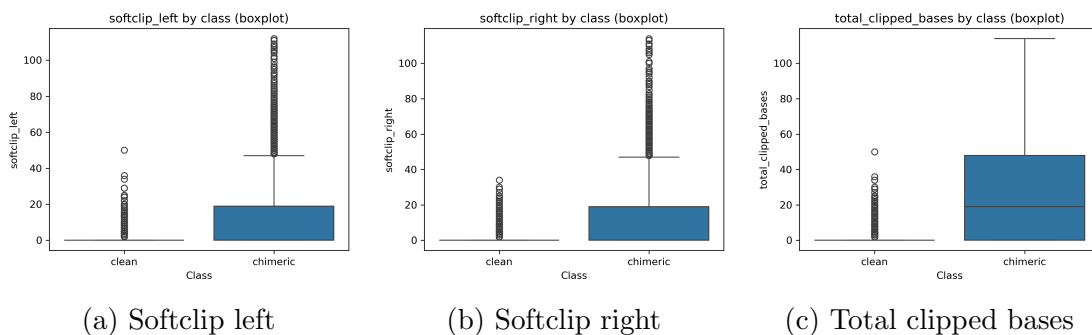


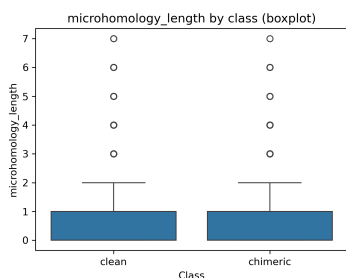
Figure B.3: Boxplots of clipping-based features by class.

1040 B.0.3 K-mer Features

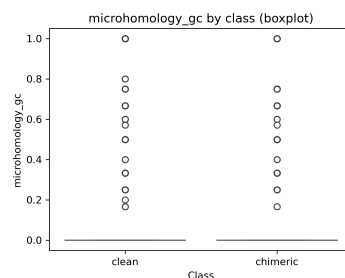


Figure B.4: Boxplots of k-mer features by class.

1041 B.0.4 Microhomology Features



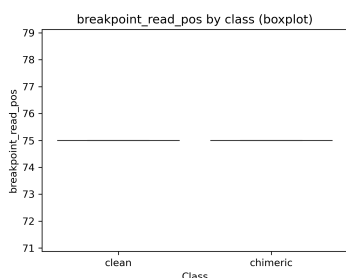
(a) Microhomology length



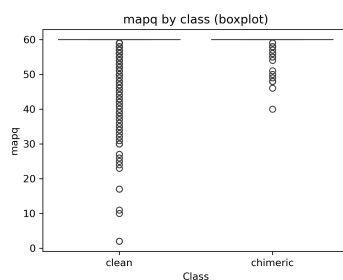
(b) Microhomology GC

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

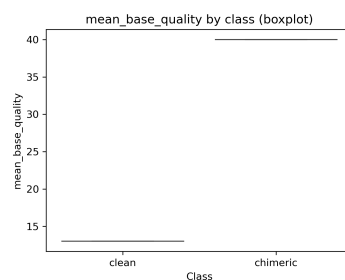
1042 B.0.5 Others



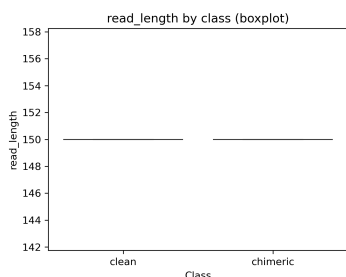
(a) Breakpoint read position



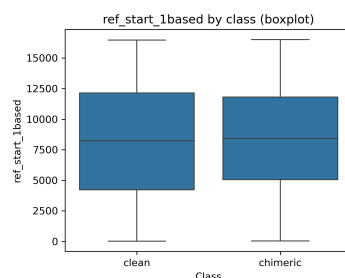
(b) MAPQ



(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*,

1062 45(4), e18. doi: 10.1093/nar/gkw955

1063 Edgar, R. C. (n.d.). *Uchime in practice*. Retrieved from [https://www.drive5](https://www.drive5.com/usearch/manual7/uchime_practical.html)

1064 [.com/usearch/manual7/uchime_practical.html](https://www.drive5.com/usearch/manual7/uchime_practical.html)

1065 Edgar, R. C. (2016). Uchime2: improved chimera prediction for amplicon se-

1066 quencing. *bioRxiv*. Retrieved from [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:88955007)

1067 [CorpusID:88955007](https://api.semanticscholar.org/CorpusID:88955007)

1068 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011).

1069 Uchime improves sensitivity and speed of chimera detection. *Bioinformatics*,

1070 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

1071 Glenn, T. C. (2011). Field guide to next-generation dna sequencers. *Molecular*

1072 *Ecology Resources*, 11(5), 759–769. doi: 10.1111/j.1755-0998.2011.03024.x

1073 Gonzalez, J. M., Zimmermann, J., & Saiz-Jimenez, C. (2004, 09). Evalu-

1074 ating putative chimeric sequences from pcr-amplified products. *Bioin-*

1075 *formatics*, 21(3), 333–337. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bti008)

1076 [bioinformatics/bti008](https://doi.org/10.1093/bioinformatics/bti008) doi: 10.1093/bioinformatics/bti008

1077 Gray, M. W. (2012). Mitochondrial evolution. *Cold Spring Harbor perspectives*

1078 *in biology*, 4. Retrieved from [https://doi.org/10.1101/cshperspect](https://doi.org/10.1101/cshperspect.a011403)

1079 [.a011403](https://doi.org/10.1101/cshperspect.a011403) doi: 10.1101/cshperspect.a011403

1080 Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial

1081 genomes directly from genomic next-generation sequencing reads—a baiting

1082 and iterative mapping approach. *Nucleic Acids Research*, 41(13), e129. doi:

1083 10.1093/nar/gkt371

1084 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,

1085 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de

1086 novo assembly of organelle genomes. *Genome Biology*, 21(1), 241. doi:

1087 10.1186/s13059-020-02154-5

- 1088 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
1089 pression of pcr-mediated recombination. *Nucleic Acids Research*, 26(7),
1090 1819–1825. doi: 10.1093/nar/26.7.1819
- 1091 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
1092 (2021). Mitochondrial dna reveals genetically structured haplogroups of
1093 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
1094 *Marine Science*, 41, 101588. doi: 10.1016/j.rsma.2020.101588
- 1095 Li, H. (2018, 05). Minimap2: pairwise alignment for nucleotide sequences. *Bioin-*
1096 *formatics*, 34(18), 3094–3100. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bty191)
1097 [bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191) doi: 10.1093/bioinformatics/bty191
- 1098 Liang, Q., Bible, P. W., Liu, Y., Zou, B., & Wei, L. (2020, 02). Deepmi-
1099 crobes: taxonomic classification for metagenomics with deep learning. *NAR*
1100 *Genomics and Bioinformatics*, 2(1), lqaa009. Retrieved from [https://](https://doi.org/10.1093/nargab/lqaa009)
1101 doi.org/10.1093/nargab/lqaa009 doi: 10.1093/nargab/lqaa009
- 1102 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*
1103 *Reviews Genetics*, 11(1), 31–46. doi: 10.1038/nrg2626
- 1104 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
1105 an ensemble classifier for chimera detection in 16s rna sequencing stud-
1106 ies. *Applied and Environmental Microbiology*, 81(5), 1573–1584. Retrieved
1107 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
1108 10.1128/AEM.02896-14
- 1109 Peccoud, J., Lequime, S., Moltini-Conclois, I., Giraud, I., Lambrechts, L., &
1110 Gilbert, C. (2018, 04). A survey of virus recombination uncovers canon-
1111 ical features of artificial chimeras generated during deep sequencing li-
1112 brary preparation. *G3 Genes—Genomes—Genetics*, 8(4), 1129–1138. Re-
1113 trieved from <https://doi.org/10.1534/g3.117.300468> doi: 10.1534/

1114 g3.117.300468

1115 Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J.
1116 (2023). Effects of error, chimera, bias, and gc content on the accuracy of
1117 amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from [https://](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)
1118 journals.asm.org/doi/abs/10.1128/msystems.01025-23 doi: 10.1128/
1119 msystems.01025-23

1120 Qiu, X., Wu, L., Huang, H., McDonel, P. E., Palumbo, A. V., Tiedje, J. M., &
1121 Zhou, J. (2001). Evaluation of pcr-generated chimeras, mutations, and het-
1122 eroduplexes with 16s rna gene-based cloning. *Applied and Environmental*
1123 *Microbiology*, 67(2), 880–887. doi: 10.1128/AEM.67.2.880-887.2001

1124 Ren, J., Song, K., Deng, C., Ahlgren, N., Fuhrman, J., Li, Y., ... Sun, F. (2020,
1125 01). Identifying viruses from metagenomic data using deep learning. *Quan-*
1126 *titative Biology*, 8. doi: 10.1007/s40484-019-0187-4

1127 Rodriguez-Martin, B., Palumbo, E., Marco-Sola, S., Griebel, T., Ribeca, P.,
1128 Alonso, G., ... Djebali, S. (2017, 01). Chimpipes: Accurate detection of
1129 fusion genes and transcription-induced chimeras from rna-seq data. *BMC*
1130 *Genomics*, 18. doi: 10.1186/s12864-016-3404-9

1131 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). Vsearch: a
1132 versatile open source tool for metagenomics. *PeerJ*, 4, e2584. doi: 10.7717/
1133 peerj.2584

1134 Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler,
1135 A., & Schatz, M. (2018, 06). Accurate detection of complex structural
1136 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10
1137 .1038/s41592-018-0001-7

1138 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
1139 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*

1140 *Sciences*, 40(11), 701-714. Retrieved from <https://www.sciencedirect>
1141 [.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)
1142 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

1143 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
1144 11). Large-scale machine learning for metagenomics sequence classifica-
1145 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
1146 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683

1147 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*
1148 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI
1149 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)
1150 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)