

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

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

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

13 by

14 Duranne Duran
15 Yvonne Lin
16 Daniella Pailden

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

19 December 5, 2025

Abstract

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

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

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

Contents

44	1 Introduction	1
45	1.1 Overview	1
46	1.2 Problem Statement	3
47	1.3 Research Objectives	4
48	1.3.1 General Objective	4
49	1.3.2 Specific Objectives	4
50	1.4 Scope and Limitations of the Research	5
51	1.5 Significance of the Research	6
52	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	24
69	3.1.3	Machine Learning Model Development	26
70	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
71		Evaluation	28
72	3.1.5	Feature Importance and Interpretation	29

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

89 List of Figures

90	3.1	Process Diagram of Special Project	20
91	4.1	Feature correlation heatmap showing relationships among alignment-	
92		derived and sequence-derived variables.	35
93	4.2	Test F1 of all baseline classifiers, showing that no single model	
94		clearly dominates and several achieve comparable performance. . .	36
95	4.3	Comparison of test F1 (left) and ROC–AUC (right) for baseline and	
96		tuned models. Hyperparameter tuning yields small but consistent	
97		gains, particularly for tree-based ensembles.	38
98	4.4	Confusion matrices for the four representative models on the held-	
99		out test set. All models show more false negatives (chimeric reads	
100		called clean) than false positives.	40
101	4.5	ROC (left) and precision–recall (right) curves for the four represen-	
102		tative models on the held-out test set. Tree-based ensembles cluster	
103		closely, with logistic regression performing slightly but consistently	
104		worse.	41

105	4.6	Permutation-based feature importance for four representative clas-	
106		sifiers. Clipping and k-mer composition features are generally the	
107		strongest predictors, whereas microhomology and other alignment	
108		metrics contribute minimally.	43
109	4.7	Aggregated feature family importance across four models. Clipping	
110		and k-mer compositional shifts are consistently the dominant con-	
111		tributors, while SA_structure, Micro_homology, and other features	
112		contribute minimally.	45
113	A.1	Histogram plots of six key features comparing clean and chimeric	
114		reads.	48

115 List of Tables

<small>116</small>	2.1 Comparison of Chimera Detection Approaches and Tools	17
<small>117</small>	3.1 Timetable of Activities	32
<small>118</small>	4.1 Performance of baseline classifiers on the held-out test set.	36
<small>119</small>	4.2 Performance of tuned classifiers on the held-out test set.	37

Chapter 1

Introduction

1.1 Overview

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

PCR chimeras form when incomplete extension products from one template

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

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

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

161 1.2 Problem Statement

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

179 1.3 Research Objectives

180 1.3.1 General Objective

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

185 1.3.2 Specific Objectives

186 Specifically, the study aims to:

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

1.4 Scope and Limitations of the Research

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

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

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

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

1.5 Significance of the Research

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

240 Chapter 2

241 Review of Related Literature

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

248 2.1 The Mitochondrial Genome

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

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

259 **2.1.1 Mitochondrial Genome Assembly**

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

318 sequence correspond to distinct high-abundance sequences.

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

324 **2.3.1 UCHIME**

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

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

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

346 **2.3.2 UCHIME2**

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

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

363 2.3.3 CATCh

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

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

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

386 2.3.4 ChimPipe

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

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

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

404 **2.4 Machine Learning Approaches for Chimera** 405 **and Sequence Quality Detection**

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

417 **2.4.1 Feature-Based Representations of Genomic Se-** 418 **quences**

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

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

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

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

442 2.5 Synthesis of Chimera Detection Approaches

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

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

458 Chapter 3

459 Research Methodology

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

465 3.1 Research Activities

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

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

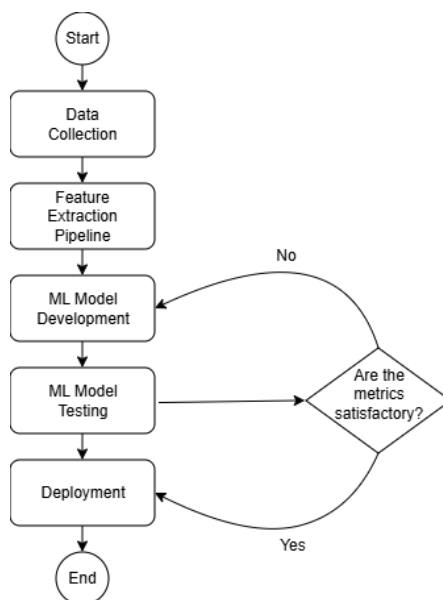


Figure 3.1: Process Diagram of Special Project

478 3.1.1 Data Collection

479 The mitochondrial genome reference sequence of *S. lemuru* was obtained from the
 480 NCBI database (accession number NC_039553.1) in FASTA format. This sequence
 481 served as the basis for generating simulated reads for model development.

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

484 proximately one (1) week.

485 Data Preprocessing

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

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

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

497 The command parameters are as follows:

- 498 • `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.
- 499 • `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability,
500 all set to a default value of 0.
- 501 • `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 502 • `-N`: number of read pairs, set to 10,000.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

549 **3.1.2 Feature Extraction Pipeline**

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

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

565 **Supplementary Alignment Flag**

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

568 ciated with reads that map to multiple genomic locations, suggesting a chimeric
569 structure. To extract this information, the script first checked whether the read
570 carried an `SA:Z` tag. If present, the tag string was parsed using the function
571 `parse_sa_tag`, yielding a structure for each alignment containing the reference
572 name, mapped position, strand, mapping 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 synthesized a summarization by aggregating metrics related to the
576 structure and reliability of the split alignments.

577 **K-mer Composition Difference**

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

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

591 Microhomology

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

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

610 3.1.3 Machine Learning Model Development

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

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

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

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

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

3.1.4 Model Benchmarking, Hyperparameter Optimization, and Evaluation

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

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

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

3.1.5 Feature Importance and Interpretation

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

For interpretability, individual features were grouped into four conceptual

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

695 **3.1.6 Validation and Testing**

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

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

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

711 **3.1.7 Documentation**

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

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

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

725 3.2 Calendar of Activities

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

Table 3.1: Timetable of Activities

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

Chapter 4

Results and Discussion

4.1 Descriptive Analysis of Features

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

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

740 4.1.1 Exploratory Data Analysis

741 An exploratory data analysis (EDA) was conducted on the extracted feature ma-
742 trix to characterize general patterns in the data and gain preliminary insight into
743 which variables might meaningfully contribute to classification. Histograms of
744 key features indicated that alignment-based variables showed clear class separa-
745 tion as chimeric reads have higher frequencies of split alignments and and no-
746 ticeably broader long-tailed distribution on soft-clipped regions (`softclip_left`
747 and `softclip_right`). In contrast, sequence-based variables such a microhomol-
748 ogy length and k-mer divergence displayed substantial overlap between classes,
749 suggesting more limited discriminative value. The complete set of histograms is
750 provided in Appendix 4.6.

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

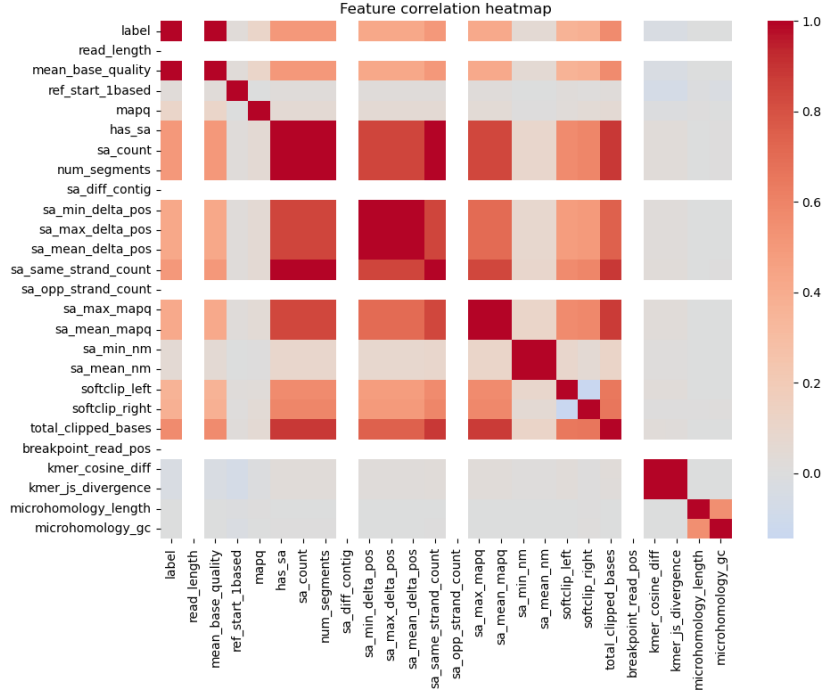


Figure 4.1: Feature correlation heatmap showing relationships among alignment-derived and sequence-derived variables.

4.2 Baseline Classification Performance

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

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

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

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

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

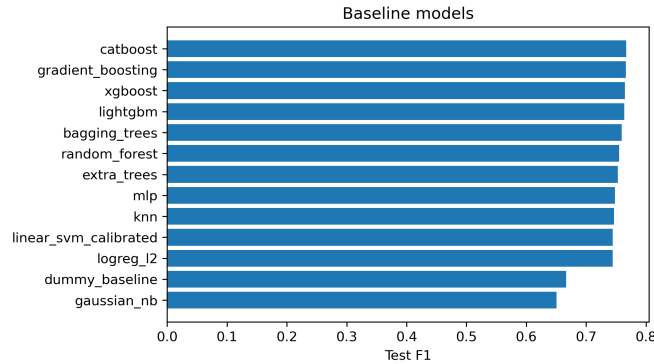


Figure 4.2: Test F1 of all baseline classifiers, showing that no single model clearly dominates and several achieve comparable performance.

4.3 Effect of Hyperparameter Tuning

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

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

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

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

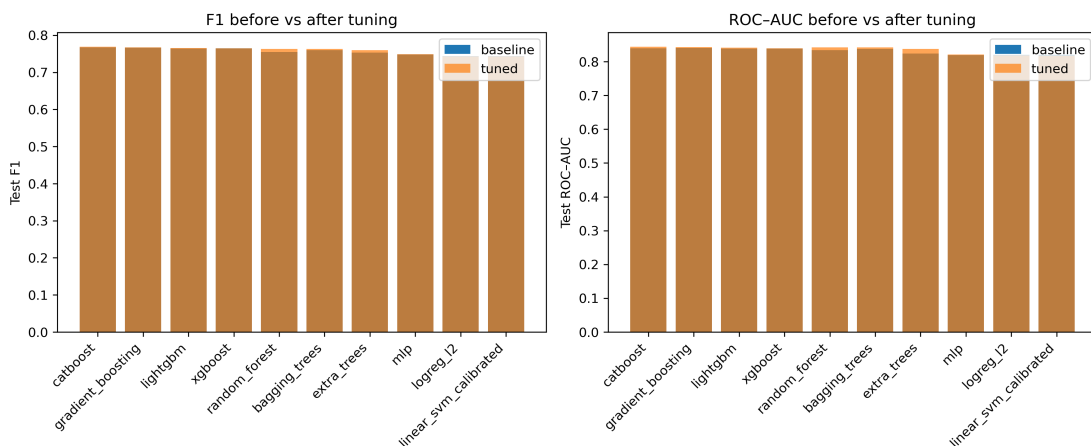


Figure 4.3: Comparison of test F1 (left) and ROC-AUC (right) for baseline and tuned models. Hyperparameter tuning yields small but consistent gains, particularly for tree-based ensembles.

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

4.4 Detailed Evaluation of Representative Models

For interpretability and diversity, four tuned models were selected for deeper analysis: CatBoost (best-performing boosted tree), scikit-learn gradient boosting (canonical gradient-boosting implementation), random forest (non-boosted ensemble baseline), and L2-regularised logistic regression (linear baseline). All models were trained on the engineered feature set and evaluated on the same held-out test data.

797 4.4.1 Confusion Matrices and Error Patterns

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

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

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

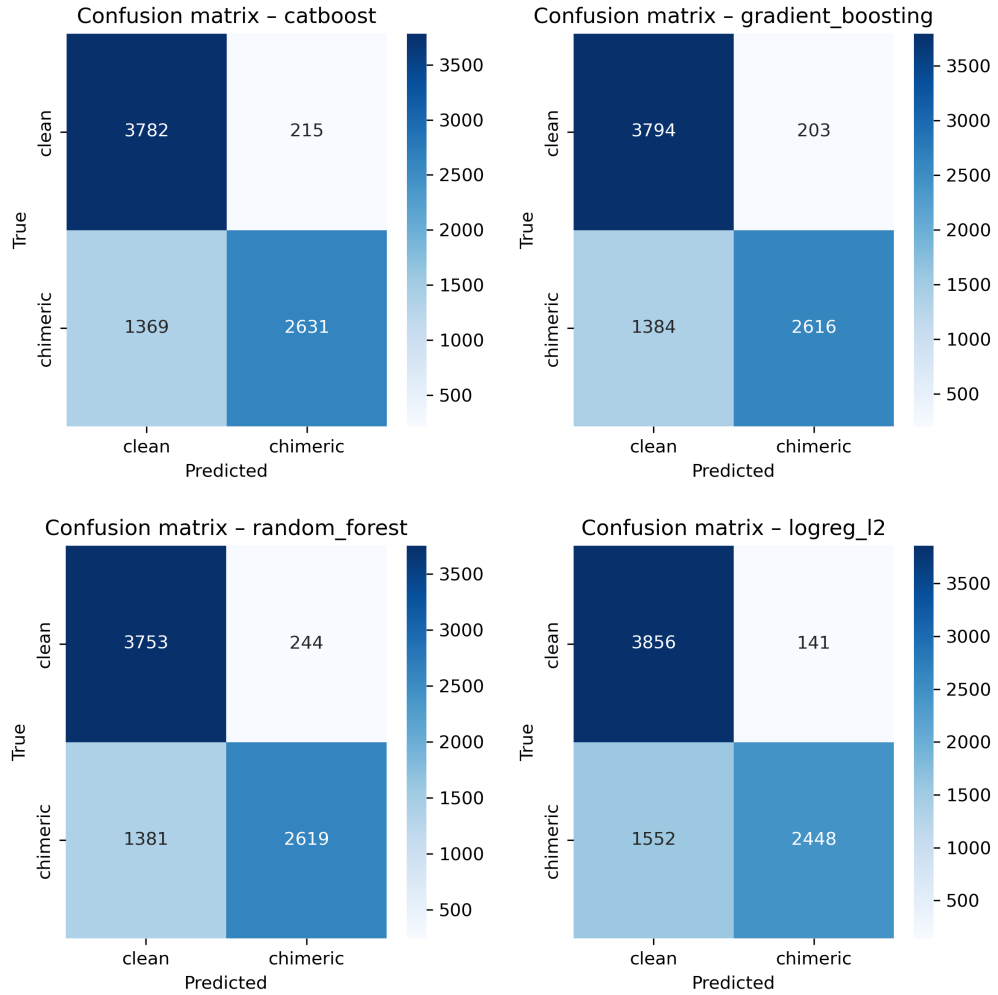


Figure 4.4: Confusion matrices for the four representative models on the held-out test set. All models show more false negatives (chimeric reads called clean) than false positives.

4.4.2 ROC and Precision–Recall Curves

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

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

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

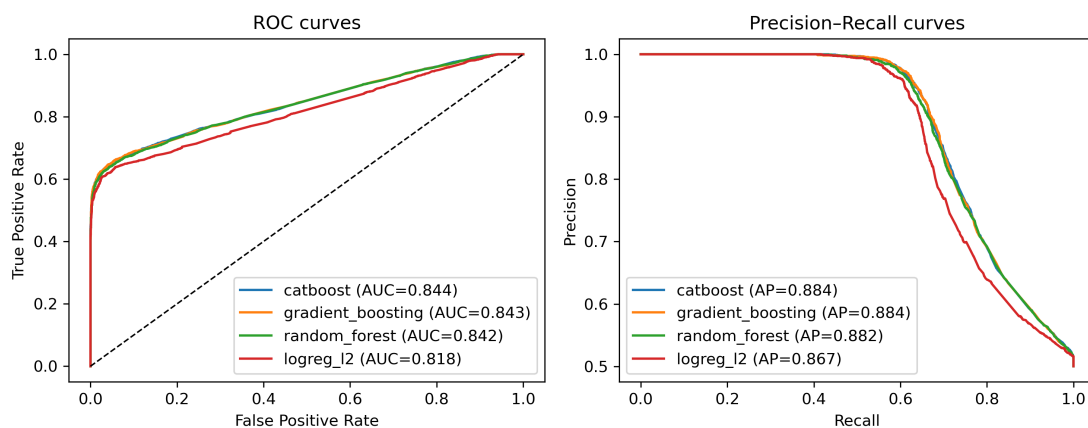


Figure 4.5: ROC (left) and precision–recall (right) curves for the four representative models on the held-out test set. Tree-based ensembles cluster closely, with logistic regression performing slightly but consistently worse.

827 4.5 Feature Importance and Biological Interpre- 828 tation

829 4.5.1 Permutation Importance of Individual Features

830 To understand how each classifier made predictions, feature importance was quan-
831 tified using permutation importance. In this approach, the values of a single fea-
832 ture are randomly shuffled, and the resulting drop in F_1 score (ΔF_1) reflects how
833 strongly the model depends on that feature. Greater decreases in F_1 indicate
834 stronger reliance on that feature. This analysis was applied to four representa-
835 tive models: CatBoost, Gradient Boosting, Random Forest, and L_2 -regularized
836 Logistic Regression.

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

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

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

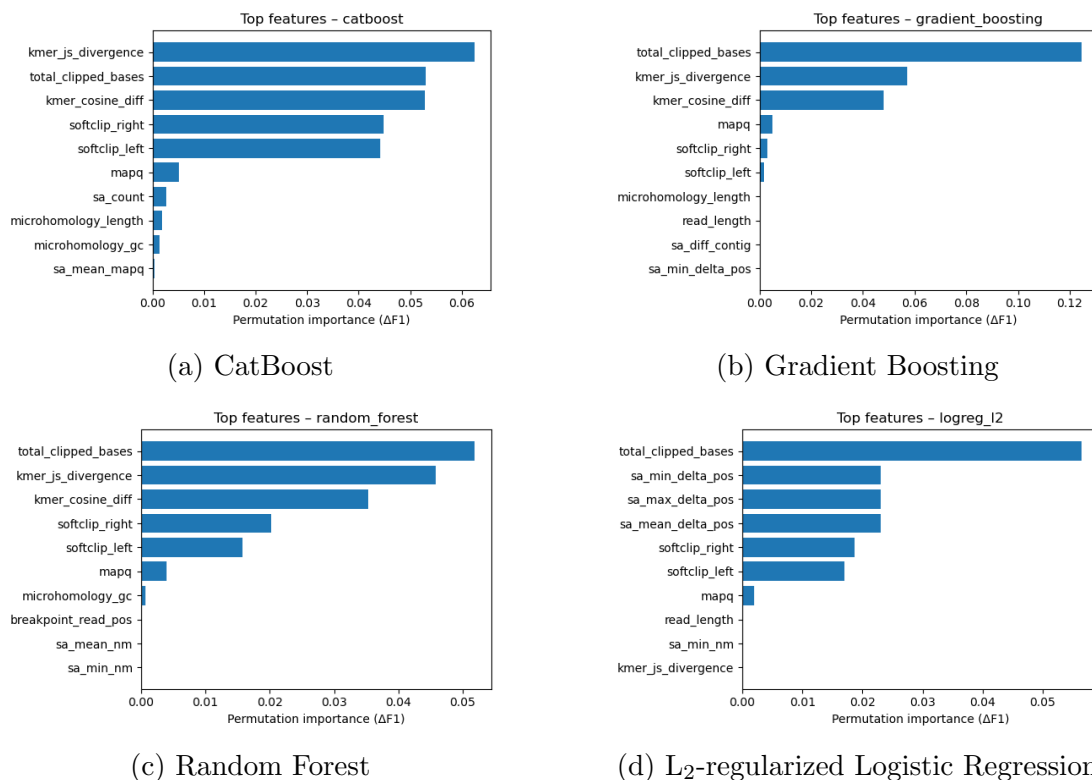


Figure 4.6: Permutation-based feature importance for four representative classifiers. Clipping and k-mer composition features are generally the strongest predictors, whereas microhomology and other alignment metrics contribute minimally.

852 4.5.2 Feature Family Importance

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

857 Kmer_jump (`kmer_cosine_diff`, `kmer_js_divergence`), `Micro_homology`, and
858 Other (e.g., `mapq`).

859 Aggregated analyses reveal consistent patterns across models. In CatBoost,
860 the Clipping family has the largest cumulative contribution (0.14), followed
861 by Kmer_jump (0.12), with Other features contributing modestly (0.005) and
862 SA_structure (0.003) and Micro_homology (0.003) providing minimal predictive
863 power. Gradient Boosting shows a similar trend, with Clipping (0.13) domi-
864 nating, Kmer_jump (0.11) secondary, and the remaining families contributing
865 negligibly. Random Forest integrates both Clipping (0.088) and Kmer_jump
866 (0.08) effectively, while SA_structure, Micro_homology, and Other remain minor
867 contributors. L₂-regularized Logistic Regression emphasizes Clipping (0.09)
868 and SA_structure (0.07), with Kmer_jump and Micro_homology having minimal
869 impact.

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

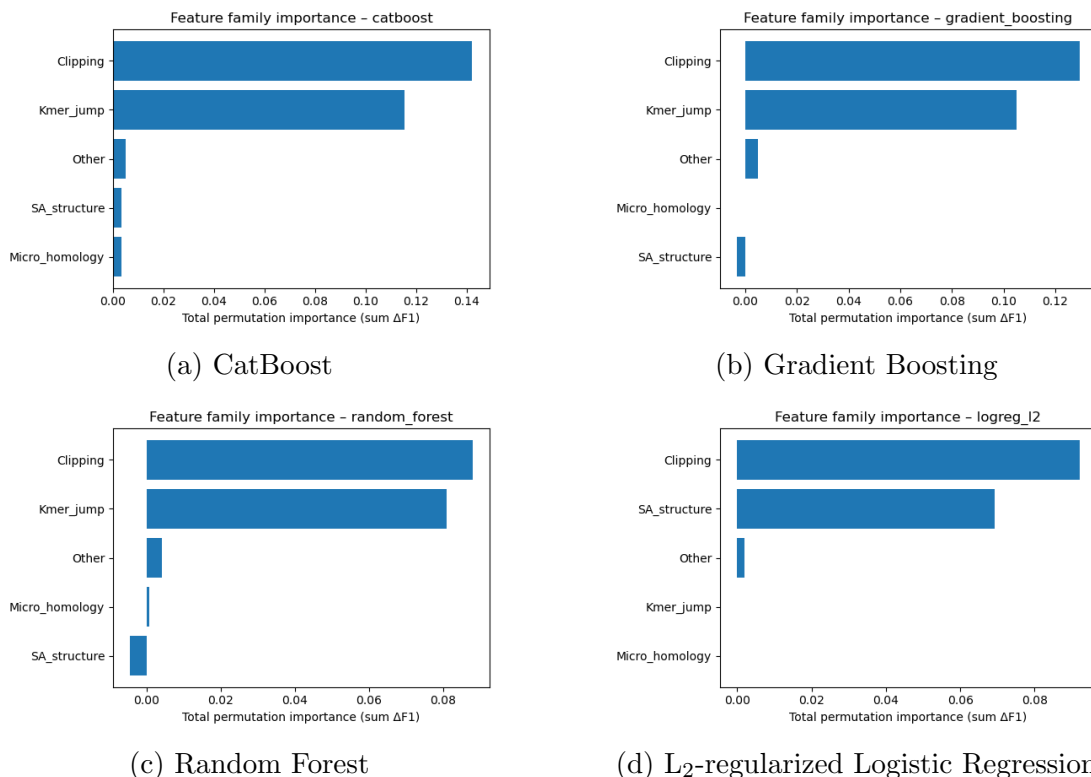


Figure 4.7: Aggregated feature family importance across four models. Clipping and k-mer compositional shifts are consistently the dominant contributors, while SA_structure, Micro_homology, and other features contribute minimally.

874 4.6 Summary of Findings

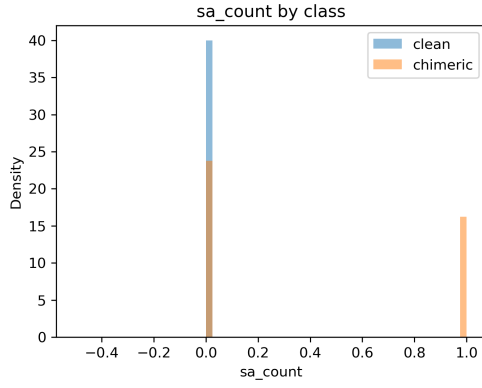
875 After removing trivially discriminative metadata, all models performed substan-
 876 tially better than the dummy baseline, with test F1-scores around 0.76 and ROC-
 877 AUC values near 0.84. Hyperparameter tuning yielded modest improvements,
 878 with boosting methods, particularly CatBoost and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

882 are costly.

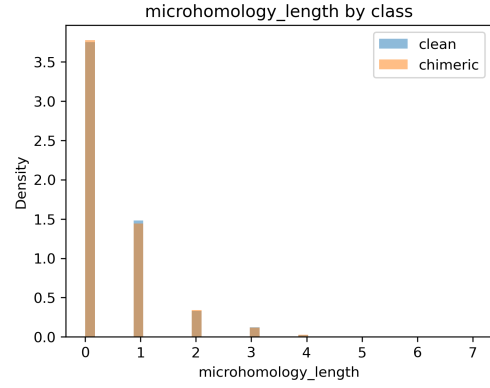
883 Feature importance analyses revealed that alignment disruptions, such as clip-
884 ping, and abrupt k-mer composition changes accounted for most predictive power.
885 In contrast, microhomology metrics and supplementary alignment descriptors con-
886 tributed minimally. These results indicate that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

890 **Appendix A**

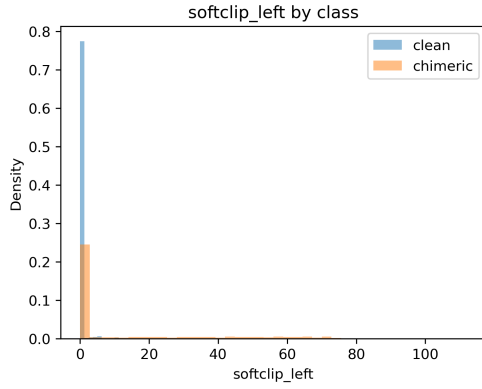
891 **Histograms of Key Features**



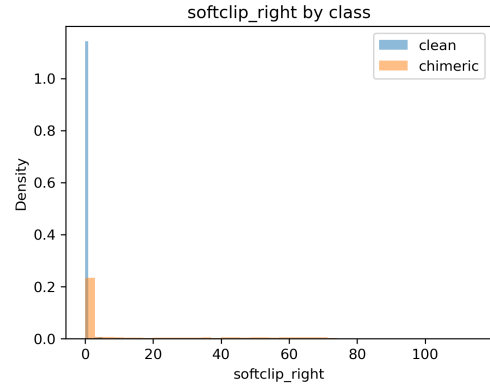
(a) sa_count



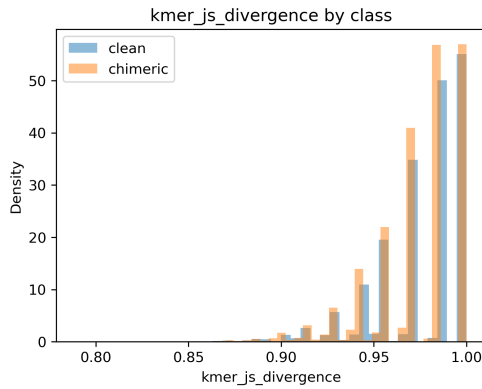
(b) Microhomology length



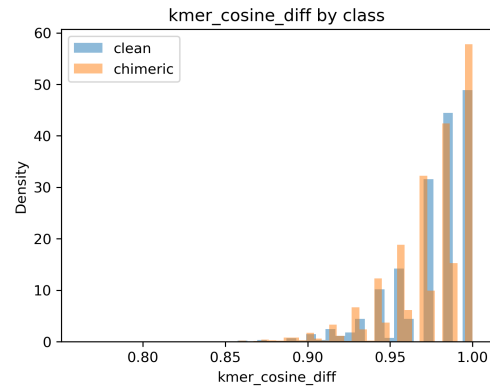
(c) softclip_left



(d) softclip_right



(e) k-mer Jensen-Shannon divergence



(f) k-mer cosine difference

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

892 References

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

911 45(4), e18. doi: 10.1093/nar/gkw955

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

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

914 CorpusID:88955007

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

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

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

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

919 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

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

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

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

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

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

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

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

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

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

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

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

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

932 10.1093/nar/gkt371

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

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

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

936 10.1186/s13059-020-02154-5

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

963 g3.117.300468

964 Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J.
 965 (2023). Effects of error, chimera, bias, and gc content on the accuracy of
 966 amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from [https://](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)
 967 journals.asm.org/doi/abs/10.1128/msystems.01025-23 doi: 10.1128/
 968 msystems.01025-23

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

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

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

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

983 Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler,
 984 A., & Schatz, M. (2018, 06). Accurate detection of complex structural
 985 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10
 986 .1038/s41592-018-0001-7

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

989 *Sciences*, 40(11), 701-714. Retrieved from <https://www.sciencedirect>
 990 [.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)
 991 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)
 992 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
 993 11). Large-scale machine learning for metagenomics sequence classifica-
 994 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
 995 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683
 996 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*
 997 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI
 998 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)
 999 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)