

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 Miag-ao, 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

Contents

21	1 Introduction	1
22	1.1 Overview	1
23	1.2 Problem Statement	3
24	1.3 Research Objectives	4
25	1.3.1 General Objective	4
26	1.3.2 Specific Objectives	4
27	1.4 Scope and Limitations of the Research	5
28	1.5 Significance of the Research	6
29	2 Review of Related Literature	7
30	2.1 The Mitochondrial Genome	7
31	2.1.1 Mitochondrial Genome Assembly	8

32	2.2	PCR Amplification and Chimera Formation	9
33	2.3	Existing Traditional Approaches for Chimera Detection	10
34	2.3.1	UCHIME	11
35	2.3.2	UCHIME2	12
36	2.3.3	CATch	13
37	2.3.4	ChimPipe	14
38	2.4	Machine Learning Approaches for Chimera and Sequence Quality	
39		Detection	15
40	2.4.1	Feature-Based Representations of Genomic Sequences . . .	16
41	2.5	Synthesis of Chimera Detection Approaches	18
42	3	Research Methodology	21
43	3.1	Research Activities	21
44	3.1.1	Data Collection	22
45	3.1.2	Bioinformatics Tools Pipeline	26
46	3.1.3	Machine Learning Model Development	29
47	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
48		Evaluation	30
49	3.1.5	Feature Importance and Interpretation	31

50	3.1.6 Validation and Testing	32
51	3.1.7 Documentation	33
52	3.2 Calendar of Activities	34
53	4 Results and Discussion	35
54	4.1 Baseline Classification Performance	35
55	4.2 Effect of Hyperparameter Tuning	37
56	4.3 Detailed Evaluation of Representative Models	38
57	4.3.1 Confusion Matrices and Error Patterns	39
58	4.3.2 ROC and Precision–Recall Curves	40

59 List of Figures

60	3.1	Process Diagram of Special Project	22
61	4.1	Test F1 of all baseline classifiers, showing that no single model	
62		clearly dominates and several achieve comparable performance. . .	36
63	4.2	Comparison of test F1 (left) and ROC–AUC (right) for baseline and	
64		tuned models. Hyperparameter tuning yields small but consistent	
65		gains, particularly for tree-based ensembles.	38
66	4.3	Confusion matrices for the four representative models on the held-	
67		out test set. All models show more false negatives (chimeric reads	
68		called clean) than false positives.	40
69	4.4	ROC (left) and precision–recall (right) curves for the four represen-	
70		tative models on the held-out test set. Tree-based ensembles cluster	
71		closely, with logistic regression performing slightly but consistently	
72		worse.	41

73 List of Tables

<small>74</small>	2.1 Comparison of Chimera Detection Methods	19
<small>75</small>	3.1 Timetable of Activities	34
<small>76</small>	4.1 Performance of baseline classifiers on the held-out test set.	36
<small>77</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

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

105 This study focuses on mitochondrial sequencing data from the genus *Sar-*
 106 *dinella*, a group of small pelagic fishes widely distributed in Philippine waters.
 107 Among them, *Sardinella lemuru* (Bali sardinella) is one of the country’s most
 108 abundant and economically important species, providing protein and livelihood
 109 to coastal communities (Labrador, Agmata, Palermo, Ravago-Gotanco, & Pante,
 110 2021; Willette, Bognot, Mutia, & Santos, 2011). Accurate mitochondrial assem-
 111 blies are critical for understanding its population genetics, stock structure, and
 112 evolutionary history. However, assembly pipelines often encounter errors or fail
 113 to complete due to undetected chimeric reads. To address this gap, this research
 114 introduces MitoChime, a machine learning pipeline designed to detect and filter
 115 PCR-induced chimeric reads using both alignment-based and sequence-derived
 116 statistical features. The tool aims to provide bioinformatics laboratories, partic-

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

119 1.2 Problem Statement

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

137 1.3 Research Objectives

138 1.3.1 General Objective

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

143 1.3.2 Specific Objectives

144 Specifically, the study aims to:

- 145 1. construct simulated *Sardinella lemuru* Illumina paired-end datasets contain-
146 ing both clean and PCR-induced chimeric reads,
- 147 2. extract alignment-based and sequence-based features such as k-mer compo-
148 sition, junction complexity, and split-alignment counts from both clean and
149 chimeric reads,
- 150 3. train, validate, and compare supervised machine-learning models for classi-
151 fying reads as clean or chimeric,
- 152 4. determine feature importance and identify indicators of PCR-induced
153 chimerism,
- 154 5. integrate the optimized classifier into a modular and interpretable pipeline
155 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: (1) 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; (2) to align the analysis with relevant *S. lemuru* sequencing projects at PGC Visayas; (3) 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 (4) 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.

179 1.5 Significance of the Research

180 This research provides both methodological and practical contributions to mi-
181 tochondrial genomics and bioinformatics. First, MitoChime filters PCR-induced
182 chimeric reads prior to genome assembly, with the goal of improving the con-
183 tiguity and correctness of *Sardinella lemuru* mitochondrial assemblies. Second,
184 it replaces informal manual curation with a documented workflow, improving au-
185 tomation and reproducibility. Third, the pipeline is designed to run on computing
186 infrastructures commonly available in regional laboratories, enabling routine use
187 at facilities such as PGC Visayas. Finally, more reliable mitochondrial assemblies
188 for *S. lemuru* provide a stronger basis for downstream applications in the field of
189 fisheries and genomics.

190 Chapter 2

191 Review of Related Literature

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

198 2.1 The Mitochondrial Genome

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

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

209 **2.1.1 Mitochondrial Genome Assembly**

210 Mitochondrial genome assembly refers to the reconstruction of the complete mito-
211 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is
212 conducted to obtain high-quality, continuous representations of the mitochondrial
213 genome that can be used for a wide range of analyses, including species identi-
214 fication, phylogenetic reconstruction, evolutionary studies, and investigations of
215 mitochondrial diseases. Because mtDNA evolves rapidly, its assembled sequence
216 provides valuable insights into population structure, lineage divergence, and adap-
217 tive evolution across taxa (Boore, 1999). Compared to nuclear genome assembly,
218 assembling the mitochondrial genome is often considered more straightforward but
219 still encounters technical challenges such as the formation of chimeric reads. Com-
220 monly used tools for mitogenome assembly such as GetOrganelle and MITObim
221 operate under the assumption of organelle genome circularity, and are vulnerable
222 when chimeric reads disrupt this circular structure, resulting in assembly errors
223 (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

268 sequence correspond to distinct high-abundance sequences.

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

274 **2.3.1 UCHIME**

275 UCHIME is one of the most widely used computational tools for detecting chimeric
276 sequences in amplicon sequencing data, as it serves as a critical quality control
277 step to prevent the misinterpretation of PCR artifacts as novel biological diversity.
278 The algorithm operates by searching for a model (M) where a query (Q) sequence
279 can be perfectly explained as a combination of two parent sequences, denoted as
280 A and B (Edgar et al., 2011).

281 In reference mode, UCHIME divides the query into four chunks and maps
282 them to a trusted chimeric-free database to identify candidate parents. It then
283 constructs a three-way alignment to calculate a score based on “votes.” A “Yes”
284 vote indicates the query aligns with parent A in one region and parent B in an-
285 other, while a “No” vote penalizes the score if the query diverges from the expected
286 chimeric model. In de novo mode, the algorithm operationalizes the abundance
287 skew principle described in Section 2.3. Instead of using an external database,
288 UCHIME dynamically treats the sample’s own high-abundance sequences as a
289 reference database, testing if lower-abundance sequences can be reconstructed as

290 mosaics of these internal ancestors (Edgar et al., 2011).

291 Despite its high sensitivity, UCHIME has inherent limitations rooted in
292 sequence divergence and database quality. The algorithm struggles to detect
293 chimeras formed from parents that are very closely related, specifically when the
294 sequence divergence between parents is less than roughly 0.8%, as the signal-to-
295 noise ratio becomes too low to distinguish a crossover event from sequencing error
296 (Edgar et al., 2011). Furthermore, in reference mode, the accuracy is strictly
297 bound by the completeness of the database; if true parents are absent, the tool
298 may fail to identify the chimera or produce false positives. Similarly, the de novo
299 mode relies on the assumption that parents are present and sufficiently more
300 abundant in the sample, which may not hold true in unevenly amplified samples
301 or complex communities.

302 **2.3.2 UCHIME2**

303 Building upon the original algorithm, UCHIME2 was developed to address the
304 nuances of high-resolution amplicon sequencing. A key contribution of the
305 UCHIME2 study was the critical re-evaluation of chimera detection benchmarks.
306 In the UCHIME2 paper (Edgar, 2016) and the UCHIME in practice website
307 (Edgar, n.d), the author has noted that the accuracy results reported in the
308 original UCHIME paper were “highly over-optimistic” because they relied on
309 unrealistic benchmark designs where parent sequences were assumed to be 100%
310 known and present. UCHIME2 introduced more rigorous testing (the CHSIMA
311 benchmark), revealing that “fake models,” where a valid biological sequence
312 perfectly mimics a chimera of two other valid sequences, are far more common

313 than previously assumed. This discovery suggests that error-free detection is
314 impossible in principle (Edgar, 2016). Another notable improvement is the in-
315 troduction of multiple application-specific modes that allow users to tailor the
316 algorithm’s performance to the characteristics of their datasets. The following
317 parameter presets: denoised, balanced, sensitive, specific, and high-confidence,
318 enable researchers to optimize the balance between sensitivity and specificity
319 according to the goals of their analysis.

320 However despite these advancements, the practical application of UCHIME2
321 requires caution. The author explicitly advises against using UCHIME2 as
322 a stand-alone tool in standard OTU clustering or denoising pipelines. Using
323 UCHIME2 as an independent filtering step in these workflows is discouraged, as
324 it often results in significantly higher error rates, increasing both false positives
325 (discarding valid sequences) and false negatives (retaining chimeras) (Edgar,
326 2016).

327 **2.3.3 CATch**

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

333 To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieurs
334 (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-

resents the first ensemble machine learning system designed for chimera detection in 16S rRNA amplicon sequencing. Rather than depending on a single detection strategy, CATCh integrates the outputs of several established tools, including UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual scores and binary decisions generated by these tools are used as input features for a supervised learning model. The algorithm employs a Support Vector Machine (SVM) with a Pearson VII Universal Kernel (PUK) to determine optimal weightings among the input features and to assign each sequence a probability of being chimeric.

Benchmarking in both reference-based and de novo modes demonstrated significant performance improvements. CATCh achieved sensitivities of approximately 85 percent in reference-based mode and 92 percent in de novo mode, with corresponding specificities of approximately 96 percent and 95 percent. These results indicate that CATCh detected 7 to 12 percent more chimeras than any individual algorithm while maintaining high precision.

2.3.4 ChimPipe

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

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

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

368 **2.4 Machine Learning Approaches for Chimera** 369 **and Sequence Quality Detection**

370 Traditional chimera detection tools rely primarily on heuristic or alignment-based
371 rules. Recent advances in machine learning (ML) have demonstrated that models
372 trained on sequence-derived features can effectively capture compositional and
373 structural patterns in biological sequences. Although most existing ML systems
374 such as those used for antibiotic resistance prediction, taxonomic classification,
375 or viral identification are not specifically designed for chimera detection, they
376 highlight how data-driven models can outperform similarity-based heuristics by
377 learning intrinsic sequence signatures. In principle, ML frameworks can integrate
378 indicators such as k-mer frequencies, GC-content variation and split-alignment

379 metrics to identify subtle anomalies that may indicate a chimeric origin (Arango
380 et al., 2018; Liang, Bible, Liu, Zou, & Wei, 2020; Ren et al., 2020).

381 **2.4.1 Feature-Based Representations of Genomic Se-** 382 **quences**

383 In genomic analysis, feature extraction converts DNA sequences into numerical
384 representations suitable for ML algorithms. A common approach is k-mer fre-
385 quency analysis, where normalized k-mer counts form the feature vector (Vervier,
386 Mahé, Tournoud, Veyrieras, & Vert, 2015). These features effectively capture lo-
387 cal compositional patterns that often differ between authentic and chimeric reads.
388 In particular, deviations in k-mer profiles between adjacent read segments can
389 serve as a compositional signature of template-switching events. Additional de-
390 scriptors such as GC content and sequence entropy can further distinguish se-
391 quence types; in metagenomic classification and virus detection, k-mer-based fea-
392 tures have shown strong performance and robustness to noise (Ren et al., 2020;
393 Vervier et al., 2015). For chimera detection specifically, abrupt shifts in GC or k-
394 mer composition along a read can indicate junctions between parental fragments.
395 Windowed feature extraction enables models to capture these discontinuities that
396 rule-based algorithms may overlook.

397 Machine learning models can also leverage alignment-derived features such as
398 the frequency of split alignments, variation in mapping quality, and local cover-
399 age irregularities. Split reads and discordant read pairs are classical indicators
400 of genomic junctions and have been formalized in probabilistic frameworks for
401 structural-variant discovery that integrate multiple evidence types (Layer, Hall, &

402 Quinlan, 2014). Similarly, long-read tools such as Sniffles employ split-alignment
 403 and coverage anomalies to accurately localize breakpoints (Sedlazeck et al., 2018).
 404 Modern aligners such as Minimap2 (Li, 2018) output supplementary (SA tags) and
 405 secondary alignments as well as chaining and alignment-score statistics that can
 406 be summarized into quantitative predictors for machine-learning models. These
 407 alignment-signal features are particularly relevant to PCR-induced mitochondrial
 408 chimeras, where template-switching events produce reads partially matching dis-
 409 tinct regions of the same or related genomes. Integrating such cues within a
 410 supervised-learning framework enables artifact detection even in datasets lacking
 411 complete or perfectly assembled references.

412 A further biologically grounded descriptor is the length of microhomology at
 413 putative junctions. Microhomology refers to short, shared sequences, often in the
 414 range of a few to tens of base pairs that are near breakpoints where template-
 415 switching events typically happen. Studies of double strand break repair and
 416 structural variation have demonstrated that the length of microhomology corre-
 417 lates with the likelihood of microhomology-mediated end joining (MMEJ) or fork-
 418 stalled template-switching pathways (Sfeir & Symington, 2015). In the context of
 419 PCR-induced chimeras, template switching during amplification often leaves short
 420 identical sequences at the junction of two concatenated fragments. Quantifying
 421 the longest exact suffix–prefix overlap at each candidate breakpoint thus provides
 422 a mechanistic signature of chimerism and complements both compositional (k-
 423 mer) and alignment (SA count) features.

424 **2.5 Synthesis of Chimera Detection Approaches**

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

Table 2.1: Comparison of Chimera Detection Methods

Methods	Approach	Limitations
Reference-based Chimera Detection	Compares query sequences against curated, non-chimeric reference databases; identifies mosaic sequences by evaluating similarity to known templates.	Depends heavily on completeness and quality of reference databases; often fails when novel taxa or missing parent sequences are present; reduced accuracy for low-divergence chimeras.
De novo Chimera Detection	Identifies chimeras using only internal dataset relationships; relies on abundance patterns and compositional similarity; reconstructs sequences as mosaics of high-abundance parents.	Assumes true sequences are more abundant—fails when amplification bias distorts abundance; struggles with evenly abundant parental sequences; can misclassify highly similar true variants.
UCHIME	Alignment-based chimera detection; segments query sequence, identifies parent candidates, performs 3-way alignment, and computes chimera scores; supports both reference-based and de novo modes.	Accuracy inflated in original benchmarks; suffers under incomplete databases; poor performance on low-divergence chimeras; sensitive to sequencing errors; misclassifies when parents are missing.
UCHIME2	Improved initial UCHIME benchmarking; offers multiple sensitivity/specificity modes; more robust with incomplete references; higher sensitivity.	Cannot achieve perfect accuracy due to “perfect fake models”; genuine variants may be indistinguishable from artificial recombinants; theoretical detection limit remains.
CATCh	First ML ensemble tool for 16S chimera detection; integrates outputs of UCHIME, ChimeraSlayer, DECIPHER, Pintail, Perseus via SVM classifier; significantly improves sensitivity and specificity.	Depends on performance of underlying tools; ML model limited to features they output; ensemble can still misclassify in datasets with extreme novelty or low coverage.
ChimPipe	Pipeline for detecting fusion genes and transcript-derived chimeras in RNA-seq; uses discordant paired-end reads and split-alignments; predicts isoforms and breakpoint coordinates.	Designed for RNA-seq, not amplicons; needs high-quality genome and annotation; computationally heavier; limited to organisms with reference genomes.

428 Across existing studies, no single approach reliably detects all forms of chimeric
429 sequences, particularly those generated by PCR-induced template switching in
430 mitochondrial genomes. Reference-based tools perform poorly when parental se-
431 quences are absent; de novo methods rely strongly on abundance assumptions;
432 alignment-based systems show reduced sensitivity to low-divergence chimeras; and
433 ensemble methods inherit the limitations of their component algorithms. RNA-
434 seq-oriented pipelines likewise do not generalize well to organelle data. Although
435 machine learning approaches offer promising feature-based detection, they are
436 rarely applied to mitochondrial genomes and are not trained specifically on PCR-
437 induced organelle chimeras. These limitations indicate a clear research gap: the
438 need for a specialized, feature-driven classifier tailored to mitochondrial PCR-
439 induced chimeras that integrates k-mer composition, split-alignment signals, and
440 micro-homology features to achieve more accurate detection than current heuristic
441 or alignment-based tools.

442 Chapter 3

443 Research Methodology

444 This chapter outlines the steps involved in completing the study, including data
445 gathering, generating simulated mitochondrial Illumina reads, preprocessing and
446 indexing the data, developing a bioinformatics pipeline to extract key features,
447 applying machine learning algorithms for chimera detection, and validating and
448 comparing model performance.

449 3.1 Research Activities

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

456 through a bioinformatics pipeline that extracted k-mer profiles, supplementary
 457 alignment (SA) features, and microhomology information to prepare the data for
 458 model construction. The machine learning model was trained using the processed
 459 input, and its precision and accuracy were assessed. It underwent tuning until it
 460 reached the desired performance threshold, after which it proceeded to validation
 461 and will undergo testing.

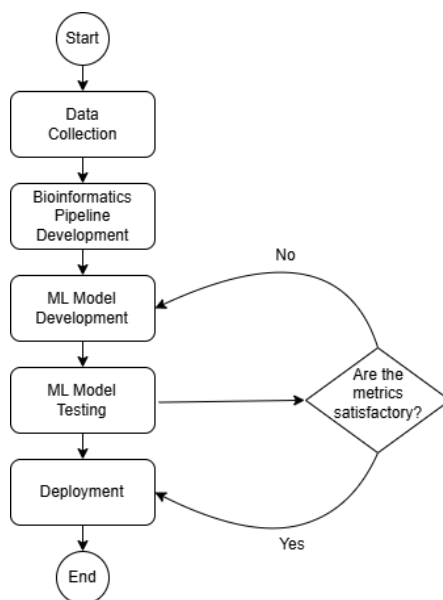


Figure 3.1: Process Diagram of Special Project

462 3.1.1 Data Collection

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

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

468 proximately one (1) week.

469 Data Preprocessing

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

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

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

481 The command parameters are as follows:

- 482 • `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.
- 483 • `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability,
484 all set to a default value of 0.
- 485 • `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 486 • `-N`: number of read pairs, set to 10,000.

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

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

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

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

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

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

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

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

510 The SAM files were then converted to BAM format, sorted, and indexed using
511 `samtools` (Version 1.20):

```
512 samtools view -bS clean.sam -o clean.bam
513 samtools view -bS chimeric.sam -o chimeric.bam
514
515 samtools sort clean.bam -o clean.sorted.bam
516 samtools index clean.sorted.bam
517
518 samtools sort chimeric.bam -o chimeric.sorted.bam
519 samtools index chimeric.sorted.bam
```

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

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

530 This whole process is scheduled to start in the second week of November 2025

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

533 **3.1.2 Bioinformatics Tools Pipeline**

534 A bioinformatics pipeline will be developed and implemented to extract the neces-
535 sary analytical features. This pipeline will function as a reproducible and modular
536 workflow that accepts FASTQ and BAM/SAM file inputs, processes them using
537 tools such as `samtools` and `jellyfish` (Version 2.3.1), and produces tabular fea-
538 ture matrices (TSV) for downstream machine learning. To ensure correctness
539 and adherence to best practices, bioinformatics experts at the PGC Visayas will
540 be consulted to validate the pipeline design, feature extraction logic, and overall
541 data integrity. This stage of the study is scheduled to begin in the first week of
542 January 2026 and conclude by the last week of February 2026, with an estimated
543 total duration of approximately two (2) months.

544 The bioinformatics pipeline focuses on three principal features from the simu-
545 lated and aligned sequencing data: (1) supplementary alignment flag (SA count),
546 (2) k-mer composition difference between read segments, and (3) microhomology
547 length at potential junctions. Each of these features captures a distinct biological
548 or computational signature associated with PCR-induced chimeras.

549 **Supplementary Alignment Flag**

550 Supplementary alignment information will be assessed using the mapped and
551 sorted BAM files (`clean.sorted.bam` and `chimeric.sorted.bam`) generated

552 from the data preprocessing stage. Alignment summaries will be checked using
553 `samtools flagstat` to obtain preliminary quality-control statistics, including
554 counts of primary, secondary, and supplementary (SA) alignments.

555 Both BAM files will be converted to SAM format for detailed inspection of
556 reads in each file:

```
557 samtools view -h clean.sorted.bam -o clean.sorted.sam
```

```
558 samtools view -h chimeric.sorted.bam -o chimeric.sorted.sam
```

559 The SAM output will be checked for reads containing the `SA:Z` flag, as it
560 denotes supplementary alignments. Reads exhibiting these or substantial soft-
561 clipped regions will be considered strong candidates for chimeric artifacts. A
562 custom Python script would be created to extract the alignment-derived features
563 and relevant metadata including mapping quality, SAM flag information, CIGAR-
564 based clipping, and alignment coordinates. These extracted attributes would then
565 be organized and compiled into a TSV (`.tsv`) file.

566 **K-mer Composition Difference**

567 Chimeric reads often comprise fragments from distinct genomic regions, resulting
568 in a compositional discontinuity between segments. Comparing k-mer frequency
569 profiles between the left and right halves of a read allows detection of such abrupt
570 compositional shifts, independent of alignment information. This will be obtained
571 using Jellyfish, a fast k-mer counting software. For each read, the sequence will
572 be divided into two segments, either at the midpoint or at empirically determined
573 breakpoints inferred from supplementary alignment data, to generate left and right

sequence segments. Jellyfish will then compute k-mer frequency profiles (with $k =$
5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized
and compared using distance metrics such as cosine similarity or Jensen–Shannon
divergence to quantify compositional disparity between the two halves of the same
read. The resulting difference scores will be stored in a structured TSV file.

Microhomology Length

The microhomology length was computed as part of the bioinformatics pipeline.
For each aligned read in the BAM files, the script first inferred a breakpoint
using the function `infer_breakpoint`, which represents a junction between two
segments. Breakpoints were determined primarily from soft-clipping patterns.
If no soft clips were present, SA tags were used to identify potential alignment
discontinuities.

Once a breakpoint was established, the script scanned a ± 40 base pair window
surrounding the breakpoint and used the function `longest_suffix_prefix_overlap`
to identify the longest exact suffix-prefix overlap between the left and right read
segments. This overlap, which represents consecutive bases shared at the junction,
was recorded as the microhomology length. Additionally, the GC content
of the overlapping sequence was calculated using the function `gc_content`, which
counts guanine (G) and cytosine (C) bases within the detected microhomology
and divides by the total length, yielding a proportion between 0 and 1.

Short microhomologies, typically 3-20 base pairs in length, are recognized signatures
of PCR-induced template switching and can promote template recombination
(Peccoud et al., 2018). Each read was annotated after capturing both the

length and GC content of microhomology.

3.1.3 Machine Learning Model Development

After feature extraction, the per-read feature matrices for clean and chimeric reads were merged into a single dataset. Each row corresponded to one paired-end read, and columns encoded alignment-structure features (e.g., supplementary alignment count and spacing between segments), CIGAR-derived soft-clipping statistics (e.g., left and right soft-clipped length, total clipped bases), k-mer composition discontinuity between read segments, and microhomology descriptors near candidate junctions. The resulting feature set was restricted to quantities that can be computed from standard BAM/FASTQ files in typical mitochondrial sequencing workflows.

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, L2-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 models for more detailed optimization. Specifically, ten model families were carried forward: L2-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 hyper-

parameters 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 bioinformatics 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 com-

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

3.1.6 Validation and Testing

Validation will involve both internal and external evaluations. Internal validation was achieved through five-fold cross-validation on the training data to verify model generalization and reduce variance due to random sampling. External validation will be achieved through testing on the 20% hold-out dataset derived from the simulated reads, which will be an unbiased benchmark to evaluate how well

689 the trained models generalized to unseen data. All feature extraction and pre-
690 processing steps were performed using the same bioinformatics pipeline to ensure
691 consistency and comparability across validation stages.

692 Comparative evaluation was performed across all candidate algorithms, in-
693 cluding a trivial dummy classifier, L2-regularized logistic regression, a calibrated
694 linear SVM, k-nearest neighbours, Gaussian Naïve Bayes, decision-tree ensembles,
695 gradient boosting methods, and a shallow MLP. This evaluation determined which
696 models demonstrated the highest predictive performance and computational effi-
697 ciency under identical data conditions. Their metrics were compared to identify
698 which algorithms were most suitable for further refinement.

699 **3.1.7 Documentation**

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

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

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

713 3.2 Calendar of Activities

714 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
 715 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	• • • •						
Bioinformatics Tools Pipeline			• • • •	• • • •			
Machine Learning Development			• •	• • • •	• • • •	• •	
Testing and Validation						• •	• • • •
Documentation	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •

716 Chapter 4

717 Results and Discussion

718 4.1 Baseline Classification Performance

719 Table 4.1 summarises the performance of eleven classifiers trained on the engi-
720 neered feature set using five-fold cross-validation and evaluated on the held-out
721 test set. All models were optimised using default hyperparameters, without ded-
722 icated tuning.

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

727 Across other models, test F1-scores clustered in a narrow band between ap-
728 proximately 0.74 and 0.77 and ROC-AUC values between 0.82 and 0.84. Gradi-
729 ent boosting, CatBoost, LightGBM, XGBoost, bagging trees, random forest, and

730 multilayer perceptron (MLP) all produced very similar scores, with CatBoost and
731 gradient boosting slightly ahead (test F1 ≈ 0.77 , ROC-AUC ≈ 0.84). Linear
732 models (logistic regression and calibrated linear SVM) performed only marginally
733 worse (test F1 ≈ 0.74), while Gaussian Naive Bayes lagged behind with substan-
734 tially lower F1 (≈ 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

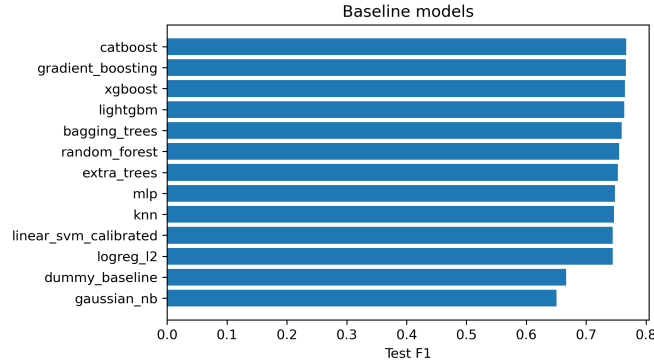


Figure 4.1: Test F1 of all baseline classifiers, showing that no single model clearly dominates and several achieve comparable performance.

4.2 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

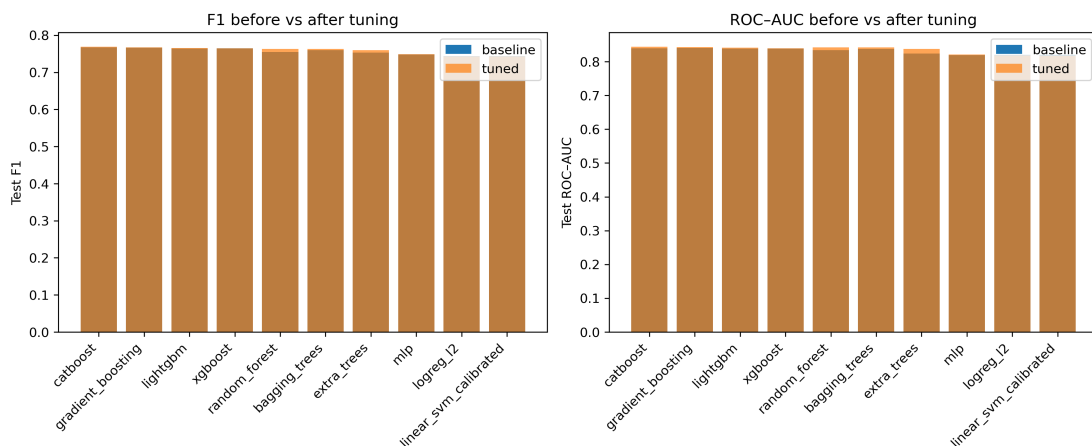


Figure 4.2: Comparison of test F1 (left) and ROC-AUC (right) for baseline and tuned models. Hyperparameter tuning yields small but consistent gains, particularly for tree-based ensembles.

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.

4.3 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.

759 4.3.1 Confusion Matrices and Error Patterns

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

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

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

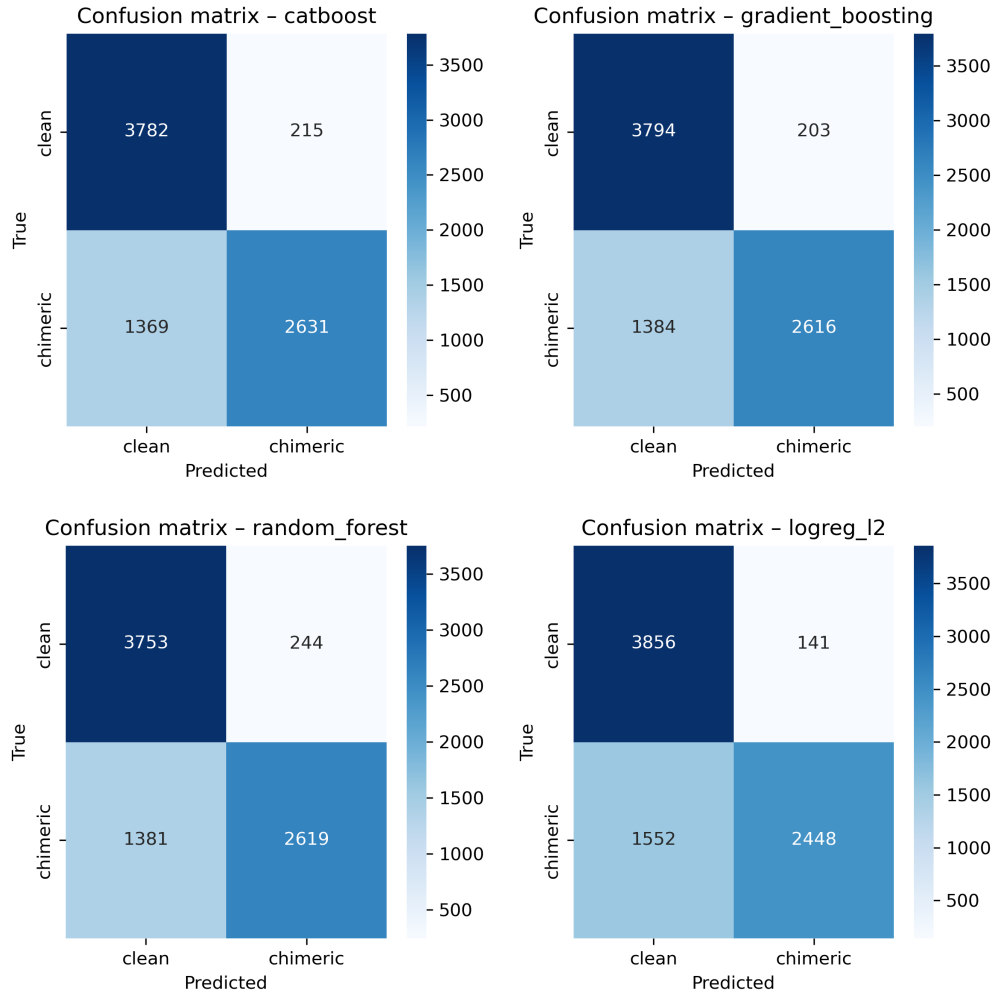


Figure 4.3: Confusion matrices for the four representative models on the held-out test set. All models show more false negatives (chimeric reads called clean) than false positives.

4.3.2 ROC and Precision–Recall Curves

Receiver operating characteristic (ROC) and precision–recall (PR) curves (Figure 4.4) 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.

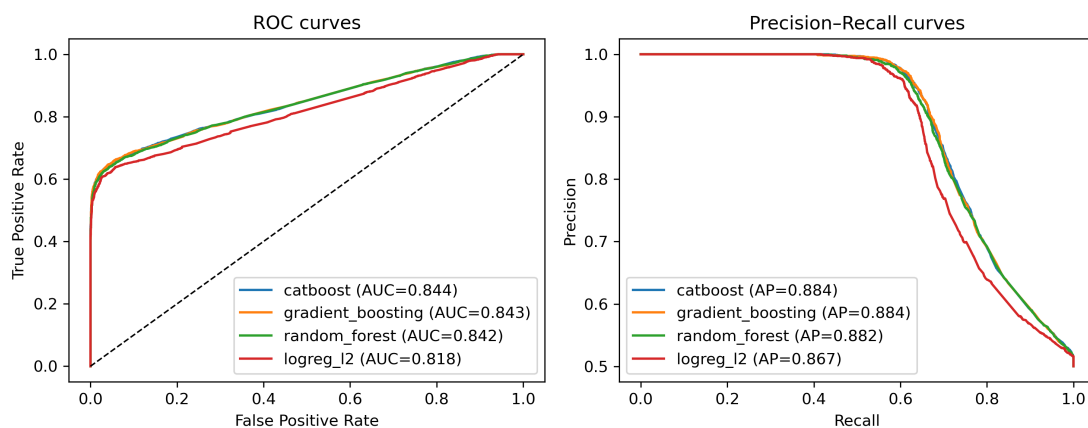


Figure 4.4: ROC (left) and precision–recall (right) curves for the four representative models on the held-out test set. Tree-based ensembles cluster closely, with logistic regression performing slightly but consistently worse.

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*,

808 45(4), e18. doi: 10.1093/nar/gkw955

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

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

811 CorpusID:88955007

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

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

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

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

816 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

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

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

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

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

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

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

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

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

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

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

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

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

829 10.1093/nar/gkt371

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

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

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

833 10.1186/s13059-020-02154-5

834 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
835 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
836 1819–1825. doi: 10.1093/nar/26.7.1819

837 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
838 (2021). Mitochondrial dna reveals genetically structured haplogroups of
839 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
840 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

841 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework
842 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-
843 -15-6-r84

844 Li, H. (2018, 05). Minimap2: pairwise alignment for nucleotide sequences. *Bioin-*
845 *formatics*, *34*(18), 3094-3100. Retrieved from [https://doi.org/10.1093/](https://doi.org/10.1093/bioinformatics/bty191)
846 [bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191) doi: 10.1093/bioinformatics/bty191

847 Liang, Q., Bible, P. W., Liu, Y., Zou, B., & Wei, L. (2020, 02). Deepmi-
848 crobes: taxonomic classification for metagenomics with deep learning. *NAR*
849 *Genomics and Bioinformatics*, *2*(1), lqaa009. Retrieved from [https://](https://doi.org/10.1093/nargab/lqaa009)
850 doi.org/10.1093/nargab/lqaa009 doi: 10.1093/nargab/lqaa009

851 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*
852 *Reviews Genetics*, *11*(1), 31–46. doi: 10.1038/nrg2626

853 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
854 an ensemble classifier for chimera detection in 16s rna sequencing stud-
855 ies. *Applied and Environmental Microbiology*, *81*(5), 1573-1584. Retrieved
856 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
857 [10.1128/AEM.02896-14](https://journals.asm.org/doi/abs/10.1128/aem.02896-14)

858 Peccoud, J., Lequime, S., Moltini-Conclois, I., Giraud, I., Lambrechts, L., &
859 Gilbert, C. (2018, 04). A survey of virus recombination uncovers canon-

ical features of artificial chimeras generated during deep sequencing library preparation. *G3 Genes—Genomes—Genetics*, 8(4), 1129-1138. Retrieved from <https://doi.org/10.1534/g3.117.300468> doi: 10.1534/g3.117.300468

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

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

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

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

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

Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler, A., & Schatz, M. (2018, 06). Accurate detection of complex structural variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10

886 .1038/s41592-018-0001-7

887 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
888 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
889 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
890 [.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)
891 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

892 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
893 11). Large-scale machine learning for metagenomics sequence classifica-
894 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
895 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683

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