

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 Dimzon**

19 October 24, 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	5
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	2.2.1	Effects of Chimeric Reads on Organelle Genome Assembly	10
34	2.3	Existing Traditional Approaches for Chimera Detection	11
35	2.3.1	UCHIME	12
36	2.3.2	UCHIME2	14
37	2.3.3	CATch	16
38	2.3.4	ChimPipe	17
39	2.4	Machine Learning Approaches for Chimera and Sequence Quality	
40		Detection	18
41	2.4.1	Feature-Based Representations of Genomic Sequences . . .	18
42	3	Research Methodology	21
43	3.1	Research Activities	21
44	3.1.1	Data Collection	22
45	3.1.2	Data Simulation	23
46	3.1.3	Bioinformatics Tools Pipeline	24
47	3.1.4	Machine-Learning Model Development	27
48	3.1.5	Validation and Testing	28

49	3.1.6 Documentation	28
50	3.2 Calendar of Activities	29

51 List of Figures

<small>52</small>	3.1 Process Diagram of Special Project	22
-------------------	--	----

53 List of Tables

<small>54</small>	3.1 Timetable of Activities	29
-------------------	---------------------------------------	----

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 mis-joined 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 and automated 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- and sequence-derived statistical features. The tool aims to provide bioinformatics laboratories, particularly the

94 Philippine Genome Center Visayas, with an efficient, interpretable, and resource-
95 optimized solution for improving mitochondrial genome reconstruction.

96 1.2 Problem Statement

97 While NGS technologies have revolutionized genomic data acquisition, the ac-
98 curacy of mitochondrial genome assembly remains limited by artifacts produced
99 during PCR amplification. These chimeric reads can distort assembly graphs and
100 cause misassemblies, with especially severe effects in small, circular mitochon-
101 drial genomes (Boore, 1999; Cameron, 2014). Existing assembly pipelines such
102 as GetOrganelle, MITObim, and NOVOPlasty assume that sequencing reads are
103 free of such artifacts (Dierckxsens et al., 2017; Hahn et al., 2013; Jin et al., 2020).
104 At the Philippine Genome Center Visayas, several mitochondrial assemblies have
105 failed or yielded incomplete contigs despite sufficient coverage, suggesting that
106 undetected chimeric reads compromise assembly reliability. Meanwhile, exist-
107 ing chimera-detection tools such as UCHIME and VSEARCH were developed
108 primarily for amplicon-based microbial community analysis and rely heavily on
109 reference or taxonomic comparisons (Edgar, Haas, Clemente, Quince, & Knight,
110 2011; Rognes, Flouri, Nichols, Quince, & Mahé, 2016). These approaches are un-
111 suitable for single-species organellar data, where complete reference genomes are
112 often unavailable. Therefore, there is a pressing need for a reference-independent,
113 data-driven tool capable of automatically detecting and filtering PCR-induced
114 chimeras in mitochondrial sequencing datasets.

115 **1.3 Research Objectives**

116 **1.3.1 General Objective**

117 To develop and evaluate a machine-learning-based pipeline (MitoChime) capable
118 of detecting PCR-induced chimeric reads in *Sardinella* mitochondrial sequencing
119 data to improve the accuracy of mitochondrial genome assembly.

120 **1.3.2 Specific Objectives**

121 Specifically, the researchers aim to:

- 122 1. Construct simulated and empirical *Sardinella* Illumina paired-end datasets
123 containing both clean and PCR-induced chimeric reads.
- 124 2. Extract alignment- and sequence-based features (e.g., k-mer composition,
125 junction complexity, split-alignment counts) from both clean and chimeric
126 reads.
- 127 3. Train, validate, and compare supervised machine-learning models (e.g., Ran-
128 dom Forest, XGBoost) for classifying reads as clean or chimeric.
- 129 4. Determine feature importance and identify the most informative indicators
130 of PCR-induced chimerism.
- 131 5. Integrate the optimized classifier into a modular and interpretable pipeline
132 deployable on standard computing environments at PGC Visayas.

133 1.4 Scope and Limitations of the Research

134 This study focuses on detecting PCR-induced chimeric reads in Illumina paired-
135 end mitochondrial sequencing data from *Sardinella* species. The work emphasizes
136 `wgsim` simulations and selected empirical data obtained from open-access genomic
137 repositories such as the National Center for Biotechnology Information (NCBI).
138 The study excludes naturally occurring chimeras, nuclear mitochondrial pseudo-
139 genes (NUMTs), and large-scale structural rearrangements in nuclear genomes.
140 Feature extraction prioritizes interpretable, shallow statistics and alignment met-
141 rics rather than deep-learning embeddings to ensure transparency and computa-
142 tional efficiency. Testing on long-read platforms (e.g., Nanopore, PacBio) and
143 other taxa lies beyond the project’s scope. The resulting pipeline will serve as a
144 foundation for future, broader chimera-detection frameworks applicable to diverse
145 organellar genomes.

146 1.5 Significance of the Research

147 This research provides both methodological and practical contributions to mi-
148 tochondrial genomics and bioinformatics. First, MitoChime enhances assembly
149 accuracy by filtering PCR-induced chimeras prior to genome assembly, thereby
150 improving the contiguity and correctness of *Sardinella* mitochondrial genomes.
151 Second, it promotes automation and reproducibility by replacing subjective man-
152 ual curation with a data-driven, machine-learning-based workflow. Third, the
153 pipeline demonstrates computational efficiency through its design, enabling im-
154 plementation on modest computing infrastructures commonly available in regional

155 laboratories. Beyond technical improvements, MitoChime contributes to local ca-
156 pacity building by strengthening expertise in bioinformatics and machine-learning
157 integration, aligning with the mission of the Philippine Genome Center Visayas.
158 Finally, accurate mitochondrial assemblies are vital for fisheries management,
159 population genetics, and biodiversity conservation, providing reliable genomic re-
160 sources for species such as *Sardinella*. Through these contributions, MitoChime
161 advances the reliability of mitochondrial genome reconstruction and supports sus-
162 tainable, data-driven research in Philippine genomics.

Chapter 2

Review of Related Literature

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

2.1 The Mitochondrial Genome

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

177 animal species, the mitochondrial genome ranges from 15–20 kb and contains 13
178 protein-coding genes, 22 tRNAs, and two rRNAs arranged compactly without in-
179 trons (Gray, 2012). In comparison to nuclear DNA the ratio of the number of
180 copies of mtDNA is higher and has relatively simple organization which make it
181 particularly suitable for genome sequencing and assembly studies (Dierckxsens et
182 al., 2017). Moreover, mitochondrial genomes provide crucial insights into evo-
183 lutionary relationships among species and are increasingly used for testing new
184 genomic assembly and analysis methods.

185 **2.1.1 Mitochondrial Genome Assembly**

186 Mitochondrial genome assembly refers to the reconstruction of the complete mito-
187 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is
188 conducted to obtain high-quality, continuous representations of the mitochondrial
189 genome that can be used for a wide range of analyses, including species identi-
190 fication, phylogenetic reconstruction, evolutionary studies, and investigations of
191 mitochondrial diseases. Because mtDNA evolves relatively rapidly and is mater-
192 nally inherited, its assembled sequence provides valuable insights into population
193 structure, lineage divergence, and adaptive evolution across taxa (Boore, 1999).
194 Compared to nuclear genome assembly, assembling the mitochondrial genome is
195 often considered more straightforward but still encounters distinct technical chal-
196 lenges such as sequencing errors, low coverage regions, and chimeric reads that can
197 distort the final assembly, leading to incomplete or misassembled genomes. These
198 errors can propagate into downstream analyses, emphasizing the need for robust
199 chimera detection and sequence validation methods in mitochondrial genome re-

200 search.

201 **2.2 PCR Amplification and Chimera Formation**

202 Polymerase Chain Reaction (PCR) plays an important role in next-generation
203 sequencing (NGS) library preparation, as it amplifies target DNA fragments for
204 downstream analysis. However, the amplification process can also introduce arti-
205 facts that affect data accuracy, one of them being the formation of chimeric se-
206 quences. Chimeras typically arise when incomplete extension occurs during a PCR
207 cycle. This causes the DNA polymerase to switch from one template to another
208 and generate hybrid recombinant molecules (Judo et al., 1998). Artificial chimeras
209 are produced through such amplification errors, whereas biological chimeras oc-
210 cur naturally through genomic rearrangements or transcriptional events. These
211 biological chimeras can have functional roles and may encode tissue-specific novel
212 proteins that link to cellular processes or diseases (Frenkel-Morgenstern et al.,
213 2012).

214 In the context of amplicon-based sequencing, PCR-induced chimeras can sig-
215 nificantly distort analytical outcomes. Their presence artificially inflates estimates
216 of genetic or microbial diversity and may cause misassemblies during genome re-
217 construction. (Qin et al., 2023) has reported that chimeric sequences may account
218 for more than 10% of raw reads in amplicon datasets. This artifact tends to be
219 most prominent among rare operational taxonomic units (OTUs) or singletons,
220 which are sometimes misinterpreted as novel diversity, which further causes the
221 complication of microbial diversity analyses (Gonzalez, Zimmermann, & Saiz-

222 Jimenez, 2004). Moreover, the likelihood of chimera formation has been found to
223 vary with the GC content of target sequences, with lower GC content generally
224 associated with a reduced rate of chimera generation (Qin et al., 2023).

225 **2.2.1 Effects of Chimeric Reads on Organelle Genome As-** 226 **sembly**

227 In mitochondrial DNA (mtDNA) assembly workflows, PCR-induced chimeras pose
228 additional challenges. Assembly tools such as GetOrganelle and MitoBeam, which
229 operate under the assumption of organelle genome circularity, are vulnerable when
230 chimeric reads disrupt this circular structure. Such disruptions can lead to assem-
231 bly errors or misassemblies (Bi et al., 2024). These artificial sequences interfere
232 with the assembly graph, which makes it more difficult to accurately reconstruct
233 mitochondrial genomes. In addition, these artifacts propagate false variants and
234 erroneous annotations in genomic data. Hence, determining and minimizing PCR-
235 induced chimera formation is vital for improving the quality of mitochondrial
236 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 microbial community studies. These methods rely on the comparison of each query sequence against a curated, high-quality database of known, non-chimeric reference sequences to determine whether the query can be more plausibly explained as a composite or a mosaic of two or more reference sequences rather than as a genuine biological variant (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. Instead of comparing each query sequence to a curated collection of known, non-chimeric sequences, de novo methods infer chimeras based on internal relationships among the sequences present within the dataset itself. This approach is particularly advantageous in studies of novel, under explored, or taxonomically diverse microbial communities where comprehensive reference databases are unavailable or incomplete (Edgar, 2016; Edgar et al., 2011). The underlying assumption on this method operates on the key biological principle that true biological sequences are generally more abundant than chimeric artifacts. During PCR amplification, authentic sequences are amplified early and tend to dominate the read pool, while

261 chimeric sequences form later resulting in the tendency to appear at lower relative
262 abundances compared to their true parental sequences. As such, the abundance
263 hierarchy is formed by treating the most abundant sequences as supposed parents
264 and testing whether less abundant sequences can be reconstructed as mosaics of
265 these dominant templates. In addition to abundance, de novo algorithms assess
266 compositional and structural similarity among sequences, examining whether cer-
267 tain regions of a candidate sequence align more closely with one high-abundance
268 sequence and other regions with a different one.

269 Both reference-based and de novo approaches are complementary rather than
270 mutually exclusive. Reference-based methods provide stability and reproducibility
271 when curated databases are available, whereas de novo methods offer flexibility
272 and independence for novel or highly diverse communities. In practice, many
273 modern bioinformatics pipelines combine both paradigms sequentially: an initial
274 de novo step identifies dataset-specific chimeras, followed by a reference-based pass
275 that removes remaining artifacts relative to established databases (Edgar, 2016).
276 These two methods of detection form the foundation of tools such as UCHIME
277 and later UCHIME2, exemplified by the dual capability of providing both modes
278 within a unified computational framework.

279 **2.3.1 UCHIME**

280 Developed by Edgar et al. (Edgar et al., 2011), UCHIME is one of the most widely
281 used computational tools for detecting chimeric sequences in amplicon sequencing
282 data. The UCHIME algorithm detects chimeras by evaluating how well a query
283 sequence (Q) can be explained as a mosaic of two parent sequences (A and B)

284 from a reference database. The query sequence is first divided into four non-
285 overlapping segments or chunks. Each chunk is independently searched against a
286 reference database that is assumed to be free of chimeras. The best matches to
287 each segment are collected, and from these results, two candidate parent sequences
288 are identified, typically the two sequences that best explain all chunks of the query.
289 Then a three-way alignment among the query (Q) and the two parent candidates
290 (A and B) is done. From this alignment, UCHIME attempts to find a chimeric
291 model (M) which is a hypothetical recombinant sequence formed by concatenating
292 fragments from A and B that best match the observed Q

293 **Chimeric Alignment and Scoring**

294 To decide whether a query is chimeric, UCHIME computes several alignment-
295 based metrics between Q, its top hit (T, the most similar known sequence), and
296 the chimeric model (M). The key differences are measured as: dQT or the number
297 of mismatches between the query and the top hit as well as dQM or the number
298 of mismatches between the query and the chimeric model. From these, a chimera
299 score is calculated to quantify how much better the chimeric model fits the query
300 compared to a single parent. If the model's similarity to Q exceeds a defined
301 threshold (typically $\geq 0.8\%$ better identity), the sequence is reported as chimeric.
302 A higher score indicates stronger evidence of chimerism, while lower scores suggest
303 that the sequence is more likely to be authentic.

304 In de novo mode, UCHIME applies an abundance-driven strategy. Only se-
305 quences at least twice as abundant as the query are considered as potential parents.
306 Non-chimeric sequences identified at each step are added iteratively to a growing

307 internal database for subsequent queries.

308 **Limitations of UCHIME**

309 Although UCHIME was a significant advancement in chimera detection, it has
310 notable limitations. According to (Edgar, 2016) and the UCHIME practical notes
311 (Edgar, n.d), many of the accuracy results reported in the original 2011 paper
312 were overly optimistic due to unrealistic benchmark designs that assumed com-
313 plete reference coverage and perfect sequence quality. In practice, UCHIME’s
314 accuracy can decline when: (1) The reference database is incomplete or contains
315 erroneous entries. (2) Low-divergence chimeras are present, as these closely resem-
316 ble genuine biological variants. (3) Sequence datasets include residual sequencing
317 errors, leading to spurious alignments or misidentification; and (4) The abundance
318 ratio between parent and chimera is distorted by amplification bias. Additionally,
319 UCHIME tends to misclassify sequences as non-chimeric when parent sequences
320 are missing from the database. These limitations motivated the development of
321 UCHIME2.

322 **2.3.2 UCHIME2**

323 To overcome the limitations of its predecessor, UCHIME2 (Edgar, 2016) intro-
324 duced several methodological and algorithmic refinements that significantly en-
325 hanced the accuracy and reliability of chimera detection. One major improve-
326 ment lies in its approach to uncertainty handling. In earlier versions, sequences
327 with limited reference support were often incorrectly classified as non-chimeric,

328 increasing the likelihood of false negatives. UCHIME2 addresses this issue by
329 designating such ambiguous sequences as “unknown,” thereby providing a more
330 conservative and reliable classification framework.

331 Another notable advancement is the introduction of multiple application-
332 specific modes that allow users to tailor the algorithm’s performance to the
333 characteristics of their datasets. The following parameter presets: denoised,
334 balanced, sensitive, specific, and high-confidence, enable researchers to optimize
335 the balance between sensitivity and specificity according to the goals of their
336 analysis.

337 In comparative evaluations, UCHIME2 demonstrated superior detection per-
338 formance, achieving sensitivity levels between 93% and 99% and lower overall
339 error rates than earlier versions or other contemporary tools such as DECIPHER
340 and ChimeraSlayer. Despite these advances, the study also acknowledged a fun-
341 damental limitation in chimera detection: complete error-free identification is
342 theoretically unattainable. This is due to the presence of “perfect fake models,”
343 wherein genuine non-chimeric sequences can be perfectly reconstructed from other
344 reference fragments. This underscore the uncertainty in differentiating authentic
345 biological sequences from artificial recombinants based solely on sequence similar-
346 ity, emphasizing the need for continued methodological refinement and cautious
347 interpretation of results.

2.3.3 CATCh

Early chimera detection programs such as UCHIME (Edgar et al., 2011) relied on alignment-based and abundance-based heuristics to identify hybrid sequences in amplicon data. However, researchers soon observed that different algorithms often produced inconsistent predictions. A sequence might be identified as chimeric by one tool but classified as non-chimeric by another, resulting in unreliable filtering outcomes across studies.

To address these inconsistencies, (Mysara, Saeys, Leys, Raes, & Monsieurs, 2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which represents 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. Integration of CATCh into amplicon-

372 processing pipelines also reduced operational taxonomic unit (OTU) inflation by
373 23 to 35 percent, producing diversity estimates that more closely reflected true
374 community composition.

375 **2.3.4 ChimPipe**

376 Among the available tools for chimera detection, ChimPipe is a bioinformat-
377 ics pipeline developed to identify chimeric sequences such as fusion genes and
378 transcription-induced chimeras from paired-end RNA sequencing data. It uses
379 both discordant paired-end reads and split-read alignments to improve the ac-
380 curacy and sensitivity of detecting fusion genes, trans-splicing events, and read-
381 through transcripts (Rodriguez-Martin et al., 2017). By combining these two
382 sources of information, ChimPipe achieves better precision than methods that
383 depend on a single type of signal.

384 The pipeline works with many eukaryotic species that have available genome
385 and annotation data, making it a versatile tool for studying chimera evolution
386 and transcriptome structure (Rodriguez-Martin et al., 2017). It can also predict
387 multiple isoforms for each gene pair and identify breakpoint coordinates that are
388 useful for reconstructing and verifying chimeric transcripts. Tests using both
389 simulated and real datasets have shown that ChimPipe maintains high accuracy
390 and reliable performance.

391 ChimPipe’s modular design lets users adjust parameters to fit different se-
392 quencing protocols or organism characteristics. Experimental results have con-
393 firmed that many chimeric transcripts detected by the tool correspond to func-

394 tional fusion proteins, showing its value for understanding chimera biology and
395 its potential applications in disease research (Rodriguez-Martin et al., 2017).

396 **2.4 Machine Learning Approaches for Chimera** 397 **and Sequence Quality Detection**

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

409 **2.4.1 Feature-Based Representations of Genomic Se-** 410 **quences**

411 In genomic analysis, feature extraction converts DNA sequences into numerical
412 representations suitable for ML algorithms. A common approach is k-mer fre-
413 quency analysis, where normalized k-mer counts form the feature vector (Vervier,

2015). These features effectively capture local compositional patterns that often differ between authentic and chimeric reads. In particular, deviations in k-mer profiles between adjacent read segments can serve as a compositional signature of template-switching events. Additional descriptors such as GC content and sequence entropy can further distinguish sequence types; in metagenomic classification and virus detection, k-mer-based features have shown strong performance and robustness to noise (Ren et al., 2020; Vervier, 2015). For chimera detection specifically, abrupt shifts in GC or k-mer composition along a read can indicate junctions between parental fragments. Windowed feature extraction enables models to capture these discontinuities that rule-based algorithms may overlook.

Machine learning models can also leverage alignment-derived features such as the frequency of split alignments, variation in mapping quality, and local coverage irregularities. Split reads and discordant read pairs are classical signatures of genomic junctions and have been formalized in probabilistic frameworks for structural-variant discovery that integrate multiple evidence types (Layer, Hall, & Quinlan, 2014). Similarly, long-read tools such as Sniffles employ split-alignment and coverage anomalies to accurately localize breakpoints (Sedlazeck et al., 2018). Modern aligners such as Minimap2 (Li, 2018) output supplementary (SA tags) and secondary alignments as well as chaining and alignment-score statistics that can be summarized into quantitative predictors for machine-learning models. These alignment-signal features are particularly relevant to PCR-induced mitochondrial chimeras, where template-switching events produce reads partially matching distinct regions of the same or related genomes. Integrating such cues within a supervised-learning framework enables artifact detection even in datasets lacking complete or perfectly assembled references.

439 A further biologically grounded descriptor is micro-homology length at puta-
440 tive junctions. Micro-homology refers to short, shared sequences (often in the
441 range of a few to tens of base pairs) that are near breakpoints and mediate
442 non-canonical repair or template-switch mechanisms. Studies of double strand
443 break repair and structural variation have demonstrated that the length of micro-
444 homology correlates with the likelihood of micro-homology-mediated end joining
445 (MMEJ) or fork-stalled template-switching pathways (Sfeir & Symington, 2015).
446 In the context of PCR-induced chimeras, template switching during amplifica-
447 tion often leaves short identical sequences at the junction of two concatenated
448 fragments. Quantifying the longest exact suffix–prefix overlap at each candidate
449 breakpoint thus provides a mechanistic signature of chimerism and complements
450 both compositional (k-mer) and alignment (SA count) features.

451 Chapter 3

452 Research Methodology

453 This chapter outlines and explains the specific steps and activities to be carried
454 out in completing the project.

455 3.1 Research Activities

456 As illustrated in Figure 3.1, the researchers will carry out a sequence of compu-
457 tational procedures designed to detect PCR-induced chimeric reads in mitochon-
458 drial genomes. The process begins with the collection of mitochondrial reference
459 sequences from the NCBI database, which will serve as the foundation for gener-
460 ating simulated chimeric reads. These datasets will then undergo bioinformatics
461 pipeline development, which includes alignment, k-mer extraction, and homology-
462 based filtering to prepare the data for model construction. The machine-learning
463 model will subsequently be trained and tested using the processed datasets to
464 assess its accuracy and reliability. Depending on the evaluation results, the model

465 will either be refined and retrained to improve performance or, if the metrics meet
 466 the desired threshold, deployed for further validation and application.

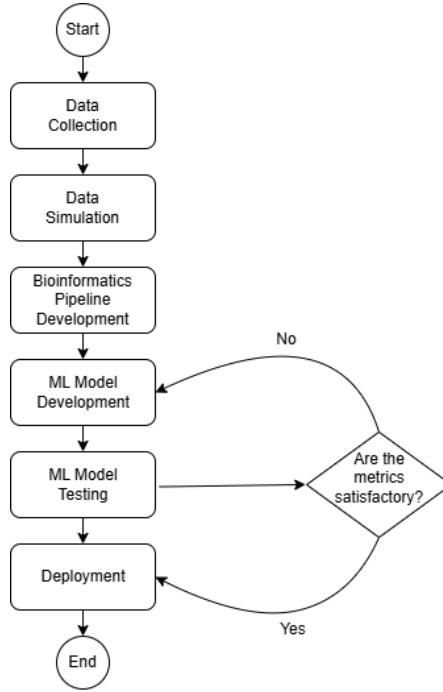


Figure 3.1: Process Diagram of Special Project

467 3.1.1 Data Collection

468 The researchers will collect mitochondrial genome reference sequences of *Sar-*
 469 *dinella lemuru* from the National Center for Biotechnology Information (NCBI)
 470 database. The downloaded files will be in FASTA format to ensure compatibility
 471 with bioinformatics tools and subsequent analysis. The gathered sequences will
 472 serve as the basis for generating simulated chimeric reads to be used in model
 473 development.

474 The expected outcome of this process is a comprehensive dataset of *Sardinella*

475 *lemuru* mitochondrial reference sequences that will serve as the foundation for
476 the succeeding stages of the study. This step is scheduled to start in the first
477 week of November 2025 and is expected to be completed by the last week of
478 November 2025, with a total duration of approximately one (1) month.

479 **3.1.2 Data Simulation**

480 The researchers will simulate sequencing data using the reference sequences col-
481 lected from NCBI. Using `wgsim`, a total of 5,000 paired-end reads (R1 and R2)
482 will be generated from the reference genome and designated as clean reads. These
483 reads will be saved in FASTQ (`.fastq`) format. From the same reference, a Bash
484 script will be created to deliberately cut and reconnect portions of the sequence,
485 introducing artificial junctions that mimic chimeric regions. The manipulated
486 reference file, saved in FASTA (`.fasta`) format, will then be processed in `wgsim`
487 to simulate an additional 5,000 paired-end chimeric reads, also stored in FASTQ
488 (`.fastq`) format. The resulting read files will be aligned to the original reference
489 genome using SAMtools, generating SAM (`.sam`) or BAM (`.bam`) alignment files.
490 During this alignment process, clean reads will be labeled as “0,” while chimeric
491 reads will be labeled as “1” in a corresponding CSV (`.tsv`) file.

492 The expected outcome of this process is a complete set of clean and chimeric
493 paired-end reads prepared for subsequent analysis and model development. This
494 step is scheduled to start in the first week of November 2025 and is expected
495 to be completed by the last week of November 2025, with a total duration of
496 approximately one (1) month.

497 **3.1.3 Bioinformatics Tools Pipeline**

498 The researchers will obtain the necessary analytical features through the devel-
499 opment and implementation of a bioinformatics pipeline. This pipeline will serve
500 as a reproducible and modular workflow that accepts FASTQ and BAM inputs,
501 processes these through a series of analytical stages, and outputs tabular feature
502 matrices (TSV) for downstream machine learning. All scripts will be version-
503 controlled through GitHub, and computational environments will be standardized
504 using Conda to ensure cross-platform reproducibility. To promote transparency
505 and replicability, the exact software versions, parameters, and command-line ar-
506 guments used in each stage will be documented. To further ensure correctness
507 and adherence to best practices, the researchers will consult with bioinformatics
508 experts in Philippine Genome Center Visayas for validation of pipeline design,
509 feature extraction logic, and overall data integrity. This stage of the study is
510 scheduled to begin in the last week of November 2025 and conclude by the last
511 week of January 2026, with an estimated total duration of approximately two (2)
512 months.

513 The bioinformatics pipeline focuses on three principal features from the sim-
514 ulated and aligned sequencing data: (1) supplementary alignment count (SA
515 count), (2) k-mer composition difference between read segments, and (3) micro-
516 homology length at potential junctions. Each of these features captures a distinct
517 biological or computational signature associated with PCR-induced chimeras.

518 Alignment and Supplementary Alignment Count

519 This will be derived through sequence alignment using Minimap2, with subsequent
520 processing performed using SAMtools and `pysam` in Python. Sequencing reads
521 will be aligned to the *Sardinella lemuru* mitochondrial reference genome using
522 Minimap2 with the `-ax sr` preset (optimized for short reads). The output will
523 be converted and sorted using SAMtools, producing an indexed BAM file which
524 will be parsed using `pysam` to count the number of supplementary alignments
525 (SA tags) per read. Each read's mapping quality, number of split segments,
526 and alignment characteristics will be recorded in a corresponding TSV file. The
527 presence of multiple alignment loci within a single read, as reflected by a nonzero
528 SA count, serves as direct computational evidence of chimerism. Reads that
529 contain supplementary alignments or soft-clipped regions are strong candidates
530 for chimeric artifacts arising from PCR template switching or improper assembly
531 during sequencing.

532 K-mer Composition Difference

533 Chimeric reads often comprise fragments from distinct genomic regions, resulting
534 in a compositional discontinuity between segments. Comparing k-mer frequency
535 profiles between the left and right halves of a read allows detection of such abrupt
536 compositional shifts, independent of alignment information. This will be obtained
537 using Jellyfish, a fast k-mer counting software. For each read, the sequence will
538 be divided into two segments, either at the midpoint or at empirically determined
539 breakpoints inferred from supplementary alignment data, to generate left and right
540 sequence segments. Jellyfish will then compute k-mer frequency profiles (with $k =$

541 5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized
542 and compared using distance metrics such as cosine similarity or Jensen–Shannon
543 divergence to quantify compositional disparity between the two halves of the same
544 read. The resulting difference scores will be stored in a structured TSV file.

545 **Micro-homology Length**

546 The micro-homology length will be computed using a custom Python script that
547 detects the longest exact suffix–prefix overlap within ± 30 base pairs surround-
548 ing a candidate breakpoint. This analysis identifies the number of consecutive
549 bases shared between the end of one segment and the beginning of another. The
550 presence and length of such micro-homology are classic molecular signatures of
551 PCR-induced template switching, where short identical regions (typically 3–15
552 base pairs) promote premature termination and recombination of DNA synthesis
553 on a different template strand. By quantifying micro-homology, the researchers
554 can assess whether the suspected breakpoint exhibits characteristics consistent
555 with PCR artifacts rather than true biological variants. Each read will therefore
556 be annotated with its corresponding micro-homology length, overlap sequence,
557 and GC content.

558 After extracting the three primary features, all resulting TSV files will be
559 joined using the read identifier as a common key to generate a unified feature ma-
560 trix. Additional read-level metadata such as read length, mean base quality, and
561 number of clipped bases will also be included to provide contextual information.
562 This consolidated dataset will serve as the input for subsequent machine-learning
563 model development and evaluation.

564 3.1.4 Machine-Learning Model Development

565 The classification component of MitoChime will employ two ensemble algo-
566 rithms—Random Forest (RF) and Extreme Gradient Boosting (XGBoost)—to
567 evaluate complementary learning paradigms. Random Forest applies bootstrap
568 aggregation (bagging) to reduce model variance and improve stability, whereas
569 XGBoost implements gradient boosting to minimize bias and capture complex
570 non-linear relationships among genomic features. Using both models enables a
571 balanced assessment of predictive performance and interpretability.

572 The dataset will be divided into training (80%) and testing (20%) subsets.
573 The training data will be used for model fitting and hyperparameter optimization
574 through five-fold cross-validation, in which the data are partitioned into five folds;
575 four folds are used for training and one for validation in each iteration. Perfor-
576 mance metrics will be averaged across folds, and the optimal parameters will be
577 selected based on mean cross-validation accuracy. The final models will then be
578 evaluated on the held-out test set to obtain unbiased performance estimates.

579 Model development and evaluation will be implemented in Python (ver-
580 sion 3.11) using the `scikit-learn` and `xgboost` libraries. Standard metrics
581 including accuracy, precision, recall, F1-score, and area under the ROC curve
582 (AUC) will be computed to quantify predictive performance. Feature-importance
583 analyses will be performed to identify the most discriminative variables contribut-
584 ing to chimera detection.

585 **3.1.5 Validation and Testing**

586 Validation will involve both internal and external evaluations. Internal validation
587 will be achieved through five-fold cross-validation on the training data to verify
588 model generalization and reduce variance due to random sampling. External
589 validation will be achieved through testing on the 20% hold-out dataset derived
590 from the simulated reads, which will serve as an unbiased benchmark to evaluate
591 how well the trained models generalize to unseen data. All feature extraction and
592 preprocessing steps will be performed using the same bioinformatics pipeline to
593 ensure consistency and comparability across validation stages.

594 Comparative evaluation between the Random Forest and XGBoost classifiers
595 will establish which model achieves superior predictive accuracy and computa-
596 tional efficiency under identical data conditions.

597 **3.1.6 Documentation**

598 Comprehensive documentation will be maintained throughout the study to en-
599 sure transparency, reproducibility, and scientific integrity. All stages of the re-
600 search—including data acquisition, preprocessing, feature extraction, model train-
601 ing, and validation—will be systematically recorded. For each analytical step, the
602 corresponding parameters, software versions, and command-line scripts will be
603 documented to enable exact replication of results.

604 Version control and collaborative management will be implemented through
605 GitHub, which will serve as the central repository for all project files, including
606 Python scripts, configuration settings, and Jupyter notebooks. The repository

607 structure will follow standard research data management practices, with clear
 608 directories for datasets, processed outputs, and analysis scripts. Changes will be
 609 tracked through commit histories to ensure traceability and accountability.

610 Computational environments will be standardized using Conda, with environ-
 611 ment files specifying dependencies and package versions to maintain consistency
 612 across systems. Experimental workflows and exploratory analyses will be con-
 613 ducted in Jupyter Notebooks, which facilitate real-time visualization, annotation,
 614 and incremental testing of results.

615 For the preparation of the final manuscript and supplementary materials,
 616 Overleaf (LaTeX) will be utilized to produce publication-quality formatting, con-
 617 sistent referencing, and reproducible document compilation. The documentation
 618 process will also include a project timeline outlining major milestones such as
 619 data collection, simulation, feature extraction, model evaluation, and reporting to
 620 ensure systematic progress and adherence to the research schedule.

621 3.2 Calendar of Activities

622 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
 623 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	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •

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
- Bi, C., Shen, F., Han, F., Qu, Y., Hou, J., Xu, K., ... Yin, T. (2024, 01).
Pmat: an efficient plant mitogenome assembly toolkit using low-coverage
hifi sequencing data. *Horticulture Research*, *11*(3), uhae023. Retrieved
from <https://doi.org/10.1093/hr/uhae023> doi: 10.1093/hr/uhae023
- 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*, 45(4), e18. doi: 10.1093/nar/gkw955

Edgar, R. C. (2016). Uchime2: improved chimera prediction for amplicon sequencing. *bioRxiv*. Retrieved from <https://api.semanticscholar.org/CorpusID:88955007>

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

Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011). Uchime improves sensitivity and speed of chimera detection. *Bioinformatics*, 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

Frenkel-Morgenstern, M., Lacroix, V., Ezkurdia, I., Levin, Y., Gabashvili, A., Prilusky, J., ... Valencia, A. (2012, 05). Chimeras taking shape: Potential functions of proteins encoded by chimeric rna transcripts. *Genome research*, 22, 1231-42. doi: 10.1101/gr.130062.111

Glenn, T. C. (2011). Field guide to next-generation dna sequencers. *Molecular Ecology Resources*, 11(5), 759–769. doi: 10.1111/j.1755-0998.2011.03024.x

Gonzalez, J. M., Zimmermann, J., & Saiz-Jimenez, C. (2004, 09). Evaluating putative chimeric sequences from pcr-amplified products. *Bioinformatics*, 21(3), 333-337. Retrieved from <https://doi.org/10.1093/bioinformatics/bti008> doi: 10.1093/bioinformatics/bti008

Gray, M. W. (2012). Mitochondrial evolution. *Cold Spring Harbor perspectives in biology*, 4. Retrieved from <https://doi.org/10.1101/cshperspect.a011403> doi: 10.1101/cshperspect.a011403

669 Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial
670 genomes directly from genomic next-generation sequencing reads—a baiting
671 and iterative mapping approach. *Nucleic Acids Research*, *41*(13), e129. doi:
672 10.1093/nar/gkt371

673 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,
674 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de
675 novo assembly of organelle genomes. *Genome Biology*, *21*(1), 241. doi:
676 10.1186/s13059-020-02154-5

677 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
678 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
679 1819–1825. doi: 10.1093/nar/26.7.1819

680 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
681 (2021). Mitochondrial dna reveals genetically structured haplogroups of
682 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
683 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

684 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework
685 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-
686 -15-6-r84

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

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

694 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*

695 *Reviews Genetics*, 11(1), 31–46. doi: 10.1038/nrg2626

696 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
 697 an ensemble classifier for chimera detection in 16s rna sequencing stud-
 698 ies. *Applied and Environmental Microbiology*, 81(5), 1573-1584. Retrieved
 699 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
 700 10.1128/AEM.02896-14

701 Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J.
 702 (2023). Effects of error, chimera, bias, and gc content on the accuracy of
 703 amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from [https://](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)
 704 journals.asm.org/doi/abs/10.1128/msystems.01025-23 doi: 10.1128/
 705 msystems.01025-23

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

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

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

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

720 Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler,

721 A., & Schatz, M. (2018, 06). Accurate detection of complex structural
 722 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10
 723 .1038/s41592-018-0001-7
 724 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
 725 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
 726 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
 727 [.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)
 728 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)
 729 Vervier, M. P. T. M. V. J. B. . V. J. P., K. (2015). Large-scale machine learning
 730 for metagenomics sequence classification. *Bioinformatics*, 32, 1023 - 1032.
 731 Retrieved from <https://api.semanticscholar.org/CorpusID:9863600>
 732 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*
 733 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI
 734 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)
 735 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)