

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

20

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 Research Methodology	7
30	2.1 Research Activities	7
31	2.1.1 Data Collection	8
32	2.1.2 Data Simulation	9
33	2.1.3 Bioinformatics Tools Pipeline	10
34	2.1.4 Machine-Learning Model Development	13
35	2.1.5 Validation and Testing	14
36	2.1.6 Documentation	14
37	2.2 Calendar of Activities	15

38 List of Figures

<small>39</small>	2.1 Process Diagram of Special Project	8
-------------------	--	---

⁴⁰ List of Tables

⁴¹	2.1 Timetable of Activities	16
---------------	---------------------------------------	----

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

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

83 1.2 Problem Statement

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

102 1.3 Research Objectives

103 1.3.1 General Objective

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

107 1.3.2 Specific Objectives

108 Specifically, the researchers aim to:

- 109 1. Construct simulated and empirical *Sardinella* Illumina paired-end datasets
110 containing both clean and PCR-induced chimeric reads.
- 111 2. Extract alignment- and sequence-based features (e.g., k-mer composition,
112 junction complexity, split-alignment counts) from both clean and chimeric
113 reads.
- 114 3. Train, validate, and compare supervised machine-learning models (e.g., Ran-
115 dom Forest, XGBoost) for classifying reads as clean or chimeric.
- 116 4. Determine feature importance and identify the most informative indicators
117 of PCR-induced chimerism.
- 118 5. Integrate the optimized classifier into a modular and interpretable pipeline
119 deployable on standard computing environments at PGC Visayas.

120 1.4 Scope and Limitations of the Research

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

133 1.5 Significance of the Research

134 This research provides both methodological and practical contributions to mi-
135 tochondrial genomics and bioinformatics. First, MitoChime enhances assembly
136 accuracy by filtering PCR-induced chimeras prior to genome assembly, thereby
137 improving the contiguity and correctness of *Sardinella* mitochondrial genomes.
138 Second, it promotes automation and reproducibility by replacing subjective man-
139 ual curation with a data-driven, machine-learning-based workflow. Third, the
140 pipeline demonstrates computational efficiency through its design, enabling im-

141 plementation on modest computing infrastructures commonly available in regional
142 laboratories. Beyond technical improvements, MitoChime contributes to local ca-
143 pacity building by strengthening expertise in bioinformatics and machine-learning
144 integration, aligning with the mission of the Philippine Genome Center Visayas.
145 Finally, accurate mitochondrial assemblies are vital for fisheries management,
146 population genetics, and biodiversity conservation, providing reliable genomic re-
147 sources for species such as *Sardinella*. Through these contributions, MitoChime
148 advances the reliability of mitochondrial genome reconstruction and supports sus-
149 tainable, data-driven research in Philippine genomics.

150 Chapter 2

151 Research Methodology

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

154 2.1 Research Activities

155 As illustrated in Figure 2.1, the researchers will carry out a sequence of compu-
156 tational procedures designed to detect PCR-induced chimeric reads in mitochon-
157 drial genomes. The process begins with the collection of mitochondrial reference
158 sequences from the NCBI database, which will serve as the foundation for gener-
159 ating simulated chimeric reads. These datasets will then undergo bioinformatics
160 pipeline development, which includes alignment, k-mer extraction, and homology-
161 based filtering to prepare the data for model construction. The machine-learning
162 model will subsequently be trained and tested using the processed datasets to
163 assess its accuracy and reliability. Depending on the evaluation results, the model

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

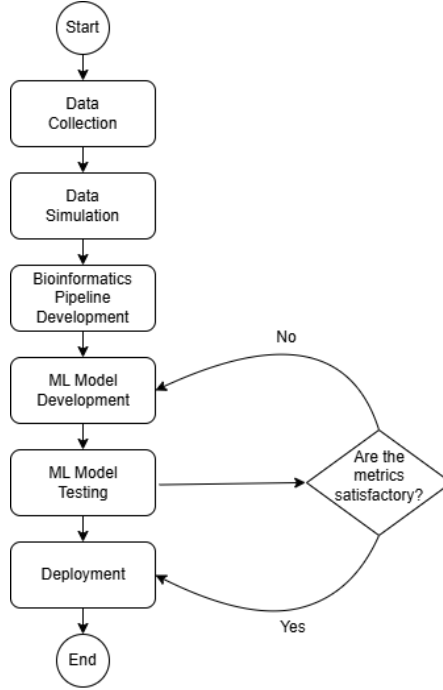


Figure 2.1: Process Diagram of Special Project

166 2.1.1 Data Collection

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

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

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

178 **2.1.2 Data Simulation**

179 The researchers will simulate sequencing data using the reference sequences col-
180 lected from NCBI. Using `wgsim`, a total of 5,000 paired-end reads (R1 and R2)
181 will be generated from the reference genome and designated as clean reads. These
182 reads will be saved in FASTQ (`.fastq`) format. From the same reference, a Bash
183 script will be created to deliberately cut and reconnect portions of the sequence,
184 introducing artificial junctions that mimic chimeric regions. The manipulated
185 reference file, saved in FASTA (`.fasta`) format, will then be processed in `wgsim`
186 to simulate an additional 5,000 paired-end chimeric reads, also stored in FASTQ
187 (`.fastq`) format. The resulting read files will be aligned to the original reference
188 genome using SAMtools, generating SAM (`.sam`) or BAM (`.bam`) alignment files.
189 During this alignment process, clean reads will be labeled as “0,” while chimeric
190 reads will be labeled as “1” in a corresponding CSV (`.csv`) file.

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

196 2.1.3 Bioinformatics Tools Pipeline

197 The researchers will obtain the necessary analytical features through the devel-
198 opment and implementation of a bioinformatics pipeline. This pipeline will serve
199 as a reproducible and modular workflow that accepts FASTQ and BAM inputs,
200 processes these through a series of analytical stages, and outputs tabular feature
201 matrices (TSV/CSV) for downstream machine learning. All scripts will be version-
202 controlled through GitHub, and computational environments will be standardized
203 using Conda to ensure cross-platform reproducibility. To promote transparency
204 and replicability, the exact software versions, parameters, and command-line ar-
205 guments used in each stage will be documented. To further ensure correctness
206 and adherence to best practices, the researchers will consult with bioinformatics
207 experts in Philippine Genome Center Visayas for validation of pipeline design,
208 feature extraction logic, and overall data integrity. This stage of the study is
209 scheduled to begin in the last week of November 2025 and conclude by the last
210 week of January 2026, with an estimated total duration of approximately two (2)
211 months.

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

217 Alignment and Supplementary Alignment Count

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

231 K-mer Composition Difference

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

240 5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized
241 and compared using distance metrics such as cosine similarity or Jensen–Shannon
242 divergence to quantify compositional disparity between the two halves of the same
243 read. The resulting difference scores will be stored in a structured TSV file.

244 **Micro-homology Length**

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

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

2.1.4 Machine-Learning Model Development

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

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

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

284 **2.1.5 Validation and Testing**

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

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

296 **2.1.6 Documentation**

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

303 Version control and collaborative management will be implemented through
304 GitHub, which will serve as the central repository for all project files, including

305 Python scripts, configuration settings, and Jupyter notebooks. The repository
306 structure will follow standard research data management practices, with clear
307 directories for datasets, processed outputs, and analysis scripts. Changes will be
308 tracked through commit histories to ensure traceability and accountability.

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

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

320 **2.2 Calendar of Activities**

321 Table 2.1 presents the project timeline in the form of a Gantt chart, where each
322 bullet point corresponds to approximately one week of planned activity.

Table 2.1: 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

- 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*, 45(4), e18. doi: 10.1093/nar/gkw955
- 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
- 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
- Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial

342 genomes directly from genomic next-generation sequencing reads—a baiting
 343 and iterative mapping approach. *Nucleic Acids Research*, *41*(13), e129. doi:
 344 10.1093/nar/gkt371

345 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,
 346 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de
 347 novo assembly of organelle genomes. *Genome Biology*, *21*(1), 241. doi:
 348 10.1186/s13059-020-02154-5

349 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
 350 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
 351 1819–1825. doi: 10.1093/nar/26.7.1819

352 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
 353 (2021). Mitochondrial dna reveals genetically structured haplogroups of
 354 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
 355 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

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

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

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

365 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*
 366 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI
 367 Technical Paper Series. Retrieved from <https://nfrdi.da.gov.ph/tpjf/>

etc/Willette%20et%20al.%20Sardines%20Review.pdf