

1 MITOCHIME: A MACHINE-LEARNING
2 PIPELINE FOR DETECTING PCR-INDUCED
3 CHIMERAS IN MITOCHONDRIAL ILLUMINA
4 READS

5 A Special Project Proposal
6 Presented to
7 the Faculty of the Division of Physical Sciences
8 and Mathematics
9 College of Arts and Sciences
10 University of the Philippines Visayas
11 Miag-ao, Iloilo

12 In Partial Fulfillment

13

of the Requirements for the Degree of

14

Bachelor of Science in Computer Science

15

by

16

DURAN, Duranne

17

LIN, Yvonne

18

PAILDEN, Daniella

19

Adviser

20

Francis DIMZON

21

October 24, 2025

22

Contents

23	0.1	Introduction	1
24	0.1.1	Overview	1
25	0.1.2	Problem Statement	2
26	0.1.3	Research Objectives	3
27	0.1.4	Scope and Limitations of the Research	4
28	0.1.5	Significance of the Research	5
29	0.2	Research Methodology	6
30	0.2.1	Research Activities	6
31	0.2.2	Calendar of Activities	14

32 List of Figures

<small>33</small>	1	Process diagram of the special project.	7
-------------------	---	---	---

³⁴ List of Tables

³⁵	1	Timetable of Activities	14
---------------	---	-----------------------------------	----

Chapter 1

0.1 Introduction

0.1.1 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

58 such artifacts (Hahn et al., 2013; Jin et al., 2020). Consequently, undetected
59 chimeras may produce fragmented assemblies or misidentified organellar bound-
60 aries. To ensure accurate reconstruction of mitochondrial genomes, a reliable
61 and automated method for detecting and filtering PCR-induced chimeras before
62 assembly is essential.

63 This study focuses on mitochondrial sequencing data from the genus *Sar-*
64 *dinella*, a group of small pelagic fishes widely distributed in Philippine waters.
65 Among them, *Sardinella lemuru* (Bali sardinella) is one of the country’s most
66 abundant and economically important species, providing protein and livelihood
67 to coastal communities (Labrador, Agmata, Palermo, Ravago-Gotanco, & Pante,
68 2021; Willette, Bognot, Mutia, & Santos, 2011). Accurate mitochondrial assem-
69 blies are critical for understanding its population genetics, stock structure, and
70 evolutionary history. However, assembly pipelines often encounter errors or fail
71 to complete due to undetected chimeric reads. To address this gap, this research
72 introduces **MitoChime**, a machine-learning pipeline designed to detect and filter
73 PCR-induced chimeric reads using both alignment- and sequence-derived statisti-
74 cal features. The tool aims to provide bioinformatics laboratories, particularly the
75 Philippine Genome Center Visayas, with an efficient, interpretable, and resource-
76 optimized solution for improving mitochondrial genome reconstruction.

77 **0.1.2 Problem Statement**

78 While NGS technologies have revolutionized genomic data acquisition, the ac-
79 curacy of mitochondrial genome assembly remains limited by artifacts produced
80 during PCR amplification. These chimeric reads can distort assembly graphs and

81 cause misassemblies, with especially severe effects in small, circular mitochon-
82 drial genomesBoore (1999); Cameron (2014). Existing assembly pipelines such as
83 GetOrganelle, MITObim, and NOVOPlasty assume that sequencing reads are free
84 of such artifactsDierckxsens et al. (2017); Hahn et al. (2013); Jin et al. (2020).
85 At the Philippine Genome Center Visayas, several mitochondrial assemblies have
86 failed or yielded incomplete contigs despite sufficient coverage, suggesting that
87 undetected chimeric reads compromise assembly reliability. Meanwhile, existing
88 chimera-detection tools such as UCHIME and VSEARCH were developed primar-
89 ily for amplicon-based microbial community analysis and rely heavily on reference
90 or taxonomic comparisonsEdgar, Haas, Clemente, Quince, and Knight (2011);
91 Rognes, Flouris, Nichols, Quince, and Mahé (2016). These approaches are un-
92 suitable for single-species organellar data, where complete reference genomes are
93 often unavailable. Therefore, there is a pressing need for a reference-independent,
94 data-driven tool capable of automatically detecting and filtering PCR-induced
95 chimeras in mitochondrial sequencing datasets.

96 **0.1.3 Research Objectives**

97 **General Objective**

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

101 **Specific Objectives**

102 Specifically, the researchers aim to:

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

114 **0.1.4 Scope and Limitations of the Research**

115 This study focuses on detecting PCR-induced chimeric reads in Illumina paired-
116 end mitochondrial sequencing data from *Sardinella* species. The work emphasizes
117 **wgsim** simulations and selected empirical data obtained from open-access genomic
118 repositories such as the National Center for Biotechnology Information (NCBI).
119 The study excludes naturally occurring chimeras, nuclear mitochondrial pseudo-
120 genes (NUMTs), and large-scale structural rearrangements in nuclear genomes.

121 Feature extraction prioritizes interpretable, shallow statistics and alignment met-
122 rics rather than deep-learning embeddings to ensure transparency and computa-
123 tional efficiency. Testing on long-read platforms (e.g., Nanopore, PacBio) and
124 other taxa lies beyond the project’s scope. The resulting pipeline will serve as a
125 foundation for future, broader chimera-detection frameworks applicable to diverse
126 organellar genomes.

127 **0.1.5 Significance of the Research**

128 This research provides both methodological and practical contributions to mi-
129 tochondrial genomics and bioinformatics. First, MitoChime enhances assembly
130 accuracy by filtering PCR-induced chimeras prior to genome assembly, thereby
131 improving the contiguity and correctness of *Sardinella* mitochondrial genomes.
132 Second, it promotes automation and reproducibility by replacing subjective man-
133 ual curation with a data-driven, machine-learning-based workflow. Third, the
134 pipeline demonstrates computational efficiency through its design, enabling im-
135 plementation on modest computing infrastructures commonly available in regional
136 laboratories. Beyond technical improvements, MitoChime contributes to local ca-
137 pacity building by strengthening expertise in bioinformatics and machine-learning
138 integration, aligning with the mission of the Philippine Genome Center Visayas.
139 Finally, accurate mitochondrial assemblies are vital for fisheries management,
140 population genetics, and biodiversity conservation, providing reliable genomic re-
141 sources for species such as *Sardinella*. Through these contributions, MitoChime
142 advances the reliability of mitochondrial genome reconstruction and supports sus-
143 tainable, data-driven research in Philippine genomics.

144 Chapter 3

145 0.2 Research Methodology

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

148 0.2.1 Research Activities

149 As illustrated in Figure 1, the researchers will carry out a sequence of compu-
150 tational procedures designed to detect PCR-induced chimeric reads in mitochon-
151 drial genomes. The process begins with the collection of mitochondrial reference
152 sequences from the NCBI database, which will serve as the foundation for gener-
153 ating simulated chimeric reads. These datasets will then undergo bioinformatics
154 pipeline development, which includes alignment, k-mer extraction, and homology-
155 based filtering to prepare the data for model construction. The machine-learning
156 model will subsequently be trained and tested using the processed datasets to
157 assess its accuracy and reliability. Depending on the evaluation results, the model
158 will either be refined and retrained to improve performance or, if the metrics meet
159 the desired threshold, deployed for further validation and application.

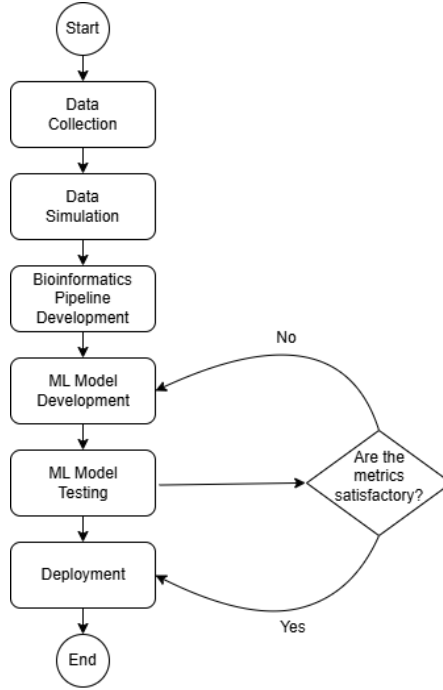


Figure 1: Process diagram of the special project.

160 Data Collection

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

167 The expected outcome of this process is a comprehensive dataset of *Sardinella*
168 *lemuru* mitochondrial reference sequences that will serve as the foundation for
169 the succeeding stages of the study. This step is scheduled to start in the first
170 week of November 2025 and is expected to be completed by the last week of

171 November 2025, with a total duration of approximately one (1) month.

172 **Data Simulation**

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

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

190 **Bioinformatics Tools Pipeline**

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

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

211 **Alignment and Supplementary Alignment Count**

212 This will be derived through sequence alignment using Minimap2, with subsequent

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

225 **K-mer Composition Difference**

226 Chimeric reads often comprise fragments from distinct genomic regions, resulting
227 in a compositional discontinuity between segments. Comparing k-mer frequency
228 profiles between the left and right halves of a read allows detection of such abrupt
229 compositional shifts, independent of alignment information. This will be obtained
230 using Jellyfish, a fast k-mer counting software. For each read, the sequence will
231 be divided into two segments, either at the midpoint or at empirically determined
232 breakpoints inferred from supplementary alignment data, to generate left and right
233 sequence segments. Jellyfish will then compute k-mer frequency profiles (with $k =$
234 5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized
235 and compared using distance metrics such as cosine similarity or Jensen–Shannon
236 divergence to quantify compositional disparity between the two halves of the same

237 read. The resulting difference scores will be stored in a structured TSV file.

238 **Micro-homology Length**

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

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

257 Machine-Learning Model Development

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

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

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

278 **Validation and Testing**

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

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

290 **Documentation**

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

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

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

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

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

314 0.2.2 Calendar of Activities

315 Table 1 presents the project timeline in the form of a Gantt chart, where each
 316 bullet point corresponds to approximately one week of planned activity.

Table 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

336 genomes directly from genomic next-generation sequencing reads—a baiting
 337 and iterative mapping approach. *Nucleic Acids Research*, *41*(13), e129. doi:
 338 10.1093/nar/gkt371

339 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,
 340 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de
 341 novo assembly of organelle genomes. *Genome Biology*, *21*(1), 241. doi:
 342 10.1186/s13059-020-02154-5

343 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
 344 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
 345 1819–1825. doi: 10.1093/nar/26.7.1819

346 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
 347 (2021). Mitochondrial dna reveals genetically structured haplogroups of
 348 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
 349 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

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

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

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

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

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