

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

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

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

13 by

14 Duranne Duran
15 Yvonne Lin
16 Daniella Pailden

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

19 December 5, 2025

Contents

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

32	2.2	PCR Amplification and Chimera Formation	9
33	2.3	Existing Traditional Approaches for Chimera Detection	10
34	2.3.1	UCHIME	11
35	2.3.2	UCHIME2	13
36	2.3.3	CATch	14
37	2.3.4	ChimPipe	15
38	2.4	Machine Learning Approaches for Chimera and Sequence Quality	
39		Detection	16
40	2.4.1	Feature-Based Representations of Genomic Sequences . . .	17
41	2.5	Synthesis of Chimera Detection Approaches	18
42	3	Research Methodology	22
43	3.1	Research Activities	22
44	3.1.1	Data Collection	23
45	3.1.2	Bioinformatics Tools Pipeline	27
46	3.1.3	Machine Learning Model Development	30
47	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
48		Evaluation	31
49	3.1.5	Feature Importance and Interpretation	32

50	3.1.6 Validation and Testing	33
51	3.1.7 Documentation	34
52	3.2 Calendar of Activities	35

⁵³ List of Figures

⁵⁴	3.1 Process Diagram of Special Project	23
---------------	--	----

55 List of Tables

56	2.1 Summary of Existing Methods and Research Gaps	20
57	3.1 Timetable of Activities	35

Chapter 1

Introduction

1.1 Overview

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

PCR chimeras form when incomplete extension products from one template

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

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

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

1.2 Problem Statement

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

117 1.3 Research Objectives

118 1.3.1 General Objective

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

123 1.3.2 Specific Objectives

124 Specifically, the study aims to:

- 125 1. construct simulated *Sardinella lemuru* Illumina paired-end datasets contain-
126 ing both clean and PCR-induced chimeric reads,
- 127 2. extract alignment-based and sequence-based features such as k-mer compo-
128 sition, junction complexity, and split-alignment counts from both clean and
129 chimeric reads,
- 130 3. train, validate, and compare supervised machine-learning models for classi-
131 fying reads as clean or chimeric,
- 132 4. determine feature importance and identify indicators of PCR-induced
133 chimerism,
- 134 5. integrate the optimized classifier into a modular and interpretable pipeline
135 deployable on standard computing environments at PGC Visayas.

1.4 Scope and Limitations of the Research

This study focuses on detecting PCR-induced chimeric reads in Illumina paired-end mitochondrial sequencing data from *Sardinella lemuru*. The decision to restrict the taxonomic scope to a single species is based on four considerations: (1) to limit interspecific variation in mitochondrial genome size, GC content, and repetitive regions so that differences in read patterns can be attributed more directly to PCR-induced chimerism; (2) to align the analysis with relevant *S. lemuru* sequencing projects at PGC Visayas; (3) to take advantage of the availability of *S. lemuru* mitochondrial assemblies and raw datasets in public repositories such as the National Center for Biotechnology Information (NCBI), which facilitates reference selection and benchmarking; and (4) to develop a tool that directly supports local studies on *S. lemuru* population structure and fisheries management.

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

159 1.5 Significance of the Research

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

170 Chapter 2

171 Review of Related Literature

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

178 2.1 The Mitochondrial Genome

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

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

189 **2.1.1 Mitochondrial Genome Assembly**

190 Mitochondrial genome assembly refers to the reconstruction of the complete mito-
191 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is
192 conducted to obtain high-quality, continuous representations of the mitochondrial
193 genome that can be used for a wide range of analyses, including species identi-
194 fication, phylogenetic reconstruction, evolutionary studies, and investigations of
195 mitochondrial diseases. Because mtDNA evolves rapidly, its assembled sequence
196 provides valuable insights into population structure, lineage divergence, and adap-
197 tive evolution across taxa (Boore, 1999). Compared to nuclear genome assembly,
198 assembling the mitochondrial genome is often considered more straightforward but
199 still encounters technical challenges such as the formation of chimeric reads. Com-
200 monly used tools for mitogenome assembly such as GetOrganelle and MITObim
201 operate under the assumption of organelle genome circularity, and are vulnerable
202 when chimeric reads disrupt this circular structure, resulting in assembly errors
203 (Hahn et al., 2013; Jin et al., 2020).

2.2 PCR Amplification and Chimera Formation

PCR plays an important role in NGS library preparation, as it amplifies target DNA fragments for downstream analysis. However as previously mentioned, the amplification process can also introduce chimeric reads which compromises the quality of the input reads supplied to sequencing or assembly workflows. Chimeras typically arise when incomplete extension occurs during a PCR cycle. This causes the DNA polymerase to switch from one template to another and generate hybrid recombinant molecules (Judo et al., 1998). Artificial chimeras are produced through such amplification errors, whereas biological chimeras occur naturally through genomic rearrangements or transcriptional events.

In the context of amplicon-based sequencing, the presence of chimeras can inflate estimates of genetic or microbial diversity and may cause misassemblies during genome reconstruction. Qin et al. (2023) has reported that chimeric sequences may account for more than 10% of raw reads in amplicon datasets. This artifact tends to be most prominent among rare operational taxonomic units (OTUs) or singletons, which are sometimes misinterpreted as novel diversity, further causing the complication of microbial diversity analyses (Gonzalez, Zimmermann, & Saiz-Jimenez, 2004). As such, determining and minimizing PCR-induced chimera formation is vital for improving the quality of mitochondrial genome assemblies, and ensuring the reliability of amplicon sequencing data.

2.3 Existing Traditional Approaches for Chimera Detection

Several computational tools have been developed to identify chimeric sequences in NGS datasets. These tools generally fall into two categories: reference-based and de novo approaches. Reference-based chimera detection, also known as database-dependent detection, is one of the earliest and most widely used computational strategies for identifying chimeric sequences in amplicon-based community studies. These methods rely on the comparison of each query sequence against a curated, high-quality database of known, non-chimeric reference sequences (Edgar et al., 2011).

On the other hand, the de novo chimera detection, also referred to as reference-free detection, represents an alternative computational paradigm that identifies chimeric sequences without reliance on external reference databases. This method infer chimeras based on internal relationships among the sequences present within the dataset itself, making it particularly advantageous in studies of under explored or taxonomically diverse communities where comprehensive reference databases are unavailable or incomplete (Edgar, 2016; Edgar et al., 2011). The underlying assumption on this method is that during PCR, true biological sequences are generally more abundant as they are amplified early and dominate the read pool, whereas chimeric sequences appear later and are generally less abundant. The de novo approach leverage this abundance hierarchy, treating the most abundant sequences as supposed parents and testing whether less abundant sequences can be reconstructed as mosaics of these templates. Compositional and structural similarity are also evaluated to check whether different regions of a candidate

248 sequence correspond to distinct high-abundance sequences.

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

254 **2.3.1 UCHIME**

255 Developed by Edgar et al. (Edgar et al., 2011), UCHIME is one of the most widely
256 used computational tools for detecting chimeric sequences in amplicon sequencing
257 data. The UCHIME algorithm detects chimeras by evaluating how well a query
258 sequence (Q) can be explained as a mosaic of two parent sequences (A and B)
259 from a reference database. The query sequence is first divided into four non-
260 overlapping segments or chunks. Each chunk is independently searched against a
261 reference database that is assumed to be free of chimeras. The best matches to
262 each segment are collected, and from these results, two candidate parent sequences
263 are identified, typically the two sequences that best explain all chunks of the query.
264 Then a three-way alignment among the query (Q) and the two parent candidates
265 (A and B) is done. From this alignment, UCHIME attempts to find a chimeric
266 model (M) which is a hypothetical recombinant sequence formed by concatenating
267 fragments from A and B that best match the observed Q

268 **Chimeric Alignment and Scoring**

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

279 In de novo mode, UCHIME applies an abundance-driven strategy. Only se-
280 quences at least twice as abundant as the query are considered as potential parents.
281 Non-chimeric sequences identified at each step are added iteratively to a growing
282 internal database for subsequent queries.

283 **Limitations of UCHIME**

284 Although UCHIME was a significant advancement in chimera detection, it has
285 notable limitations. According to (Edgar, 2016) and the UCHIME practical notes
286 (Edgar, n.d), many of the accuracy results reported in the original 2011 paper
287 were overly optimistic due to unrealistic benchmark designs that assumed com-
288 plete reference coverage and perfect sequence quality. In practice, UCHIME's
289 accuracy can decline when (1) the reference database is incomplete or contains

erroneous entries; (2) low-divergence chimeras are present, as these closely resemble genuine biological variants; (3) sequence datasets include residual sequencing errors, leading to spurious alignments or misidentification; and (4) the abundance ratio between parent and chimera is distorted by amplification bias. Additionally, UCHIME tends to misclassify sequences as non-chimeric when parent sequences are missing from the database. These limitations motivated the development of UCHIME2.

2.3.2 UCHIME2

To overcome the limitations of its predecessor, UCHIME2 (Edgar, 2016) introduced several methodological and algorithmic refinements that significantly enhanced the accuracy and reliability of chimera detection. One major improvement lies in its approach to uncertainty handling. In earlier versions, sequences with limited reference support were often incorrectly classified as non-chimeric, increasing the likelihood of false negatives. UCHIME2 addresses this issue by designating such ambiguous sequences as “unknown,” thereby providing a more conservative and reliable classification framework.

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

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

323 **2.3.3 CATch**

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

329 To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieurs
330 (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-
331 resents the first ensemble machine learning system designed for chimera detection
332 in 16S rRNA amplicon sequencing. Rather than depending on a single detec-
333 tion strategy, CATCh integrates the outputs of several established tools, includ-
334 ing UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual

335 scores and binary decisions generated by these tools are used as input features for
336 a supervised learning model. The algorithm employs a Support Vector Machine
337 (SVM) with a Pearson VII Universal Kernel (PUK) to determine optimal weight-
338 ings among the input features and to assign each sequence a probability of being
339 chimeric.

340 Benchmarking in both reference-based and de novo modes demonstrated signif-
341 ificant performance improvements. CATCh achieved sensitivities of approximately
342 85 percent in reference-based mode and 92 percent in de novo mode, with corre-
343 sponding specificities of approximately 96 percent and 95 percent. These results
344 indicate that CATCh detected 7 to 12 percent more chimeras than any individual
345 algorithm while maintaining high precision.

346 **2.3.4 ChimPipe**

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

353 The pipeline works with many eukaryotic species that have available genome
354 and annotation data (Rodriguez-Martin et al., 2017). It can also predict multiple
355 isoforms for each gene pair and identify breakpoint coordinates that are useful
356 for reconstructing and verifying chimeric transcripts. Tests using both simulated

357 and real datasets have shown that ChimPipe maintains high accuracy and reliable
358 performance.

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

364 **2.4 Machine Learning Approaches for Chimera** 365 **and Sequence Quality Detection**

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

377 2.4.1 Feature-Based Representations of Genomic Se- 378 quences

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

393 Machine learning models can also leverage alignment-derived features such as
394 the frequency of split alignments, variation in mapping quality, and local cover-
395 age irregularities. Split reads and discordant read pairs are classical indicators
396 of genomic junctions and have been formalized in probabilistic frameworks for
397 structural-variant discovery that integrate multiple evidence types (Layer, Hall, &
398 Quinlan, 2014). Similarly, long-read tools such as Sniffles employ split-alignment
399 and coverage anomalies to accurately localize breakpoints (Sedlazeck et al., 2018).
400 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.

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

2.5 Synthesis of Chimera Detection Approaches

To provide an integrated overview of the literature discussed in this chapter, Table 2.1 summarizes the major chimera detection studies, their methodological

⁴²³ approaches, and their known limitations.

Table 2.1: Summary of Existing Methods and Research Gaps

Method/Study	Scope/Approach	Limitations
Reference-based Chimera Detection	Compares query sequences against curated, non-chimeric reference databases; identifies mosaic sequences by evaluating similarity to known templates.	Depends heavily on completeness and quality of reference databases; often fails when novel taxa or missing parent sequences are present; reduced accuracy for low-divergence chimeras.
De novo Chimera Detection	Identifies chimeras using only internal dataset relationships; relies on abundance patterns and compositional similarity; reconstructs sequences as mosaics of high-abundance parents.	Assumes true sequences are more abundant—fails when amplification bias distorts abundance; struggles with evenly abundant parental sequences; can misclassify highly similar true variants.
UCHIME	Alignment-based chimera detection; segments query sequence, identifies parent candidates, performs 3-way alignment, and computes chimera scores; supports both reference-based and de novo modes.	Accuracy inflated in original benchmarks; suffers under incomplete databases; poor performance on low-divergence chimeras; sensitive to sequencing errors; misclassifies when parents are missing.
UCHIME2	Improved uncertainty handling; classifies ambiguous sequences as unknown; offers multiple sensitivity/specificity modes; more robust with incomplete references; higher sensitivity (93–99%).	Cannot achieve perfect accuracy due to “perfect fake models”; genuine variants may be indistinguishable from artificial recombinants; theoretical detection limit remains.
CATCh	First ML ensemble tool for 16S chimera detection; integrates outputs of UCHIME, ChimeraSlayer, DECIPHER, Pintail, Perseus via SVM classifier; significantly improves sensitivity and specificity.	Depends on performance of underlying tools; ML model limited to features they output; ensemble can still misclassify in datasets with extreme novelty or low coverage.
ChimPipe	Pipeline for detecting fusion genes and transcript-derived chimeras in	Designed for RNA-seq, not amplicons; needs high-quality genome

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

Chapter 3

Research Methodology

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

3.1 Research Activities

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

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

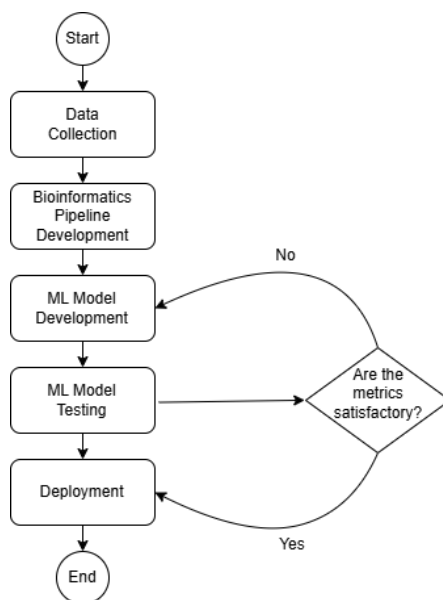


Figure 3.1: Process Diagram of Special Project

458 3.1.1 Data Collection

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

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

464 proximately one (1) week.

465 Data Preprocessing

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

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

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

477 The command parameters are as follows:

- 478 • `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.
- 479 • `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability,
480 all set to a default value of 0.
- 481 • `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 482 • `-N`: number of read pairs, set to 10,000.

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

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

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

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

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

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

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

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

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

```
508 samtools view -bS clean.sam -o clean.bam
509 samtools view -bS chimeric.sam -o chimeric.bam
510
511 samtools sort clean.bam -o clean.sorted.bam
512 samtools index clean.sorted.bam
513
514 samtools sort chimeric.bam -o chimeric.sorted.bam
515 samtools index chimeric.sorted.bam
```

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

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

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

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

529 **3.1.2 Bioinformatics Tools Pipeline**

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

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

545 **Supplementary Alignment Flag**

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

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

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

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

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

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

562 **K-mer Composition Difference**

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

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

Microhomology Length

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

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

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

length and GC content of microhomology.

3.1.3 Machine Learning Model Development

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

The labelled dataset was randomly partitioned into training (80%) and test (20%) subsets using stratified sampling to preserve the 1:1 ratio of clean to chimeric reads. Model development and evaluation were implemented in Python (Version 3.11) using the `scikit-learn`, `xgboost`, `lightgbm`, and `catboost` libraries. A broad panel of classification algorithms was then benchmarked on the training data to obtain a fair comparison of different model families under identical feature conditions. The panel included: a trivial dummy classifier, L2-regularized logistic regression, a calibrated linear support vector machine (SVM), k -nearest neighbours, Gaussian Naïve Bayes, decision-tree ensembles (Random Forest, Extremely Randomized Trees, and Bagging with decision trees), gradient boosting methods (Gradient Boosting, XGBoost, LightGBM, and CatBoost), and a shallow multilayer perceptron (MLP).

For each model, five-fold stratified cross-validation was performed on the training set. In every fold, four-fifths of the data were used for fitting and the remaining one-fifth for validation. Mean cross-validation accuracy, precision, recall, F1-score for the chimeric class, and area under the receiver operating characteristic curve (ROC-AUC) were computed to summarize performance and rank candidate methods. This baseline screen allowed comparison of linear, probabilistic, neural, and ensemble-based approaches and identified tree-based ensemble and boosting models as consistently strong performers relative to simpler baselines.

3.1.4 Model Benchmarking, Hyperparameter Optimization, and Evaluation

Model selection and refinement proceeded in two stages. First, the cross-validation results from the broad panel were used to identify a subset of competitive models for more detailed optimization. Specifically, ten model families were carried forward: L2-regularized logistic regression, calibrated linear SVM, Random Forest, ExtraTrees, Gradient Boosting, XGBoost, LightGBM, CatBoost, Bagging with decision trees, and a shallow MLP. This subset spans both linear and non-linear decision boundaries, but emphasizes ensemble and boosting methods, which showed superior F1 and ROC-AUC in the initial benchmark.

Second, hyperparameter optimization was conducted for each of the ten selected models using randomized search with five-fold stratified cross-validation (`RandomizedSearchCV`). For tree-based ensembles, the search space included the number of trees, maximum depth, minimum samples per split and leaf, and the fraction of features considered at each split. For boosting methods, key hyper-

parameters such as the number of boosting iterations, learning rate, tree depth, subsampling rate, and column subsampling rate were tuned. For the MLP, the number and size of hidden layers, learning rate, and L_2 regularization strength were varied. In all cases, the primary optimisation criterion was the F1-score of the chimeric class, averaged across folds.

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

3.1.5 Feature Importance and Interpretation

To relate model decisions to biologically meaningful signals, feature-importance analyses were performed on the best-performing tree-based models. Two complementary approaches were used. First, built-in importance measures from ensemble methods (e.g., split-based importances in Random Forest and Gradient Boosting) were examined to obtain an initial ranking of features based on their contribution to reducing impurity. Second, model-agnostic permutation importance was com-

puted on the test set by repeatedly permuting each feature column while keeping all others fixed and measuring the resulting decrease in F1-score. Features whose permutation led to a larger performance drop were interpreted as more influential for chimera detection.

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

3.1.6 Validation and Testing

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

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

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

695 **3.1.7 Documentation**

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

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

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

709 3.2 Calendar of Activities

710 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
 711 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	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •

712 References

- 713 Anderson, S., Bankier, A., Barrell, B., Bruijn, M., Coulson, A., Drouin, J., ...
714 Young, I. (1981, 04). Sequence and organization of the human mitochondrial
715 genome. *Nature*, 290, 457-465. doi: 10.1038/290457a0
- 716 Arango, G., Garner, E., Pruden, A., Heath, L., Vikesland, P., & Zhang, L. (2018,
717 02). Deeparg: A deep learning approach for predicting antibiotic resistance
718 genes from metagenomic data. *Microbiome*, 6. doi: 10.1186/s40168-018
719 -0401-z
- 720 Bentley, D. R., Balasubramanian, S., Swerdlow, H. P., Smith, G. P., Milton, J.,
721 Brown, C. G., ... Smith, A. J. (2008). Accurate whole human genome
722 sequencing using reversible terminator chemistry. *Nature*, 456(7218), 53–
723 59. doi: 10.1038/nature07517
- 724 Boore, J. L. (1999). Animal mitochondrial genomes. *Nucleic Acids Research*,
725 27(8), 1767–1780. doi: 10.1093/nar/27.8.1767
- 726 Cameron, S. L. (2014). Insect mitochondrial genomics: Implications for evolution
727 and phylogeny. *Annual Review of Entomology*, 59, 95–117. doi: 10.1146/
728 annurev-ento-011613-162007
- 729 Dierckxsens, N., Mardulyn, P., & Smits, G. (2017). Novoplasty: de novo assembly
730 of organelle genomes from whole genome data. *Nucleic Acids Research*,

731 45(4), e18. doi: 10.1093/nar/gkw955

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

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

734 CorpusID:88955007

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

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

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

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

739 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

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

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

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

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

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

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

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

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

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

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

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

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

752 10.1093/nar/gkt371

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

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

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

756 10.1186/s13059-020-02154-5

757 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
758 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
759 1819–1825. doi: 10.1093/nar/26.7.1819

760 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
761 (2021). Mitochondrial dna reveals genetically structured haplogroups of
762 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
763 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

764 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework
765 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-
766 -15-6-r84

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

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

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

776 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
777 an ensemble classifier for chimera detection in 16s rna sequencing stud-
778 ies. *Applied and Environmental Microbiology*, *81*(5), 1573–1584. Retrieved
779 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
780 10.1128/AEM.02896-14

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

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

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

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

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

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

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

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

809 .1038/s41592-018-0001-7

810 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
 811 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
 812 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
 813 [.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)
 814 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

815 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
 816 11). Large-scale machine learning for metagenomics sequence classifica-
 817 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
 818 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683

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