

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 2, 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	Validation and Testing	31
48	3.1.5	Documentation	31

49	3.2 Calendar of Activities	32
----	--------------------------------------	----

50 List of Figures

<small>51</small>	3.1 Process Diagram of Special Project	23
-------------------	--	----

52 List of Tables

<small>53</small>	2.1 Summary of Existing Methods and Research Gaps	20
<small>54</small>	3.1 Timetable of Activities	32

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.

114 1.3 Research Objectives

115 1.3.1 General Objective

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

120 1.3.2 Specific Objectives

121 Specifically, the study aims to:

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

156 1.5 Significance of the Research

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

Chapter 2

Review of Related Literature

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

2.1 The Mitochondrial Genome

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

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

186 **2.1.1 Mitochondrial Genome Assembly**

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

245 sequence correspond to distinct high-abundance sequences.

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

251 **2.3.1 UCHIME**

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

265 **Chimeric Alignment and Scoring**

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

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

280 **Limitations of UCHIME**

281 Although UCHIME was a significant advancement in chimera detection, it has
282 notable limitations. According to (Edgar, 2016) and the UCHIME practical notes
283 (Edgar, n.d), many of the accuracy results reported in the original 2011 paper
284 were overly optimistic due to unrealistic benchmark designs that assumed com-
285 plete reference coverage and perfect sequence quality. In practice, UCHIME's
286 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.

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

320 **2.3.3 CATch**

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

326 To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieus
327 (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-
328 resents the first ensemble machine learning system designed for chimera detection
329 in 16S rRNA amplicon sequencing. Rather than depending on a single detec-
330 tion strategy, CATCh integrates the outputs of several established tools, includ-
331 ing UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual

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

337 Benchmarking in both reference-based and de novo modes demonstrated signif-
338 ificant performance improvements. CATCh achieved sensitivities of approximately
339 85 percent in reference-based mode and 92 percent in de novo mode, with corre-
340 sponding specificities of approximately 96 percent and 95 percent. These results
341 indicate that CATCh detected 7 to 12 percent more chimeras than any individual
342 algorithm while maintaining high precision.

343 **2.3.4 ChimPipe**

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

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

354 and real datasets have shown that ChimPipe maintains high accuracy and reliable
355 performance.

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

361 **2.4 Machine Learning Approaches for Chimera** 362 **and Sequence Quality Detection**

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

374 2.4.1 Feature-Based Representations of Genomic Se- 375 quences

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

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

⁴²⁰ 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

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

449 formatics pipeline that includes k-mer extraction and homology-based filtering to
 450 prepare the data for model construction. The machine learning model will subse-
 451 quently be trained and tested using the processed datasets to assess its precision
 452 and accuracy. The model will undergo refinement and retraining until it meets the
 453 required performance threshold, after which it will proceed to validation, testing,
 454 and deployment.

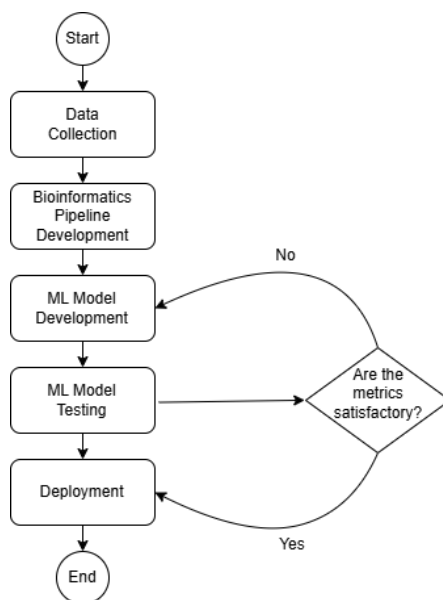


Figure 3.1: Process Diagram of Special Project

455 3.1.1 Data Collection

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

459 This step is scheduled to begin in the first week of November 2025 and is
 460 expected to be completed by the last week of November 2025, with a total duration

461 of approximately one (1) month.

462 Data Preprocessing

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

467 Sequencing data were simulated from the NCBI reference genome using **wgsim**
468 (Version 1.13). First, 10,000 paired-end reads (R1 and R2) were generated from
469 the original reference (`original_reference.fasta`) and designated as clean reads
470 using the command:

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

473 The command parameters are as follows:

- 474 • **-1** and **-2**: read lengths of 150 base pairs for each paired-end read.
- 475 • **-r**, **-R**, **-X**: mutation rate, fraction of indels, and indel extension probability,
476 all set to a default value of 0.
- 477 • **-e**: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 478 • **-N**: number of read pairs, set to 10,000.

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

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

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

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

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

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

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

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

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

```
504 samtools view -bS clean.sam -o clean.bam
505 samtools view -bS chimeric.sam -o chimeric.bam
506
507 samtools sort clean.bam -o clean.sorted.bam
508 samtools index clean.sorted.bam
509
510 samtools sort chimeric.bam -o chimeric.sorted.bam
511 samtools index chimeric.sorted.bam
```

512 BAM files are the compressed binary version of SAM files, which enables faster
513 processing and reduced storage. Sorting will arrange reads by genomic coordi-
514 nates, and indexing will allow detection of supplementary alignments (SA) as a
515 feature for the machine learning model.

516 The results of this process will be used for feature extraction. Once the primary
517 features have been extracted from the reads, a custom script will be created to
518 merge them into a single dataset and assign class labels: clean reads as “0” and
519 chimeric reads as “1”. The final dataset will contain 10,000 clean reads and 10,000
520 chimeric reads (a total of 20,000 entries) to ensure equal representation of both
521 classes. The merged dataset will be saved in TSV (**.tsv**) format and subsequently
522 split so that 80% are used for training and 20% for testing.

523 This whole process is scheduled to start in the first week of November 2025
524 and is expected to be completed by the last week of November 2025, with a total
525 duration of approximately one (1) month.

526 **3.1.2 Bioinformatics Tools Pipeline**

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

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

542 **Supplementary Alignment Flag**

543 Supplementary alignment information will be assessed using the mapped and
544 sorted BAM files (`clean.sorted.bam` and `chimeric.sorted.bam`) generated
545 from the data preprocessing stage. Alignment summaries will be checked using
546 `samtools flagstat` to obtain preliminary quality-control statistics, including
547 counts of primary, secondary, and supplementary (SA) alignments.

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

```
550 samtools view -h clean.sorted.bam -o clean.sorted.sam  
551 samtools view -h chimeric.sorted.bam -o chimeric.sorted.sam
```

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

559 **K-mer Composition Difference**

560 Chimeric reads often comprise fragments from distinct genomic regions, resulting
561 in a compositional discontinuity between segments. Comparing k-mer frequency
562 profiles between the left and right halves of a read allows detection of such abrupt

563 compositional shifts, independent of alignment information. This will be obtained
564 using Jellyfish, a fast k-mer counting software. For each read, the sequence will
565 be divided into two segments, either at the midpoint or at empirically determined
566 breakpoints inferred from supplementary alignment data, to generate left and right
567 sequence segments. Jellyfish will then compute k-mer frequency profiles (with $k =$
568 5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized
569 and compared using distance metrics such as cosine similarity or Jensen–Shannon
570 divergence to quantify compositional disparity between the two halves of the same
571 read. The resulting difference scores will be stored in a structured TSV file.

572 **Micro-homology Length**

573 The micro-homology length will be computed using a custom Python script that
574 detects the longest exact suffix–prefix overlap within ± 30 base pairs surround-
575 ing a candidate breakpoint. This analysis identifies the number of consecutive
576 bases shared between the end of one segment and the beginning of another. The
577 presence and length of such micro-homology are classic molecular signatures of
578 PCR-induced template switching, where short identical regions (typically 3–15
579 base pairs) promote premature termination and recombination of DNA synthesis
580 on a different template strand. Quantifying micro-homology allows assessment of
581 whether the suspected breakpoint reflects PCR artifacts or true biological variants.
582 Each read will therefore be annotated with its corresponding micro-homology
583 length, overlap sequence, and GC content.

584 After extracting the three primary features, all resulting TSV files will be
585 joined using the read identifier as a common key to generate a unified feature ma-

586 trix. Additional read-level metadata such as read length, mean base quality, and
587 number of clipped bases will also be included to provide contextual information.
588 This consolidated dataset will serve as the input for subsequent machine-learning
589 model development and evaluation.

590 **3.1.3 Machine Learning Model Development**

591 This study will explore multiple machine-learning approaches to detect PCR-
592 induced chimeras from mitochondrial Illumina reads: Support Vector Machines
593 (SVM) to separate reads with complex patterns, decision trees to capture hier-
594 archical interactions among SA count, k-mer composition, and micro-homology
595 length, logistic regression as a linear baseline, Random Forest (RF) to improve
596 stability and reduce variance, and gradient boosting (e.g., XGBoost) to model
597 non-linear relationships among the extracted features. Using these approaches
598 enables a balanced assessment of predictive performance and interpretability.

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

606 Model development and evaluation will be implemented in Python (v3.11)
607 using the `scikit-learn` and `xgboost` libraries. Standard metrics including ac-

608 curacy, precision, recall, F1-score, and area under the ROC curve (AUC) will be
609 computed to quantify predictive performance.

610 **3.1.4 Validation and Testing**

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

619 Comparative evaluation across all candidate algorithms, including SVM, de-
620 cision trees, logistic regression, Random Forest, gradient boosting, and others,
621 will determine which models demonstrate the highest predictive performance and
622 computational efficiency under identical data conditions. Their metrics will be
623 compared to identify the which algorithms are most suitable for further refine-
624 ment.

625 **3.1.5 Documentation**

626 Comprehensive documentation will be maintained throughout the study to en-
627 sure transparency and reproducibility. All stages of the research, including data
628 gathering, preprocessing, feature extraction, model training, and validation, will

629 be systematically recorded in a `.README` file in the GitHub repository. For each
 630 analytical step, the corresponding parameters, software versions, and command
 631 line scripts will be documented to enable exact replication of results.

632 The repository structure will follow standard research data management
 633 practices, with clear directories for datasets and scripts. Computational
 634 environments will be standardized using Conda, with an environment file
 635 (`environment.arm.yml`) specifying dependencies and package versions to main-
 636 tain consistency across systems.

637 For manuscript preparation and supplementary materials, Overleaf (L^AT_EX)
 638 will be used to produce publication-quality formatting and consistent referencing.

639 3.2 Calendar of Activities

640 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
 641 bullet point corresponds to approximately one week of planned activity.

Table 3.1: Timetable of Activities

Activities (2025)	Nov	Dec	Jan	Feb	Mar	Apr	May
Data Collection and Simulation	• • • •						
Bioinformatics Tools Pipeline			• • • •	• • • •			
Machine Learning Development			• •	• • • •	• • • •	• •	
Testing and Validation						• •	• • • •
Documentation	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •	• • • •

References

- Anderson, S., Bankier, A., Barrell, B., Bruijn, M., Coulson, A., Drouin, J., ...
Young, I. (1981, 04). Sequence and organization of the human mitochondrial
genome. *Nature*, 290, 457-465. doi: 10.1038/290457a0
- Arango, G., Garner, E., Pruden, A., Heath, L., Vikesland, P., & Zhang, L. (2018,
02). Deeparg: A deep learning approach for predicting antibiotic resistance
genes from metagenomic data. *Microbiome*, 6. doi: 10.1186/s40168-018
-0401-z
- Bentley, D. R., Balasubramanian, S., Swerdlow, H. P., Smith, G. P., Milton, J.,
Brown, C. G., ... Smith, A. J. (2008). Accurate whole human genome
sequencing using reversible terminator chemistry. *Nature*, 456(7218), 53–
59. doi: 10.1038/nature07517
- 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*,

661 45(4), e18. doi: 10.1093/nar/gkw955

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

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

664 CorpusID:88955007

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

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

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

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

669 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

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

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

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

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

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

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

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

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

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

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

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

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

682 10.1093/nar/gkt371

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

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

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

686 10.1186/s13059-020-02154-5

687 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
688 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
689 1819–1825. doi: 10.1093/nar/26.7.1819

690 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
691 (2021). Mitochondrial dna reveals genetically structured haplogroups of
692 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
693 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

694 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework
695 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-
696 -15-6-r84

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

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

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

706 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
707 an ensemble classifier for chimera detection in 16s rrna sequencing stud-
708 ies. *Applied and Environmental Microbiology*, *81*(5), 1573-1584. Retrieved
709 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
710 [10.1128/AEM.02896-14](https://journals.asm.org/doi/abs/10.1128/aem.02896-14)

711 Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J.
712 (2023). Effects of error, chimera, bias, and gc content on the accuracy of

713 amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from [https://](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)
 714 journals.asm.org/doi/abs/10.1128/msystems.01025-23 doi: 10.1128/
 715 [msystems.01025-23](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)
 716 Qiu, X., Wu, L., Huang, H., McDonel, P. E., Palumbo, A. V., Tiedje, J. M., &
 717 Zhou, J. (2001). Evaluation of pcr-generated chimeras, mutations, and het-
 718 eroduplexes with 16s rna gene-based cloning. *Applied and Environmental*
 719 *Microbiology*, 67(2), 880–887. doi: 10.1128/AEM.67.2.880-887.2001
 720 Ren, J., Song, K., Deng, C., Ahlgren, N., Fuhrman, J., Li, Y., ... Sun, F. (2020,
 721 01). Identifying viruses from metagenomic data using deep learning. *Quan-*
 722 *titative Biology*, 8. doi: 10.1007/s40484-019-0187-4
 723 Rodriguez-Martin, B., Palumbo, E., Marco-Sola, S., Griebel, T., Ribeca, P.,
 724 Alonso, G., ... Djebali, S. (2017, 01). Chimpipes: Accurate detection of
 725 fusion genes and transcription-induced chimeras from rna-seq data. *BMC*
 726 *Genomics*, 18. doi: 10.1186/s12864-016-3404-9
 727 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). Vsearch: a
 728 versatile open source tool for metagenomics. *PeerJ*, 4, e2584. doi: 10.7717/
 729 [peerj.2584](https://doi.org/10.7717/peerj.2584)
 730 Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler,
 731 A., & Schatz, M. (2018, 06). Accurate detection of complex structural
 732 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10
 733 [.1038/s41592-018-0001-7](https://doi.org/10.1038/s41592-018-0001-7)
 734 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
 735 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
 736 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
 737 [.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)
 738 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

739 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
740 11). Large-scale machine learning for metagenomics sequence classifica-
741 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
742 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683
743 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*
744 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI
745 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)
746 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)