# MitoChime: A Machine-Learning Pipeline for Detecting PCR-Induced Chimeras in Mitochondrial Illumina Reads

A Special Project Proposal

Presented to

the Faculty of the Division of Physical Sciences and Mathematics

College of Arts and Sciences

University of the Philippines Visayas

Miag-ao, Iloilo

In Partial Fulfillment

of the Requirements for the Degree of

Bachelor of Science in Computer Science

by

Duranne Duran

Yvonne Lin

Daniella Pailden

Adviser

Francis D. Dimzon, Ph.D.

December 5, 2025

# Contents

# List of Figures

# List of Tables

# 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

1

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.

## 1.3 Research Objectives

### 1.3.1 General Objective

This study aims to develop and evaluate a machine learning-based pipeline (MitoChime) that detects PCR-induced chimeric reads in *Sardinella lemuru* mitochondrial sequencing data in order to improve the quality and reliability of downstream mitochondrial genome assemblies.

### 1.3.2 Specific Objectives

Specifically, the study aims to:

1. construct simulated *Sardinella lemuru* Illumina paired-end datasets containing both clean and PCR-induced chimeric reads,

2. extract alignment-based and sequence-based features such as k-mer composition, junction complexity, and split-alignment counts from both clean and chimeric reads,

3. train, validate, and compare supervised machine-learning models for classifying reads as clean or chimeric,

4. determine feature importance and identify indicators of PCR-induced chimerism,

5. integrate the optimized classifier into a modular and interpretable pipeline 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.

## 1.5 Significance of the Research

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

7

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

### 2.1.1   Mitochondrial Genome Assembly

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

8

## 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

<sub>248</sub> sequence correspond to distinct high-abundance sequences.

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

### <sub>254</sub> 2.3.1   UCHIME

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

11

## Chimeric Alignment and Scoring

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

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

## Limitations of UCHIME

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

13

In comparative evaluations, UCHIME2 demonstrated superior detection performance, achieving sensitivity levels between 93% and 99% and lower overall error rates than earlier versions or other contemporary tools such as DECIPHER and ChimeraSlayer. Despite these advances, the study also acknowledged a fundamental limitation in chimera detection: complete error-free identification is theoretically unattainable. This is due to the presence of "perfect fake models," wherein genuine non-chimeric sequences can be perfectly reconstructed from other reference fragments. This underscore the uncertainty in differentiating authentic biological sequences from artificial recombinants based solely on sequence similarity, emphasizing the need for continued methodological refinement and cautious interpretation of results.

### 2.3.3 CATch

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

To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieurs (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which represents the first ensemble machine learning system designed for chimera detection in 16S rRNA amplicon sequencing. Rather than depending on a single detection strategy, CATCh integrates the outputs of several established tools, including UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual

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

Benchmarking in both reference-based and de novo modes demonstrated significant performance improvements. CATCh achieved sensitivities of approximately 85 percent in reference-based mode and 92 percent in de novo mode, with corresponding specificities of approximately 96 percent and 95 percent. These results indicate that CATCh detected 7 to 12 percent more chimeras than any individual algorithm while maintaining high precision.

## 2.3.4   ChimPipe

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

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

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

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

## 2.4 Machine Learning Approaches for Chimera and Sequence Quality Detection

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

## 2.4.1 Feature-Based Representations of Genomic Sequences

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

Machine learning models can also leverage alignment-derived features such as the frequency of split alignments, variation in mapping quality, and local coverage irregularities. Split reads and discordant read pairs are classical indicators of genomic junctions and have been formalized in probabilistic frameworks for structural-variant discovery that integrate multiple evidence types (Layer, Hall, & Quinlan, 2014). Similarly, long-read tools such as Sniffles employ split-alignment and coverage anomalies to accurately localize breakpoints (Sedlazeck et al., 2018). Modern aligners such as Minimap2 (Li, 2018) output supplementary (SA tags) and

17

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

18

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 |

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

22

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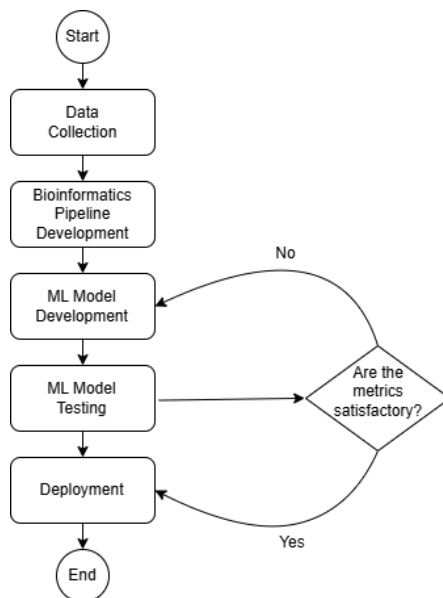


Figure 3.1: Process Diagram of Special Project

### 3.1.1 Data Collection

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

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

proximately one (1) week.

## Data Preprocessing

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

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

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

The command parameters are as follows:

- `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.

- `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability, all set to a default value of 0.

- `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.

- `-N`: number of read pairs, set to 10,000.

24

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

25

threads. The resulting clean and chimeric SAM files contain the alignment positions of each read relative to the original reference genome.

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

```
samtools view -bS clean.sam -o clean.bam
samtools view -bS chimeric.sam -o chimeric.bam

samtools sort clean.bam -o clean.sorted.bam
samtools index clean.sorted.bam

samtools sort chimeric.bam -o chimeric.sorted.bam
samtools index chimeric.sorted.bam
```

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

The total number of simulated reads was expected to be 40,000. The final collection of reads contained 19,984 clean reads and 20,000 chimeric reads (39,984 entries in total), providing a roughly balanced distribution between the two classes. After alignment with `minimap2`, only 19,984 clean reads remained because unmapped reads were not included in the BAM file. Some sequences failed to align due to the 5% error rate defined during `wgsim` simulation, which produced mismatches that caused certain reads to fall below the aligner's matching threshold.

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

26

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

## 3.1.2 Bioinformatics Tools Pipeline

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

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

**Supplementary Alignment Flag**

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

27

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

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

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

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

**K-mer Composition Difference**

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

28

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.

## Micro-homology Length

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

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

### 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

30

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 nonlinear 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 hyperparameters 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 computed 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

32

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 the trained models generalized to unseen data. All feature extraction and pre-processing steps were performed using the same bioinformatics pipeline to ensure

33

consistency and comparability across validation stages.

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

### 3.1.7   Documentation

Comprehensive documentation was maintained throughout the study to ensure transparency and reproducibility. All stages of the research, including data gathering, preprocessing, feature extraction, model training, and validation, were systematically recorded in a `.README` file in the GitHub repository. For each analytical step, the corresponding parameters, software versions, and command line scripts were documented to enable exact replication of results.

The repository structure followed standard research data management practices, with clear directories for datasets and scripts. Computational environments were standardized using Conda, with an environment file (`environment.arm.yml`) specifying dependencies and package versions to maintain consistency across systems.

For manuscript preparation and supplementary materials, Overleaf (LaTeX) was used to produce publication-quality formatting and consistent referencing.

## 3.2 Calendar of Activities

Table 3.1 presents the project timeline in the form of a Gantt chart, where each 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 | ● ● ●● | ● ● ●● | ● ● ●● | ● ● ●● | ● ● ●● | ● ● ●● | ● ● ●● |

35

# 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*,

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

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

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

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

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

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

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

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

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

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

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

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

Li, H. (2018, 05). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, *34*(18), 3094-3100. Retrieved from `https://doi.org/10.1093/` `bioinformatics/bty191` doi: 10.1093/bioinformatics/bty191

Liang, Q., Bible, P. W., Liu, Y., Zou, B., & Wei, L. (2020, 02). Deepmicrobes: taxonomic classification for metagenomics with deep learning. *NAR Genomics and Bioinformatics*, *2*(1), lqaa009. Retrieved from `https://` `doi.org/10.1093/nargab/lqaa009` doi: 10.1093/nargab/lqaa009

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

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

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 rrna 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). Chimpipe: 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 .1038/s41592-018-0001-7

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

808 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,

809    11). Large-scale machine learning for metagenomics sequence classifica-

810    tion. *Bioinformatics*, *32*(7), 1023-1032. Retrieved from `https://doi.org/`

811    `10.1093/bioinformatics/btv683`  doi: 10.1093/bioinformatics/btv683

812 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*

813    *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI

814    Technical Paper Series. Retrieved from `https://nfrdi.da.gov.ph/tpjf/`

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