

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                             November 29, 2025

# Contents

|    |   |          |
|----|---|----------|
| 21 | <b>1 Introduction</b>                               | <b>1</b> |
| 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 | <b>2 Review of Related Literature</b>               | <b>7</b> |
| 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 | 2.2.1    | Effects of Chimeric Reads on Organelle Genome Assembly          | 10        |
| 34 | 2.3      | Existing Traditional Approaches for Chimera Detection . . . . . | 11        |
| 35 | 2.3.1    | UCHIME . . . . .  | 12        |
| 36 | 2.3.2    | UCHIME2 . . . . .   | 14        |
| 37 | 2.3.3    | CATch . . . . .   | 16        |
| 38 | 2.3.4    | ChimPipe . . . . .  | 17        |
| 39 | 2.4      | Machine Learning Approaches for Chimera and Sequence Quality    |           |
| 40 |          | Detection . . . . .   | 18        |
| 41 | 2.4.1    | Feature-Based Representations of Genomic Sequences . . .        | 18        |
| 42 | 2.5      | Synthesis of Chimera Detection Approaches . . . . .             | 20        |
| 43 | <b>3</b> | <b>Research Methodology</b>                                     | <b>23</b> |
| 44 | 3.1      | Research Activities . . . . .                                   | 23        |
| 45 | 3.1.1    | Data Collection . . . . .                                       | 24        |
| 46 | 3.1.2    | Bioinformatics Tools Pipeline . . . . .                         | 26        |
| 47 | 3.1.3    | Machine-Learning Model Development . . . . .                    | 29        |
| 48 | 3.1.4    | Validation and Testing . . . . .                                | 30        |

|    |                                      |    |
|----|--------------------------------------|----|
| 49 | 3.1.5 Documentation . . . . .        | 30 |
| 50 | 3.2 Calendar of Activities . . . . . | 31 |

# 51 List of Figures

|                   |  |    |
|-------------------|--|----|
| <small>52</small> | 3.1 Process Diagram of Special Project . . . . . | 24 |
|-------------------|--|----|

# 53 List of Tables

|                   |   |    |
|-------------------|---|----|
| <small>54</small> | 2.1 Summary of Existing Methods and Research Gaps . . . . . | 21 |
| <small>55</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 mis-joined reads can reduce assembly contiguity and introduce false junctions during organelle genome reconstruction (Dierckxsens, Mardulyn, & Smits, 2017; Hahn, Bachmann, & Chevreux, 2013; Jin et al., 2020). Existing assembly tools such as GetOrganelle and MITObim assume that input reads are largely free of such artifacts (Hahn et al., 2013; Jin et al., 2020). Consequently, undetected chimeras may produce fragmented assemblies or misidentified organellar boundaries. To ensure accurate reconstruction of mitochondrial genomes, a reliable and automated method for detecting and filtering PCR-induced chimeras before assembly is essential.

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



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

## 97 1.2 Problem Statement

98 While NGS technologies have revolutionized genomic data acquisition, the ac-  
99 curacy of mitochondrial genome assembly remains limited by artifacts produced  
100 during PCR amplification. These chimeric reads can distort assembly graphs and  
101 cause misassemblies, with especially severe effects in small, circular mitochon-  
102 drial genomes (Boore, 1999; Cameron, 2014). Existing assembly pipelines such  
103 as GetOrganelle, MITObim, and NOVOPlasty assume that sequencing reads are  
104 free of such artifacts (Dierckxsens et al., 2017; Hahn et al., 2013; Jin et al., 2020).  
105 At the Philippine Genome Center Visayas, several mitochondrial assemblies have  
106 failed or yielded incomplete contigs despite sufficient coverage, suggesting that  
107 undetected chimeric reads compromise assembly reliability. Meanwhile, exist-  
108 ing chimera-detection tools such as UCHIME and VSEARCH were developed  
109 primarily for amplicon-based microbial community analysis and rely heavily on  
110 reference or taxonomic comparisons (Edgar, Haas, Clemente, Quince, & Knight,  
111 2011; Rognes, Flouri, Nichols, Quince, & Mahé, 2016). These approaches are un-  
112 suitable for single-species organellar data, where complete reference genomes are  
113 often unavailable. Therefore, there is a pressing need for a reference-independent,  
114 data-driven tool capable of automatically detecting and filtering PCR-induced  
115 chimeras in mitochondrial sequencing datasets.

## 116 1.3 Research Objectives

### 117 1.3.1 General Objective

118 This study aims to develop and evaluate a machine-learning-based pipeline (Mi-  
119 toChime) that detects PCR-induced chimeric reads in *Sardinella lemuru* mito-  
120 chondrial sequencing data in order to improve mitochondrial genome assembly  
121 accuracy.

### 122 1.3.2 Specific Objectives

123 Specifically, the study aims to:

- 124 1. construct empirical and simulated *Sardinella lemuru* Illumina paired-end  
125 datasets containing both clean and PCR-induced chimeric reads,
- 126 2. extract alignment-based and sequence-based features such as k-mer compo-  
127 sition, junction complexity, and split-alignment counts from both clean and  
128 chimeric reads,
- 129 3. train, validate, and compare supervised machine-learning models for classi-  
130 fying reads as clean or chimeric,
- 131 4. determine feature importance and identify the most informative indicators  
132 of PCR-induced chimerism,
- 133 5. integrate the optimized classifier into a modular and interpretable pipeline  
134 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 ongoing *S. lemuru* sequencing projects at the Philippine Genome Center Visayas; (3) to make use of existing *S. lemuru* mitochondrial assemblies and raw datasets available in public repositories such as the National Center for Biotechnology Information (NCBI); 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 structural rearrangements in nuclear genomes. Feature extraction is restricted to low-dimensional, hand-crafted alignment and sequence statistics—such as k-mer frequency profiles, GC content, read length, soft-clipping and split-alignment counts, and mapping quality—rather than high-dimensional deep-learning embeddings, so that model behaviour remains interpretable and the pipeline can be executed on standard workstations at PGC Visayas. Testing on long-read platforms (e.g., Nanopore, PacBio) and on other taxa lies beyond the scope of this project; the implemented pipeline is evaluated only on short-read *S. lemuru* datasets.

## 158 1.5 Significance of the Research

159 This research provides both methodological and practical contributions to mi-  
160 tochondrial genomics and bioinformatics. First, MitoChime filters PCR-induced  
161 chimeric reads prior to genome assembly, with the goal of improving the conti-  
162 guity and correctness of *Sardinella lemuru* mitochondrial assemblies. Second, it  
163 replaces ad hoc manual curation with a documented, data-driven workflow, im-  
164 proving automation and reproducibility. Third, the pipeline is designed to run  
165 on modest computing infrastructures commonly available in regional laborato-  
166 ries, enabling routine use at the Philippine Genome Center Visayas. In addition,  
167 MitoChime supports local capacity building by providing a concrete example of  
168 machine-learning integration into mitochondrial genome analysis, consistent with  
169 the mandate of the Philippine Genome Center Visayas. Finally, more reliable  
170 mitochondrial assemblies for *S. lemuru* provide a stronger basis for downstream  
171 applications such as fisheries management, population genetics, and biodiversity  
172 assessments.

## 173 Chapter 2

### 174 Review of Related Literature

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

#### 181 2.1 The Mitochondrial Genome

182 Mitochondrial genome (mtDNA) is a small, typically circular molecule found in  
183 most eukaryotes. It encodes essential genes involved in oxidative phosphorylation  
184 and energy metabolism. Because of its conserved structure and maternal inher-  
185 itance, mtDNA has become a valuable genetic marker for studies in evolution,  
186 population genetics, and phylogenetics (Anderson et al., 1981; Boore, 1999). In

187 animal species, the mitochondrial genome ranges from 15–20 kilobase and contains  
188 13 protein-coding genes, 22 tRNAs, and two rRNAs arranged compactly without  
189 introns (Gray, 2012). In comparison to nuclear DNA the ratio of the number  
190 of copies of mtDNA is higher and has relatively simple organization which make  
191 it particularly suitable for genome sequencing and assembly studies (Dierckxsens  
192 et al., 2017). Moreover, mitochondrial genomes provide crucial insights into evo-  
193 lutionary relationships among species and are increasingly used for testing new  
194 genomic assembly and analysis methods.

### 195 **2.1.1 Mitochondrial Genome Assembly**

196 Mitochondrial genome assembly refers to the reconstruction of the complete mito-  
197 chondrial DNA (mtDNA) sequence from raw or fragmented sequencing reads. It is  
198 conducted to obtain high-quality, continuous representations of the mitochondrial  
199 genome that can be used for a wide range of analyses, including species identi-  
200 fication, phylogenetic reconstruction, evolutionary studies, and investigations of  
201 mitochondrial diseases. Because mtDNA evolves relatively rapidly and is mater-  
202 nally inherited, its assembled sequence provides valuable insights into population  
203 structure, lineage divergence, and adaptive evolution across taxa (Boore, 1999).  
204 Compared to nuclear genome assembly, assembling the mitochondrial genome is  
205 often considered more straightforward but still encounters distinct technical chal-  
206 lenges such as sequencing errors, low coverage regions, and chimeric reads that can  
207 distort the final assembly, leading to incomplete or misassembled genomes. These  
208 errors can propagate into downstream analyses, emphasizing the need for robust  
209 chimera detection and sequence validation methods in mitochondrial genome re-

210 search.

## 211 **2.2 PCR Amplification and Chimera Formation**

212 Polymerase Chain Reaction (PCR) plays an important role in next-generation  
213 sequencing (NGS) library preparation, as it amplifies target DNA fragments for  
214 downstream analysis. However, the amplification process can also introduce arti-  
215 facts that affect data accuracy, one of them being the formation of chimeric se-  
216 quences. Chimeras typically arise when incomplete extension occurs during a PCR  
217 cycle. This causes the DNA polymerase to switch from one template to another  
218 and generate hybrid recombinant molecules (Judo et al., 1998). Artificial chimeras  
219 are produced through such amplification errors, whereas biological chimeras oc-  
220 cur naturally through genomic rearrangements or transcriptional events. These  
221 biological chimeras can have functional roles and may encode tissue-specific novel  
222 proteins that link to cellular processes or diseases (Frenkel-Morgenstern et al.,  
223 2012).

224 In the context of amplicon-based sequencing, PCR-induced chimeras can sig-  
225 nificantly distort analytical outcomes. Their presence artificially inflates estimates  
226 of genetic or microbial diversity and may cause misassemblies during genome re-  
227 construction. (Qin et al., 2023) has reported that chimeric sequences may account  
228 for more than 10% of raw reads in amplicon datasets. This artifact tends to be  
229 most prominent among rare operational taxonomic units (OTUs) or singletons,  
230 which are sometimes misinterpreted as novel diversity, which further causes the  
231 complication of microbial diversity analyses (Gonzalez, Zimmermann, & Saiz-

232 Jimenez, 2004). Moreover, the likelihood of chimera formation has been found to  
233 vary with the GC content of target sequences, with lower GC content generally  
234 associated with a reduced rate of chimera generation (Qin et al., 2023).

### 235 **2.2.1 Effects of Chimeric Reads on Organelle Genome As-** 236 **sembly**

237 In mitochondrial DNA (mtDNA) assembly workflows, PCR-induced chimeras pose  
238 additional challenges. Assembly tools such as GetOrganelle and MitoBeam, which  
239 operate under the assumption of organelle genome circularity, are vulnerable when  
240 chimeric reads disrupt this circular structure. Such disruptions can lead to assem-  
241 bly errors or misassemblies (Bi et al., 2024). These artificial sequences interfere  
242 with the assembly graph, which makes it more difficult to accurately reconstruct  
243 mitochondrial genomes. In addition, these artifacts propagate false variants and  
244 erroneous annotations in genomic data. Hence, determining and minimizing PCR-  
245 induced chimera formation is vital for improving the quality of mitochondrial  
246 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 microbial community studies. These methods rely on the comparison of each query sequence against a curated, high-quality database of known, non-chimeric reference sequences to determine whether the query can be more plausibly explained as a composite or a mosaic of two or more reference sequences rather than as a genuine biological variant (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. Instead of comparing each query sequence to a curated collection of known, non-chimeric sequences, de novo methods infer chimeras based on internal relationships among the sequences present within the dataset itself. This approach is particularly advantageous in studies of novel, under explored, or taxonomically diverse microbial communities where comprehensive reference databases are unavailable or incomplete (Edgar, 2016; Edgar et al., 2011). The underlying assumption on this method operates on the key biological principle that true biological sequences are generally more abundant than chimeric artifacts. During PCR amplification, authentic sequences are amplified early and tend to dominate the read pool, while

271 chimeric sequences form later resulting in the tendency to appear at lower relative  
272 abundances compared to their true parental sequences. As such, the abundance  
273 hierarchy is formed by treating the most abundant sequences as supposed parents  
274 and testing whether less abundant sequences can be reconstructed as mosaics of  
275 these dominant templates. In addition to abundance, de novo algorithms assess  
276 compositional and structural similarity among sequences, examining whether cer-  
277 tain regions of a candidate sequence align more closely with one high-abundance  
278 sequence and other regions with a different one.

279 Both reference-based and de novo approaches are complementary rather than  
280 mutually exclusive. Reference-based methods provide stability and reproducibility  
281 when curated databases are available, whereas de novo methods offer flexibility  
282 and independence for novel or highly diverse communities. In practice, many  
283 modern bioinformatics pipelines combine both paradigms sequentially: an initial  
284 de novo step identifies dataset-specific chimeras, followed by a reference-based pass  
285 that removes remaining artifacts relative to established databases (Edgar, 2016).  
286 These two methods of detection form the foundation of tools such as UCHIME  
287 and later UCHIME2, exemplified by the dual capability of providing both modes  
288 within a unified computational framework.

### 289 **2.3.1 UCHIME**

290 Developed by Edgar et al. (Edgar et al., 2011), UCHIME is one of the most widely  
291 used computational tools for detecting chimeric sequences in amplicon sequencing  
292 data. The UCHIME algorithm detects chimeras by evaluating how well a query  
293 sequence (Q) can be explained as a mosaic of two parent sequences (A and B)

294 from a reference database. The query sequence is first divided into four non-  
295 overlapping segments or chunks. Each chunk is independently searched against a  
296 reference database that is assumed to be free of chimeras. The best matches to  
297 each segment are collected, and from these results, two candidate parent sequences  
298 are identified, typically the two sequences that best explain all chunks of the query.  
299 Then a three-way alignment among the query (Q) and the two parent candidates  
300 (A and B) is done. From this alignment, UCHIME attempts to find a chimeric  
301 model (M) which is a hypothetical recombinant sequence formed by concatenating  
302 fragments from A and B that best match the observed Q

### 303 **Chimeric Alignment and Scoring**

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

314 In de novo mode, UCHIME applies an abundance-driven strategy. Only se-  
315 quences at least twice as abundant as the query are considered as potential parents.  
316 Non-chimeric sequences identified at each step are added iteratively to a growing

317 internal database for subsequent queries.

## 318 **Limitations of UCHIME**

319 Although UCHIME was a significant advancement in chimera detection, it has  
320 notable limitations. According to (Edgar, 2016) and the UCHIME practical notes  
321 (Edgar, n.d), many of the accuracy results reported in the original 2011 paper  
322 were overly optimistic due to unrealistic benchmark designs that assumed com-  
323 plete reference coverage and perfect sequence quality. In practice, UCHIME’s  
324 accuracy can decline when (1) the reference database is incomplete or contains  
325 erroneous entries; (2) low-divergence chimeras are present, as these closely resem-  
326 ble genuine biological variants; (3) sequence datasets include residual sequencing  
327 errors, leading to spurious alignments or misidentification; and (4) the abundance  
328 ratio between parent and chimera is distorted by amplification bias. Additionally,  
329 UCHIME tends to misclassify sequences as non-chimeric when parent sequences  
330 are missing from the database. These limitations motivated the development of  
331 UCHIME2.

### 332 **2.3.2 UCHIME2**

333 To overcome the limitations of its predecessor, UCHIME2 (Edgar, 2016) intro-  
334 duced several methodological and algorithmic refinements that significantly en-  
335 hanced the accuracy and reliability of chimera detection. One major improve-  
336 ment lies in its approach to uncertainty handling. In earlier versions, sequences  
337 with limited reference support were often incorrectly classified as non-chimeric,

338 increasing the likelihood of false negatives. UCHIME2 addresses this issue by  
339 designating such ambiguous sequences as “unknown,” thereby providing a more  
340 conservative and reliable classification framework.

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

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

### 358 2.3.3 CATCh

359 Early chimera detection programs such as UCHIME (Edgar et al., 2011) relied on  
360 alignment-based and abundance-based heuristics to identify hybrid sequences in  
361 amplicon data. However, researchers soon observed that different algorithms often  
362 produced inconsistent predictions. A sequence might be identified as chimeric by  
363 one tool but classified as non-chimeric by another, resulting in unreliable filtering  
364 outcomes across studies.

365 To address these inconsistencies, (Mysara, Saeys, Leys, Raes, & Monsieurs,  
366 2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-  
367 resents the first ensemble machine learning system designed for chimera detection  
368 in 16S rRNA amplicon sequencing. Rather than depending on a single detec-  
369 tion strategy, CATCh integrates the outputs of several established tools, includ-  
370 ing UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus. The individual  
371 scores and binary decisions generated by these tools are used as input features for  
372 a supervised learning model. The algorithm employs a Support Vector Machine  
373 (SVM) with a Pearson VII Universal Kernel (PUK) to determine optimal weight-  
374 ings among the input features and to assign each sequence a probability of being  
375 chimeric.

376 Benchmarking in both reference-based and de novo modes demonstrated signif-  
377 icant performance improvements. CATCh achieved sensitivities of approximately  
378 85 percent in reference-based mode and 92 percent in de novo mode, with corre-  
379 sponding specificities of approximately 96 percent and 95 percent. These results  
380 indicate that CATCh detected 7 to 12 percent more chimeras than any individual  
381 algorithm while maintaining high precision. Integration of CATCh into amplicon-

382 processing pipelines also reduced operational taxonomic unit (OTU) inflation by  
383 23 to 35 percent, producing diversity estimates that more closely reflected true  
384 community composition.

### 385 **2.3.4 ChimPipe**

386 Among the available tools for chimera detection, ChimPipe is a bioinformat-  
387 ics pipeline developed to identify chimeric sequences such as fusion genes and  
388 transcription-induced chimeras from paired-end RNA sequencing data. It uses  
389 both discordant paired-end reads and split-read alignments to improve the ac-  
390 curacy and sensitivity of detecting fusion genes, trans-splicing events, and read-  
391 through transcripts (Rodriguez-Martin et al., 2017). By combining these two  
392 sources of information, ChimPipe achieves better precision than methods that  
393 depend on a single type of signal.

394 The pipeline works with many eukaryotic species that have available genome  
395 and annotation data, making it a versatile tool for studying chimera evolution  
396 and transcriptome structure (Rodriguez-Martin et al., 2017). It can also predict  
397 multiple isoforms for each gene pair and identify breakpoint coordinates that are  
398 useful for reconstructing and verifying chimeric transcripts. Tests using both  
399 simulated and real datasets have shown that ChimPipe maintains high accuracy  
400 and reliable performance.

401 ChimPipe’s modular design lets users adjust parameters to fit different se-  
402 quencing protocols or organism characteristics. Experimental results have con-  
403 firmed that many chimeric transcripts detected by the tool correspond to func-

404 tional fusion proteins, showing its value for understanding chimera biology and  
405 its potential applications in disease research (Rodriguez-Martin et al., 2017).

## 406 **2.4 Machine Learning Approaches for Chimera** 407 **and Sequence Quality Detection**

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

### 419 **2.4.1 Feature-Based Representations of Genomic Se-** 420 **quences**

421 In genomic analysis, feature extraction converts DNA sequences into numerical  
422 representations suitable for ML algorithms. A common approach is k-mer fre-  
423 quency analysis, where normalized k-mer counts form the feature vector (Vervier,



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, 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 signatures 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 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.

449 A further biologically grounded descriptor is micro-homology length at puta-  
450 tive junctions. Micro-homology refers to short, shared sequences (often in the  
451 range of a few to tens of base pairs) that are near breakpoints and mediate  
452 non-canonical repair or template-switch mechanisms. Studies of double strand  
453 break repair and structural variation have demonstrated that the length of micro-  
454 homology correlates with the likelihood of micro-homology-mediated end joining  
455 (MMEJ) or fork-stalled template-switching pathways (Sfeir & Symington, 2015).  
456 In the context of PCR-induced chimeras, template switching during amplifica-  
457 tion often leaves short identical sequences at the junction of two concatenated  
458 fragments. Quantifying the longest exact suffix-prefix overlap at each candidate  
459 breakpoint thus provides a mechanistic signature of chimerism and complements  
460 both compositional (k-mer) and alignment (SA count) features.

## 461 2.5 Synthesis of Chimera Detection Approaches

462 To provide an integrated overview of the literature discussed in this chapter, Ta-  
463 ble 2.1 summarizes the major chimera detection studies, their methodological  
464 approaches, and their known limitations. This consolidated comparison brings to-  
465 gether reference-based approaches, de novo strategies, alignment-driven tools, en-  
466 semble machine-learning systems, and general ML-based sequence-quality frame-  
467 works. Presenting these methods side-by-side clarifies their performance bound-  
468 aries and highlights the unresolved challenges that persist in mitochondrial genome  
469 analysis and chimera detection.

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   |

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

## Chapter 3

# Research Methodology

### 3.1 Research Activities

As illustrated in Figure 3.1, this study will carry out a sequence of computational procedures designed to detect PCR-induced chimeric reads in mitochondrial genomes. The process begins with the collection of mitochondrial reference sequences from the NCBI database, which will serve as the foundation for generating simulated chimeric reads. These datasets will then undergo bioinformatics pipeline development, which includes alignment, k-mer extraction, and homology-based filtering to prepare the data for model construction. The machine-learning model will subsequently be trained and tested using the processed datasets to assess its accuracy and reliability. Depending on the evaluation results, the model will either be refined and retrained to improve performance or, if the metrics meet the desired threshold, deployed for further validation and application.

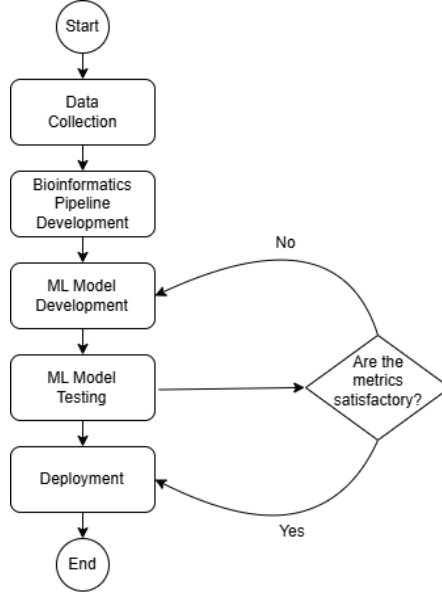


Figure 3.1: Process Diagram of Special Project

### 3.1.1 Data Collection

The mitochondrial genome reference sequences of *Sardinella lemuru* will be collected from the National Center for Biotechnology Information (NCBI) database. The downloaded files will be in FASTA format to ensure compatibility with bioinformatics tools and subsequent analysis. The gathered sequences will serve as the basis for generating simulated chimeric reads to be used in model development.

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

## 509 Data Preprocessing

510 Sequencing data will be simulated using the reference sequences collected from  
511 NCBI. Using `wgsim`, a total of 10,000 paired-end reads (R1 and R2) will be gen-  
512 erated from the reference genome and designated as clean reads. These reads will  
513 be saved in FASTQ (`.fastq`) format. From the same reference, a Bash script will  
514 be created to deliberately cut and reconnect portions of the sequence, introducing  
515 artificial junctions that mimic chimeric regions. The manipulated reference file,  
516 saved in FASTA (`.fasta`) format, will then be processed in `wgsim` to simulate  
517 an additional 10,000 paired-end chimeric reads, also stored in FASTQ (`.fastq`)  
518 format. The resulting read files will be aligned to the original reference genome  
519 using SAMtools, generating SAM (`.sam`) or BAM (`.bam`) alignment files. During  
520 this alignment process, clean reads will be labeled as “0,” while chimeric reads will  
521 be labeled as “1” in a corresponding CSV (`.tsv`) file. This results in a balanced  
522 dataset with an equal number of clean and chimeric reads. This is important to  
523 prevent model bias and ensure that the machine learning classifiers can learn to  
524 detect chimeras accurately.

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

### 530 3.1.2 Bioinformatics Tools Pipeline

531 A bioinformatics pipeline will be developed and implemented to extract the nec-  
532 essary analytical features. This pipeline will serve as a reproducible and modular  
533 workflow that accepts FASTQ and BAM inputs, processes these through a series  
534 of analytical stages, and outputs tabular feature matrices (TSV) for downstream  
535 machine learning. All scripts will be version-controlled through GitHub, and  
536 computational environments will be standardized using Conda to ensure cross-  
537 platform reproducibility. To promote transparency and replicability, the exact  
538 software versions, parameters, and command-line arguments used in each stage  
539 will be documented. To further ensure correctness and adherence to best practices,  
540 bioinformatics experts at the Philippine Genome Center Visayas will be consulted  
541 to validate the pipeline design, feature extraction logic, and overall data integrity.  
542 This stage of the study is scheduled to begin in the last week of November 2025  
543 and conclude by the last week of January 2026, with an estimated total duration  
544 of approximately two (2) months.

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



## 550 Alignment and Supplementary Alignment Count

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

## 564 K-mer Composition Difference

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

573 5 or 6) for each segment. The resulting k-mer frequency vectors will be normalized  
574 and compared using distance metrics such as cosine similarity or Jensen–Shannon  
575 divergence to quantify compositional disparity between the two halves of the same  
576 read. The resulting difference scores will be stored in a structured TSV file.

## 577 **Micro-homology Length**

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

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

### 595 3.1.3 Machine-Learning Model Development

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

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

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

### 617 **3.1.4 Validation and Testing**

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

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

### 632 **3.1.5 Documentation**

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

639 Version control and collaborative management will be implemented through  
640 GitHub, which will serve as the central repository for all project files, including  
641 Python scripts, configuration settings, and Jupyter notebooks. The repository  
642 structure will follow standard research data management practices, with clear  
643 directories for datasets, processed outputs, and analysis scripts. Changes will be  
644 tracked through commit histories to ensure traceability and accountability.

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

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

## 656 **3.2 Calendar of Activities**

657 Table 3.1 presents the project timeline in the form of a Gantt chart, where each  
658 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
- Bi, C., Shen, F., Han, F., Qu, Y., Hou, J., Xu, K., ... Yin, T. (2024, 01).  
Pmat: an efficient plant mitogenome assembly toolkit using low-coverage  
hifi sequencing data. *Horticulture Research*, 11(3), uhae023. Retrieved  
from <https://doi.org/10.1093/hr/uhae023> doi: 10.1093/hr/uhae023
- 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](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

Frenkel-Morgenstern, M., Lacroix, V., Ezkurdia, I., Levin, Y., Gabashvili, A., Prilusky, J., ... Valencia, A. (2012, 05). Chimeras taking shape: Potential functions of proteins encoded by chimeric rna transcripts. *Genome research*, 22, 1231-42. doi: 10.1101/gr.130062.111

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



704 Hahn, C., Bachmann, L., & Chevreux, B. (2013). Reconstructing mitochondrial  
705 genomes directly from genomic next-generation sequencing reads—a baiting  
706 and iterative mapping approach. *Nucleic Acids Research*, *41*(13), e129. doi:  
707 10.1093/nar/gkt371

708 Jin, J.-J., Yu, W.-B., Yang, J., Song, Y., dePamphilis, C. W., Yi, T.-S., & Li,  
709 D.-Z. (2020). Getorganelle: a fast and versatile toolkit for accurate de  
710 novo assembly of organelle genomes. *Genome Biology*, *21*(1), 241. doi:  
711 10.1186/s13059-020-02154-5

712 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-  
713 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),  
714 1819–1825. doi: 10.1093/nar/26.7.1819

715 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.  
716 (2021). Mitochondrial dna reveals genetically structured haplogroups of  
717 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*  
718 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

719 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework  
720 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-  
721 -15-6-r84

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

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

729 Metzker, M. L. (2010). Sequencing technologies — the next generation. *Nature*

730 *Reviews Genetics*, 11(1), 31–46. doi: 10.1038/nrg2626

731 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,  
 732 an ensemble classifier for chimera detection in 16s rna sequencing stud-  
 733 ies. *Applied and Environmental Microbiology*, 81(5), 1573-1584. Retrieved  
 734 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:  
 735 10.1128/AEM.02896-14

736 Qin, Y., Wu, L., Zhang, Q., Wen, C., Nostrand, J. D. V., Ning, D., ... Zhou, J.  
 737 (2023). Effects of error, chimera, bias, and gc content on the accuracy of  
 738 amplicon sequencing. *mSystems*, 8(6), e01025-23. Retrieved from [https://](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23)  
 739 [journals.asm.org/doi/abs/10.1128/msystems.01025-23](https://journals.asm.org/doi/abs/10.1128/msystems.01025-23) doi: 10.1128/  
 740 msystems.01025-23

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

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

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

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

755 Sedlazeck, F., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., von Haeseler,

756 A., & Schatz, M. (2018, 06). Accurate detection of complex structural  
 757 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10  
 758 .1038/s41592-018-0001-7  
 759 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A  
 760 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*  
 761 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)  
 762 [.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)  
 763 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)  
 764 Vervier, M. P. T. M. V. J. B. . V. J. P., K. (2015). Large-scale machine learning  
 765 for metagenomics sequence classification. *Bioinformatics*, 32, 1023 - 1032.  
 766 Retrieved from <https://api.semanticscholar.org/CorpusID:9863600>  
 767 Willette, D., Bognot, E., Mutia, M. T., & Santos, M. (2011). *Biology and ecology*  
 768 *of sardines in the philippines: A review* (Vol. 13; Tech. Rep. No. 1). NFRDI  
 769 Technical Paper Series. Retrieved from [https://nfrdi.da.gov.ph/tpjf/](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)  
 770 [etc/Willette%20et%20al.%20Sardines%20Review.pdf](https://nfrdi.da.gov.ph/tpjf/etc/Willette%20et%20al.%20Sardines%20Review.pdf)