

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

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

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

13 by

14 Duranne Duran
15 Yvonne Lin
16 Daniella Pailden

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

19 December 5, 2025

Contents

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

32	2.2	PCR Amplification and Chimera Formation	9
33	2.3	Existing Traditional Approaches for Chimera Detection	10
34	2.3.1	UCHIME	11
35	2.3.2	UCHIME2	12
36	2.3.3	CATch	13
37	2.3.4	ChimPipe	14
38	2.4	Machine Learning Approaches for Chimera and Sequence Quality	
39		Detection	15
40	2.4.1	Feature-Based Representations of Genomic Sequences . . .	16
41	2.5	Synthesis of Chimera Detection Approaches	18
42	3	Research Methodology	21
43	3.1	Research Activities	21
44	3.1.1	Data Collection	22
45	3.1.2	Bioinformatics Tools Pipeline	26
46	3.1.3	Machine Learning Model Development	29
47	3.1.4	Model Benchmarking, Hyperparameter Optimization, and	
48		Evaluation	30
49	3.1.5	Feature Importance and Interpretation	31

50	3.1.6 Validation and Testing	32
51	3.1.7 Documentation	33
52	3.2 Calendar of Activities	34

53 List of Figures

<small>54</small>	3.1 Process Diagram of Special Project	22
-------------------	--	----

55 List of Tables

56	2.1 Comparison of Chimera Detection Methods	19
57	3.1 Timetable of Activities	34

Chapter 1

Introduction

1.1 Overview

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

PCR chimeras form when incomplete extension products from one template

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

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

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

1.2 Problem Statement

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

117 1.3 Research Objectives

118 1.3.1 General Objective

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

123 1.3.2 Specific Objectives

124 Specifically, the study aims to:

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

1.4 Scope and Limitations of the Research

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

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

159 1.5 Significance of the Research

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

170 Chapter 2

171 Review of Related Literature

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

178 2.1 The Mitochondrial Genome

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

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

189 **2.1.1 Mitochondrial Genome Assembly**

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

2.2 PCR Amplification and Chimera Formation

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

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

2.3 Existing Traditional Approaches for Chimera Detection

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

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

248 sequence correspond to distinct high-abundance sequences.

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

254 **2.3.1 UCHIME**

255 UCHIME is one of the most widely used computational tools for detecting chimeric
256 sequences in amplicon sequencing data, as it serves as a critical quality control
257 step to prevent the misinterpretation of PCR artifacts as novel biological diversity.
258 The algorithm operates by searching for a model (M) where a query (Q) sequence
259 can be perfectly explained as a combination of two parent sequences, denoted as
260 A and B (Edgar et al., 2011).

261 In reference mode, UCHIME divides the query into four chunks and maps
262 them to a trusted chimeric-free database to identify candidate parents. It then
263 constructs a three-way alignment to calculate a score based on “votes.” A “Yes”
264 vote indicates the query aligns with parent A in one region and parent B in an-
265 other, while a “No” vote penalizes the score if the query diverges from the expected
266 chimeric model. In de novo mode, the algorithm operationalizes the abundance
267 skew principle described in Section 2.3. Instead of using an external database,
268 UCHIME dynamically treats the sample’s own high-abundance sequences as a
269 reference database, testing if lower-abundance sequences can be reconstructed as

270 mosaics of these internal ancestors (Edgar et al., 2011).

271 Despite its high sensitivity, UCHIME has inherent limitations rooted in
272 sequence divergence and database quality. The algorithm struggles to detect
273 chimeras formed from parents that are very closely related, specifically when the
274 sequence divergence between parents is less than roughly 0.8%, as the signal-to-
275 noise ratio becomes too low to distinguish a crossover event from sequencing error
276 (Edgar et al., 2011). Furthermore, in reference mode, the accuracy is strictly
277 bound by the completeness of the database; if true parents are absent, the tool
278 may fail to identify the chimera or produce false positives. Similarly, the de novo
279 mode relies on the assumption that parents are present and sufficiently more
280 abundant in the sample, which may not hold true in unevenly amplified samples
281 or complex communities.

282 **2.3.2 UCHIME2**

283 Building upon the original algorithm, UCHIME2 was developed to address the
284 nuances of high-resolution amplicon sequencing. A key contribution of the
285 UCHIME2 study was the critical re-evaluation of chimera detection benchmarks.
286 In the UCHIME2 paper (Edgar, 2016) and the UCHIME in practice website
287 (Edgar, n.d), the author has noted that the accuracy results reported in the
288 original UCHIME paper were “highly over-optimistic” because they relied on
289 unrealistic benchmark designs where parent sequences were assumed to be 100%
290 known and present. UCHIME2 introduced more rigorous testing (the CHSIMA
291 benchmark), revealing that “fake models,” where a valid biological sequence
292 perfectly mimics a chimera of two other valid sequences, are far more common

293 than previously assumed. This discovery suggests that error-free detection is
294 impossible in principle (Edgar, 2016). Another notable improvement is the in-
295 troduction of multiple application-specific modes that allow users to tailor the
296 algorithm’s performance to the characteristics of their datasets. The following
297 parameter presets: denoised, balanced, sensitive, specific, and high-confidence,
298 enable researchers to optimize the balance between sensitivity and specificity
299 according to the goals of their analysis.

300 However despite these advancements, the practical application of UCHIME2
301 requires caution. The author explicitly advises against using UCHIME2 as
302 a stand-alone tool in standard OTU clustering or denoising pipelines. Using
303 UCHIME2 as an independent filtering step in these workflows is discouraged, as
304 it often results in significantly higher error rates, increasing both false positives
305 (discarding valid sequences) and false negatives (retaining chimeras) (Edgar,
306 2016).

307 **2.3.3 CATch**

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

313 To address these inconsistencies, Mysara, Saeys, Leys, Raes, and Monsieurs
314 (2015) developed the Classifier for Amplicon Tool Chimeras (CATCh), which rep-

resents 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.

337 The pipeline works with many eukaryotic species that have available genome
338 and annotation data (Rodriguez-Martin et al., 2017). It can also predict multiple
339 isoforms for each gene pair and identify breakpoint coordinates that are useful
340 for reconstructing and verifying chimeric transcripts. Tests using both simulated
341 and real datasets have shown that ChimPipe maintains high accuracy and reliable
342 performance.

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

348 **2.4 Machine Learning Approaches for Chimera** 349 **and Sequence Quality Detection**

350 Traditional chimera detection tools rely primarily on heuristic or alignment-based
351 rules. Recent advances in machine learning (ML) have demonstrated that models
352 trained on sequence-derived features can effectively capture compositional and
353 structural patterns in biological sequences. Although most existing ML systems
354 such as those used for antibiotic resistance prediction, taxonomic classification,
355 or viral identification are not specifically designed for chimera detection, they
356 highlight how data-driven models can outperform similarity-based heuristics by
357 learning intrinsic sequence signatures. In principle, ML frameworks can integrate
358 indicators such as k-mer frequencies, GC-content variation and split-alignment

359 metrics to identify subtle anomalies that may indicate a chimeric origin (Arango
360 et al., 2018; Liang, Bible, Liu, Zou, & Wei, 2020; Ren et al., 2020).

361 **2.4.1 Feature-Based Representations of Genomic Se-** 362 **quences**

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

377 Machine learning models can also leverage alignment-derived features such as
378 the frequency of split alignments, variation in mapping quality, and local cover-
379 age irregularities. Split reads and discordant read pairs are classical indicators
380 of genomic junctions and have been formalized in probabilistic frameworks for
381 structural-variant discovery that integrate multiple evidence types (Layer, Hall, &

382 Quinlan, 2014). Similarly, long-read tools such as Sniffles employ split-alignment
383 and coverage anomalies to accurately localize breakpoints (Sedlazeck et al., 2018).
384 Modern aligners such as Minimap2 (Li, 2018) output supplementary (SA tags) and
385 secondary alignments as well as chaining and alignment-score statistics that can
386 be summarized into quantitative predictors for machine-learning models. These
387 alignment-signal features are particularly relevant to PCR-induced mitochondrial
388 chimeras, where template-switching events produce reads partially matching dis-
389 tinct regions of the same or related genomes. Integrating such cues within a
390 supervised-learning framework enables artifact detection even in datasets lacking
391 complete or perfectly assembled references.

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

404 **2.5 Synthesis of Chimera Detection Approaches**

405 To provide an integrated overview of the literature discussed in this chapter, Ta-
406 ble 2.1 summarizes the major chimera detection studies, their methodological
407 approaches, and their known limitations.

Table 2.1: Comparison of Chimera Detection Methods

Methods	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 initial UCHIME benchmarking; offers multiple sensitivity/specificity modes; more robust with incomplete references; higher sensitivity.	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 RNA-seq; uses discordant paired-end reads and split-alignments; predicts isoforms and breakpoint coordinates.	Designed for RNA-seq, not amplicons; needs high-quality genome and annotation; computationally heavier; limited to organisms with reference genomes.

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

422 Chapter 3

423 Research Methodology

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

429 3.1 Research Activities

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

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

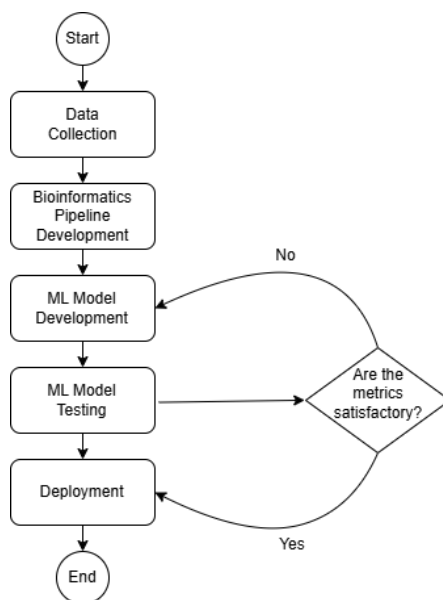


Figure 3.1: Process Diagram of Special Project

442 3.1.1 Data Collection

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

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

448 proximately one (1) week.

449 Data Preprocessing

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

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

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

461 The command parameters are as follows:

- 462 • `-1` and `-2`: read lengths of 150 base pairs for each paired-end read.
- 463 • `-r`, `-R`, `-X`: mutation rate, fraction of indels, and indel extension probability,
464 all set to a default value of 0.
- 465 • `-e`: base error rate, set to 0.001 to simulate realistic sequencing errors.
- 466 • `-N`: number of read pairs, set to 10,000.

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

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

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

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

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

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

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

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

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

```
492 samtools view -bS clean.sam -o clean.bam
493 samtools view -bS chimeric.sam -o chimeric.bam
494
495 samtools sort clean.bam -o clean.sorted.bam
496 samtools index clean.sorted.bam
497
498 samtools sort chimeric.bam -o chimeric.sorted.bam
499 samtools index chimeric.sorted.bam
```

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

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

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

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

513 **3.1.2 Bioinformatics Tools Pipeline**

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

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

529 **Supplementary Alignment Flag**

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

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

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

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

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

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

546 **K-mer Composition Difference**

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

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

Microhomology Length

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

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

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

length and GC content of microhomology.

3.1.3 Machine Learning Model Development

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

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

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

3.1.4 Model Benchmarking, Hyperparameter Optimization, and Evaluation

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

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

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

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

3.1.5 Feature Importance and Interpretation

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

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

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

3.1.6 Validation and Testing

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

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

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

693 3.2 Calendar of Activities

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

Table 3.1: Timetable of Activities

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

References

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

715 45(4), e18. doi: 10.1093/nar/gkw955

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

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

718 CorpusID:88955007

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

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

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

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

723 27(16), 2194–2200. doi: 10.1093/bioinformatics/btr381

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

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

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

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

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

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

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

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

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

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

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

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

736 10.1093/nar/gkt371

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

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

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

740 10.1186/s13059-020-02154-5

741 Judo, M. S. B., Wedel, W. R., & Wilson, B. H. (1998). Stimulation and sup-
742 pression of pcr-mediated recombination. *Nucleic Acids Research*, *26*(7),
743 1819–1825. doi: 10.1093/nar/26.7.1819

744 Labrador, K., Agmata, A., Palermo, J. D., Ravago-Gotanco, R., & Pante, M. J.
745 (2021). Mitochondrial dna reveals genetically structured haplogroups of
746 bali sardinella (*sardinella lemuru*) in philippine waters. *Regional Studies in*
747 *Marine Science*, *41*, 101588. doi: 10.1016/j.rsma.2020.101588

748 Layer, R., Hall, I., & Quinlan, A. (2014, 10). Lumpy: A probabilistic framework
749 for structural variant discovery. *Genome Biology*, *15*. doi: 10.1186/gb-2014-
750 -15-6-r84

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

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

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

760 Mysara, M., Saeys, Y., Leys, N., Raes, J., & Monsieurs, P. (2015). Catch,
761 an ensemble classifier for chimera detection in 16s rna sequencing stud-
762 ies. *Applied and Environmental Microbiology*, *81*(5), 1573-1584. Retrieved
763 from <https://journals.asm.org/doi/abs/10.1128/aem.02896-14> doi:
764 [10.1128/AEM.02896-14](https://journals.asm.org/doi/abs/10.1128/aem.02896-14)

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

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

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

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

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

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

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

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

793 .1038/s41592-018-0001-7

794 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
 795 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
 796 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
 797 [.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)
 798 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

799 Vervier, K., Mahé, P., Tournoud, M., Veyrieras, J.-B., & Vert, J.-P. (2015,
 800 11). Large-scale machine learning for metagenomics sequence classifica-
 801 tion. *Bioinformatics*, 32(7), 1023-1032. Retrieved from [https://doi.org/](https://doi.org/10.1093/bioinformatics/btv683)
 802 [10.1093/bioinformatics/btv683](https://doi.org/10.1093/bioinformatics/btv683) doi: 10.1093/bioinformatics/btv683

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