

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 26, 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	2.2.1	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	3	Research Methodology	25
44	3.1	Research Activities	25
45	3.1.1	Data Collection	26
46	3.1.2	Data Simulation	27
47	3.1.3	Bioinformatics Tools Pipeline	28
48	3.1.4	Machine-Learning Model Development	31

49	3.1.5 Validation and Testing	32
50	3.1.6 Documentation	32
51	3.2 Calendar of Activities	33

52 List of Figures

<small>53</small>	3.1 Process Diagram of Special Project	26
-------------------	--------------------------------------------------	----

54 List of Tables

<small>55</small>	2.1 Summary of Existing Methods and Research Gaps	21
<small>56</small>	3.1 Timetable of Activities	33

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

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

98 1.2 Problem Statement

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

123 1.3.2 Specific Objectives

124 Specifically, the study aims to:

- 125 1. construct empirical and simulated *Sardinella lemuru* Illumina paired-end
126 datasets containing 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 the most informative indicators
133 of PCR-induced 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 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.

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 conti-
163 guity and correctness of *Sardinella lemuru* mitochondrial assemblies. Second, it
164 replaces ad hoc manual curation with a documented, data-driven workflow, im-
165 proving automation and reproducibility. Third, the pipeline is designed to run
166 on modest computing infrastructures commonly available in regional laborato-
167 ries, enabling routine use at the Philippine Genome Center Visayas. In addition,
168 MitoChime supports local capacity building by providing a concrete example of
169 machine-learning integration into mitochondrial genome analysis, consistent with
170 the mandate of the Philippine Genome Center Visayas. Finally, more reliable
171 mitochondrial assemblies for *S. lemuru* provide a stronger basis for downstream
172 applications such as fisheries management, population genetics, and biodiversity
173 assessments.

174 Chapter 2

175 Review of Related Literature

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

182 2.1 The Mitochondrial Genome

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

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

196 **2.1.1 Mitochondrial Genome Assembly**

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

211 search.

212 **2.2 PCR Amplification and Chimera Formation**

213 Polymerase Chain Reaction (PCR) plays an important role in next-generation
214 sequencing (NGS) library preparation, as it amplifies target DNA fragments for
215 downstream analysis. However, the amplification process can also introduce arti-
216 facts that affect data accuracy, one of them being the formation of chimeric se-
217 quences. Chimeras typically arise when incomplete extension occurs during a PCR
218 cycle. This causes the DNA polymerase to switch from one template to another
219 and generate hybrid recombinant molecules (Judo et al., 1998). Artificial chimeras
220 are produced through such amplification errors, whereas biological chimeras oc-
221 cur naturally through genomic rearrangements or transcriptional events. These
222 biological chimeras can have functional roles and may encode tissue-specific novel
223 proteins that link to cellular processes or diseases (Frenkel-Morgenstern et al.,
224 2012).

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

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

236 **2.2.1 Effects of Chimeric Reads on Organelle Genome As-** 237 **sembly**

238 In mitochondrial DNA (mtDNA) assembly workflows, PCR-induced chimeras pose
239 additional challenges. Assembly tools such as GetOrganelle and MitoBeam, which
240 operate under the assumption of organelle genome circularity, are vulnerable when
241 chimeric reads disrupt this circular structure. Such disruptions can lead to assem-
242 bly errors or misassemblies (Bi et al., 2024). These artificial sequences interfere
243 with the assembly graph, which makes it more difficult to accurately reconstruct
244 mitochondrial genomes. In addition, these artifacts propagate false variants and
245 erroneous annotations in genomic data. Hence, determining and minimizing PCR-
246 induced chimera formation is vital for improving the quality of mitochondrial
247 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

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

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

290 **2.3.1 UCHIME**

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

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

304 **Chimeric Alignment and Scoring**

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

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

318 internal database for subsequent queries.

319 **Limitations of UCHIME**

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

333 **2.3.2 UCHIME2**

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

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

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

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

359 2.3.3 CATch

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

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

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

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

386 **2.3.4 ChimPipe**

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

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

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

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

407 **2.4 Machine Learning Approaches for Chimera** 408 **and Sequence Quality Detection**

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

420 **2.4.1 Feature-Based Representations of Genomic Se-** 421 **quences**

422 In genomic analysis, feature extraction converts DNA sequences into numerical
423 representations suitable for ML algorithms. A common approach is k-mer fre-
424 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.

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

462 2.5 Synthesis of Chimera Detection Approaches

463 To provide an integrated overview of the literature discussed in this chapter, Ta-
464 ble 2.1 summarizes the major chimera detection studies, their methodological
465 approaches, and their known limitations. This consolidated comparison brings to-
466 gether reference-based approaches, de novo strategies, alignment-driven tools, en-
467 semble machine-learning systems, and general ML-based sequence-quality frame-
468 works. Presenting these methods side-by-side clarifies their performance bound-
469 aries and highlights the unresolved challenges that persist in mitochondrial genome
470 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.

Method/Study	Scope/Approach	Limitations
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.

Method/Study	Scope/Approach	Limitations
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.
Machine-Learning Sequence Quality & Chimera Detection (general)	Uses k-mer profiles, GC content shifts, entropy, split-read statistics, mapping quality variation, and micro-homology signatures as predictive features; identifies subtle artifacts missed by heuristics.	Requires labeled training data; model performance depends on feature engineering; may capture dataset-specific biases; limited generalization if training data is narrow or unrepresentative.

Across existing studies, no single approach reliably detects all forms of chimeric sequences, particularly those generated by PCR-induced template switching in mitochondrial genomes. Reference-based tools perform poorly when parental sequences are absent; de novo methods rely strongly on abundance assumptions; alignment-based systems show reduced sensitivity to low-divergence chimeras; and ensemble methods inherit the limitations of their component algorithms. RNA-seq-oriented pipelines likewise do not generalize well to organelle data. Although machine learning approaches offer promising feature-based detection, they are rarely applied to mitochondrial genomes and are not trained specifically on PCR-

480 induced organelle chimeras. These limitations indicate a clear research gap: the
481 need for a specialized, feature-driven classifier tailored to mitochondrial PCR-
482 induced chimeras that integrates k-mer composition, split-alignment signals, and
483 micro-homology features to achieve more accurate detection than current heuristic
484 or alignment-based tools.

Chapter 3

Research Methodology

This chapter outlines and explains the specific steps and activities to be carried out in completing the project.

3.1 Research Activities

As illustrated in Figure 3.1, the researchers 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

499 will either be refined and retrained to improve performance or, if the metrics meet
500 the desired threshold, deployed for further validation and application.

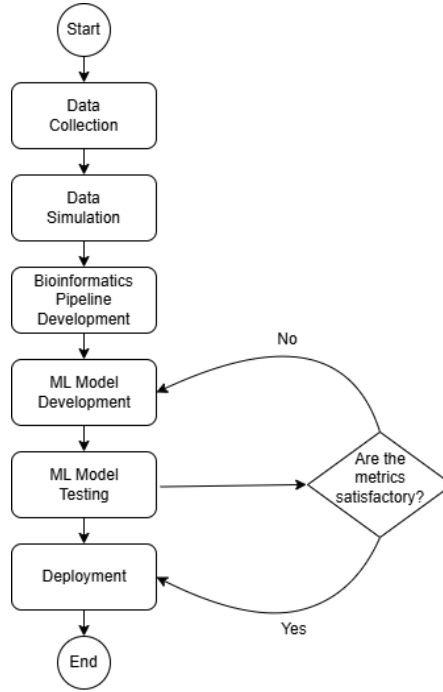


Figure 3.1: Process Diagram of Special Project

501 3.1.1 Data Collection

502 The researchers will collect mitochondrial genome reference sequences of *Sar-*
503 *dinella lemuru* from the National Center for Biotechnology Information (NCBI)
504 database. The downloaded files will be in FASTA format to ensure compatibility
505 with bioinformatics tools and subsequent analysis. The gathered sequences will
506 serve as the basis for generating simulated chimeric reads to be used in model
507 development.

508 The expected outcome of this process is a comprehensive dataset of *Sardinella*

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

513 3.1.2 Data Simulation

514 The researchers will simulate sequencing data using the reference sequences col-
515 lected from NCBI. Using `wgsim`, a total of 5,000 paired-end reads (R1 and R2)
516 will be generated from the reference genome and designated as clean reads. These
517 reads will be saved in FASTQ (`.fastq`) format. From the same reference, a Bash
518 script will be created to deliberately cut and reconnect portions of the sequence,
519 introducing artificial junctions that mimic chimeric regions. The manipulated
520 reference file, saved in FASTA (`.fasta`) format, will then be processed in `wgsim`
521 to simulate an additional 5,000 paired-end chimeric reads, also stored in FASTQ
522 (`.fastq`) format. The resulting read files will be aligned to the original reference
523 genome using SAMtools, generating SAM (`.sam`) or BAM (`.bam`) alignment files.
524 During this alignment process, clean reads will be labeled as “0,” while chimeric
525 reads will be labeled as “1” in a corresponding CSV (`.tsv`) file.

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

531 3.1.3 Bioinformatics Tools Pipeline

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

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

552 Alignment and Supplementary Alignment Count

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

566 K-mer Composition Difference

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

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

579 **Micro-homology Length**

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

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

598 3.1.4 Machine-Learning Model Development

599 The classification component of MitoChime will employ two ensemble algo-
600 rithms—Random Forest (RF) and Extreme Gradient Boosting (XGBoost)—to
601 evaluate complementary learning paradigms. Random Forest applies bootstrap
602 aggregation (bagging) to reduce model variance and improve stability, whereas
603 XGBoost implements gradient boosting to minimize bias and capture complex
604 non-linear relationships among genomic features. Using both models enables a
605 balanced assessment of predictive performance and interpretability.

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

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

619 **3.1.5 Validation and Testing**

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

628 Comparative evaluation between the Random Forest and XGBoost classifiers
629 will establish which model achieves superior predictive accuracy and computa-
630 tional efficiency under identical data conditions.

631 **3.1.6 Documentation**

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

638 Version control and collaborative management will be implemented through
639 GitHub, which will serve as the central repository for all project files, including
640 Python scripts, configuration settings, and Jupyter notebooks. The repository

641 structure will follow standard research data management practices, with clear
 642 directories for datasets, processed outputs, and analysis scripts. Changes will be
 643 tracked through commit histories to ensure traceability and accountability.

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

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

655 3.2 Calendar of Activities

656 Table 3.1 presents the project timeline in the form of a Gantt chart, where each
 657 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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

755 A., & Schatz, M. (2018, 06). Accurate detection of complex structural
 756 variations using single-molecule sequencing. *Nature Methods*, 15. doi: 10
 757 .1038/s41592-018-0001-7

758 Sfeir, A., & Symington, L. S. (2015). Microhomology-mediated end joining: A
 759 back-up survival mechanism or dedicated pathway? *Trends in Biochemical*
 760 *Sciences*, 40(11), 701-714. Retrieved from [https://www.sciencedirect](https://www.sciencedirect.com/science/article/pii/S0968000415001589)
 761 [.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)
 762 [10.1016/j.tibs.2015.08.006](https://doi.org/10.1016/j.tibs.2015.08.006)

763 Vervier, M. P. T. M. V. J. B. . V. J. P., K. (2015). Large-scale machine learning
 764 for metagenomics sequence classification. *Bioinformatics*, 32, 1023 - 1032.
 765 Retrieved from <https://api.semanticscholar.org/CorpusID:9863600>

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