

# Machine Learning Pipeline for Detecting PCR-Induced Chimeric Reads

MitoChime: Organellar Chimera Detection from Per-Read Features

Duran, Lin, Pailden

University of the Philippines Visayas  
Philippine Genome Center Visayas

December 9, 2025

# Outline

- 1 Introduction
- 2 Problem Statement & Proposed Solution
- 3 Objectives
- 4 Scope and Limitations
- 5 Methodology

## Next Generation Sequencing (NGS)

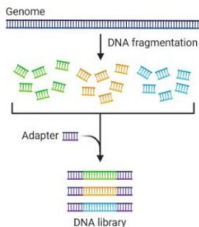


*Source: University of the Philippines  
Visayas, 2022*

## Illumina Seq Workflow

### Step 1. Library Preparation

#### ① Library preparation



Source: *Microbe Notes*, 2024

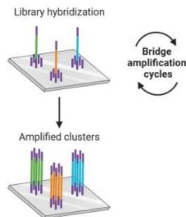


Source: *Philippine Genome Center Visayas*, 2025

## Illumina Seq Workflow

### Step 2. Library Bridge Amplification (PCR)

#### ② DNA library bridge amplification



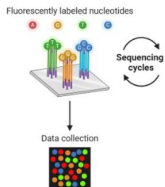
Source: *Microbe Notes*, 2024



Source: Philippine Genome Center  
Visayas, 2025

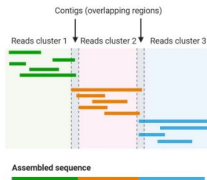
# Illumina Seq Workflow

### Step 3. Sequencing and Alignment



**Example of a short-read:**

@HS001123:45:11576680X:1:11101:12345:1000 1:N:0:ATCG  
 GATTGACTCCATGGTACCGTAATGCGTAGGATCTATGCGTACCGTATG  
 \*  
 AAAAAAAAAA



**Example of an assembled sequence:**

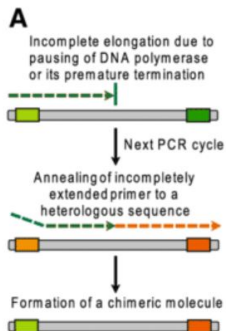
assembled\_mitochondrial\_genome

```

GATCACAGGTTATACCCCTATTACCACTCACGGGAGCTTCATCATG
TGTTGATTTCCTGCTGGGGGAGGACACGCGATAGCATTCGAGACGCTG
AGCCCTCTTAAACACAAAGCAGGAGAACTGATCATGATGACCTTTGGCT
TAGGGTCAGGTAGAGGAGACGCTGGTAGCTCAGTAGGACTTGCTCTTGT
GGCTATTATTCAGGAGCACTACTCTCAACTGAGCTAGTGGTCTATGGA
ACACTGAGCTCATGACACTAGCAGTCTACCGTACCAAGTATGCTGATAG
TCTGGTAGCTCGGCTCTGTCAGTCAAGTAGCTACCAAGTACCTCTACTG

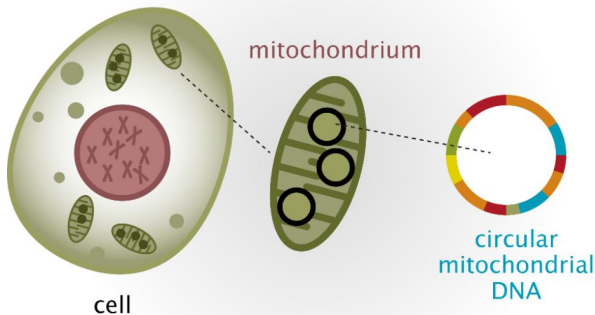
```

## PCR-Chimera Formation



*Source: Omelina et al., 2024*

## The Mitochondrial Genome

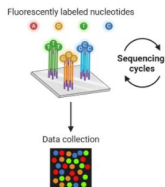


*Source: UZ Brussel., 2020*



## Disrupts Genome Assembly

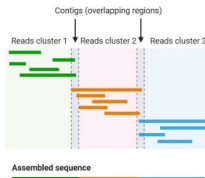
### ③ DNA library sequencing



**Example of a short-read:**

@HIS00123:45:H5TWLBCXX:1:11101:12345:1000 1:N:0:ATCG  
GATTGACTCCATGGTACCGTAATGGCTAGGTATCTATGGGTACCGTATG  
+  
AAAAFFFFF111111111111111111111111111111111111

#### ④ Alignment and data analysis



**Example of an assembled sequence:**

```
>assembled_mitochondrial_genome
GATCACAGGCTCTATCACCTATTAAACCATCAGGGAGCTTCTCCATGCA
TTGGTATTTTTCGTCTGGGGGTTATGCACCGCATGTCATGCGAGACGCTG
AGCCCTCTTAAACCAACAGCGAGAACTCTGATCATGATGACCTTTTGCTC
TAGGGTCAGGTAGAGAGACGCTGTGATCTCAGTAGGACTTGTCTTGT
GGGCTATTTCTCGCAGGACATCTCTCAACCTGAGCCCTATGCTCATGGA
ACACTGAAGCTAGCAGCTAGCAGCTACCTCTACCGTAGACCTTAGTGATAG
TCTCGTAGCTCGCTTCTAGTCAGTTAGTACCTACCTACCTCTACCTCT
```

# Existing Approaches

Table 2.1: Comparison of Chimera Detection Approaches and Tools

Method / Tool	Core Approach	Key Limitations
<b>Reference-based Detection</b>	Compares each query sequence against curated databases of verified, non-chimeric sequences; evaluates segment similarity to identify mosaic patterns.	Accuracy depends on database completeness; performs poorly for novel taxa or missing parents; limited sensitivity for low-divergence chimeras.
<b>De novo Detection</b>	Identifies chimeras using only internal dataset structure; leverages abundance hierarchy and compositional similarity to infer whether low-abundance sequences can be reconstructed from abundant parents.	Assumes true sequences are more abundant; fails when amplification bias distorts abundances; struggles when parental sequences are similarly abundant or highly similar.
<b>UCHIME</b>	Alignment-based model that partitions the query into segments, identifies parent candidates, and computes a chimera score via a three-way alignment; supports reference and de novo modes.	Reduced accuracy for very closely related parents ( $<0.8\%$ divergence); sensitive to incomplete databases; de novo mode fails if parents are absent or not sufficiently more abundant.
<b>UCHIME2</b>	Updated UCHIME with improved benchmarking (CHSIMA) and multiple sensitivity/specificity presets; better handles incomplete references and dataset variability.	"Fake models" limit theoretical accuracy; genuine variants may mimic chimeras; not recommended as a standalone step in OTU or denoising pipelines due to increased false positives/negatives.
<b>CATCh</b>	First ensemble ML model for 16S chimera detection; integrates outputs of UCHIME, ChimeraSlayer, DECIPHER, Pintail, and Perseus using an SVM to boost overall prediction accuracy.	Performance constrained by underlying tools; ML model cannot capture features not present in component algorithms; may misclassify in highly novel or low-coverage datasets.
<b>ChimPipe</b>	Pipeline for detecting biological chimeras in RNA-seq using discordant paired-end reads and split-read alignments; identifies isoforms and breakpoint coordinates.	Requires high-quality genome and annotation; tailored to RNA-seq rather than amplicons; computationally intensive; limited to organisms with available reference genomes.

# Problem Statement & Proposed Solution

- **Problem Statement:** Chimeric sequencing reads can disrupt mitochondrial genome assembly, but current assembly pipelines assume artifact-free input and existing chimera detection tools are not designed specifically for organellar, particularly mitochondrial datasets, leaving assemblies vulnerable to undetected artifacts.
- **Proposed Solution:** A machine-learning pipeline designed to detect PCR-induced chimeric reads using both alignment-based and sequence-derived features to improve the quality and reliability of downstream mitochondrial genome assemblies.

# General Objective

- Develop and evaluate a machine-learning pipeline (MitoChime) to detect PCR-induced chimeric reads in *S. lemur* mitochondrial sequencing data to improve downstream assembly quality.

# Specific Objectives

- 1 Construct simulated *Sardinella lemuru* Illumina paired-end datasets containing both clean and PCR-induced chimeric reads.
- 2 Extract alignment-based and sequence-based features such as k-mer composition, microhomology, and split-alignment counts from both clean and chimeric reads
- 3 Train, validate, and compare supervised machine learning models for classifying reads as clean or chimeric.
- 4 Determine feature importance and identify indicators of PCR-induced chimerism.
- 5 Integrate the optimized classifier into a modular and interpretable pipeline deployable on standard computing environments at PGC Visayas.

# Scope of the Study

- Focuses on PCR-induced chimeric reads in *Sardinella lemuru* mitochondrial sequencing data to:
  - 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
  - to align the analysis with relevant *S. lemuru* sequencing projects at PGC Visayas
  - 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
  - to develop a tool that directly supports local studies on *S. lemuru* population structure and fisheries management produce tools applicable to local population and fisheries studies

# Scope of the Study

- Uses wgsim-based simulations and selected empirical mitochondrial datasets
- Analysis targets low-dimensional alignment and sequence features (k-mers, GC content, clipping, split alignments) to maintain interpretability and computational accessibility
- Long-read platforms and other taxa are not included

# Key Exclusions

- Naturally occurring chimeras
- NUMTs
- Large-scale nuclear genome rearrangements
- High-dimensional deep learning embeddings



# Other Limitations

- No simulations with variable sequencing error rates
- No testing of alternative parameter settings (k-mer length, microhomology windows)
- Reliance on supervised machine learning may limit detection of novel/unknown chimeric patterns

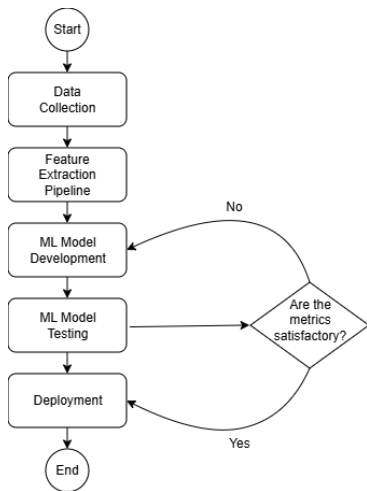


Figure: Process Diagram of the Special Project

The *S. lemur* mitochondrial reference genome (NCBI: NC\_039553.1) was downloaded in FASTA format and used as the basis for generating simulated reads.

- A Python script was used to generate the reads.
- Clean reads were produced with wgsim from the reference genome.
- A chimeric reference was created by creating a custom script to combine non-adjacent segments with microhomology
- Chimeric reads were simulated with wgsim.
- All reads were mapped with minimap2 to extract alignment information.
- SAM/BAM files were converted, sorted, and indexed with samtools.

- Final dataset: 40k reads, roughly balanced between clean and chimeric (19,984 clean reads and 20,000 chimeric).
- Some of the clean reads failed to align due to the set error rate.

# Data Preprocessing

```
NC_039553.1_3_540_8:0:0_6:0:0_ef2 163 NC_039553.1 3 60 150M = 391 538
TGGTGTAGCTTAAACAAGCATAAAGCTGAAGATGTTACGATGGGCGGTGATAAGCCCAACGACCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCGGAATTTACACACCGAGGCTCCGCGGCCGCTGAGGATGGCTCA
..... NM:i:8 ms:i:220
AS:i:220 nn:i:0 tp:A:P cm:i:8 s1:i:164 s2:i:0 de:f:0.0533 rl:i:0
NC_039553.1_4_430_13:0:0_11:0:0_243d 163 NC_039553.1 4 60 150M = 281 427
GGTGTAGCTTAAACAAGCATAAAGCTGAGATGATCCGCTGGGCGGTGATAAGCCGACGAGGAGTGAAGTTTGGTCCAGGCTTTATTATCAGCTTTACCCCAATTTACACATGCGAGCTCCGCGGCCGCTGAGGATGGCTCAG
..... NM:i:13 ms:i:170
AS:i:170 nn:i:0 tp:A:P cm:i:9 s1:i:135 s2:i:0 de:f:0.0867 rl:i:0
NC_039553.1_5_495_6:0:0_11:0:0_1d49 163 NC_039553.1 5 60 150M = 346 491
GTGTAGCTTACACAAGCATAAAGCTGAAGATGTTAAGATGGGCGGTGATCAGCCCAACGACCTGAAAGGTTAGGTCCTGGCTTTATTATCAGGTTTCCCCCAATTTACACATGCGAGCTCCGCGGCCGCTGAGGATGGCTCAGC
..... NM:i:6 ms:i:240
AS:i:240 nn:i:0 tp:A:P cm:i:12 s1:i:148 s2:i:0 de:f:0.04 rl:i:0
NC_039553.1_6_523_6:0:0_9:0:0_82c 163 NC_039553.1 6 60 150M = 374 518
TGTAGCTTAAACAAGCATAAAGCTGAAGATGTTAAGATGGGCGGTGATAAGCCCAACGACCTGCAAGGTTTGGTCTGGCTATATTACAGCTTTACCCCAATTTACACATGCGGCTCCGCGGCCGCTGAGGATGGCTCAGCC
..... NM:i:6 ms:i:240
AS:i:240 nn:i:0 tp:A:P cm:i:10 s1:i:157 s2:i:0 de:f:0.04 rl:i:0
NC_039553.1_9_574_7:0:0_7:0:0_181b 163 NC_039553.1 9 60 150M = 425 566
AGCTTAAACAAGCATAAAGCTGAGATGTTAAGCTGGGCGGTGATAAGCCCAACGACCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCGCAATTTACACATGCGAGCTCCGCGGCCGCTGAGGCTGCCCTCCGCTCC
..... NM:i:7 ms:i:230
AS:i:230 nn:i:0 tp:A:P cm:i:12 s1:i:176 s2:i:0 de:f:0.0467 rl:i:0
NC_039553.1_10_391_9:0:0_8:0:0_256b 99 NC_039553.1 10 60 150M = 242 382
GCTTAAACAAGCATAAAGCTGAGATGTTAAGATGGGCGGTGATAAGCCCAACGACCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCCCAATTTAGACATGCGAGCTCCGCGGCCGCTGATGCTGGCTCAGCTCCC
..... NM:i:9 ms:i:210
AS:i:210 nn:i:0 tp:A:P cm:i:15 s1:i:156 s2:i:0 de:f:0.06 rl:i:0
NC_039553.1_11_509_6:0:0_11:0:0_a19 99 NC_039553.1 11 60 150M = 360 499
CTTCAACAAGCATAAAGCTGAAGATGTTAAGTGGGCGGTATAGCCCGACAAGCCTGAAAGGTTAGGTCCTGGCTTTATTATGAGCTTTACCCCAATTTACACATGCGATCTCCGCGGCCGCTGAGGATGCCCTCAGCTCCCG
..... NM:i:6 ms:i:242
AS:i:242 nn:i:0 tp:A:P cm:i:10 s1:i:150 s2:i:0 de:f:0.04 rl:i:0
NC_039553.1_12_427_9:0:0_9:0:0_157 163 NC_039553.1 12 60 150M = 278 416
TTAAACAAGCATAAAGCTGAAGATTTAGATGGGCGGTGATAAGCCCAACGACCTGAAAGTTTGGTCTGGCTTTATTATCAGCTTTACCCCAATTTACACATGCGAGCTCCGCGGCCGCTGAGGATGCCCTCCGCTCCCGT
..... NM:i:9 ms:i:210
AS:i:210 nn:i:0 tp:A:P cm:i:8 s1:i:150 s2:i:0 de:f:0.06 rl:i:0
```

Figure: SAM File of Clean Reads

# Data Preprocessing

[illegible]

### Figure: SAM File of Chimeric Reads

# Feature Extraction Pipeline

- BAM files were processed with a Python script (`extract_features.py`) to build a TSV feature matrix.
- Used Pysam for parsing alignments and NumPy for computation.



# Feature Extraction Pipeline

- Focused on three features linked to PCR-induced chimeras:
  - ① **Supplementary Alignment (SA)**: Detects split alignments; counts and metrics extracted from SA tags
  - ② **K-mer Composition Difference**: Breakpoints inferred; left/right segments compared using cosine and JS metrics.
  - ③ **Microhomology**: Overlap at junction quantified (length + GC content) within a defined window.
- Pipeline design and outputs to be validated by experts.

# Feature Extraction Pipeline

read_id	label	read_length	mean_baseq	name	ref_start	1strand	mapq	cigar	has_sa	sa_count	num_seg	sa_diff	co	sa_min	dk	sa_max	d	sa_mean	sa_same	sa_opp	st	sa_max	r	sa_mean	sa_min	r	sa_mean	softclip_l	softclip_r	total_clip	breakpoint	kmer	cool	kmer_jc	d	microhom	microhom	
NC_039502	0	150	13	NC_039502	3	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.9726	0.97143	1	0				
NC_039502	0	150	13	NC_039502	4	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98591	0.98571	1	0				
NC_039502	0	150	13	NC_039502	5	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95887	0.95714	0	0				
NC_039502	0	150	13	NC_039502	6	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97183	0.97143	1	1				
NC_039502	0	150	13	NC_039502	9	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98664	0.98571	0	0				
NC_039502	0	150	13	NC_039502	10	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97296	0.97143	0	0				
NC_039502	0	150	13	NC_039502	11	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	0				
NC_039502	0	150	13	NC_039502	12	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	1				
NC_039502	0	150	13	NC_039502	12	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98640	0.98571	1	1				
NC_039502	0	150	13	NC_039502	12	0	24	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95889	0.95714	1	1				
NC_039502	0	150	13	NC_039502	14	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	0				
NC_039502	0	150	13	NC_039502	15	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98649	0.98571	1	0				
NC_039502	0	150	13	NC_039502	17	0	60	148M4S	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	4	148	0	0.5	0	0				
NC_039502	0	150	13	NC_039502	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98649	0.98571	3	0				
NC_039502	0	150	13	NC_039502	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97221	0.97143	3	0				
NC_039502	0	150	13	NC_039502	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98629	0.98571	3	0				
NC_039502	0	150	13	NC_039502	19	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	3	0				
NC_039502	0	150	13	NC_039502	20	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97221	0.97143	0	0				
NC_039502	0	150	13	NC_039502	21	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0				
NC_039502	0	150	13	NC_039502	23	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98667	0.98571	0	0				
NC_039502	0	150	13	NC_039502	25	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98629	0.98571	0	0				
NC_039502	0	150	13	NC_039502	28	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98603	0.98571	1	0				
NC_039502	0	150	13	NC_039502	32	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97258	0.97143	2	1				
NC_039502	0	150	13	NC_039502	34	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	0				
NC_039502	0	150	13	NC_039502	34	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	0	0				
NC_039502	0	150	13	NC_039502	35	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0				
NC_039502	0	150	13	NC_039502	36	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98648	0.98571	0	0				
NC_039502	0	150	13	NC_039502	38	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0				
NC_039502	0	150	13	NC_039502	39	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98684	0.98571	0	0				
NC_039502	0	150	13	NC_039502	41	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97296	0.97143	2	0.5				
NC_039502	0	150	13	NC_039502	43	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	0	0				

Figure: TSV Dataset showing Clean Reads

# Feature Extraction Pipeline

1985	read_1	label	1	read_1	mean	ref_nst	ref_sta	strand	1	cgis	mas_sa	sa_cis	mas_nst	sa_dft	sa_mst	sa_mai	sa_snt	sa_snp	sa_posp	sa_mai	sa_mst	sa_mst	sa_mst	softsc_z	softsc_g	total_c	brnaskp	chr	l	kmr	macro_c	macro_g	kitopg	g		
1985	chmerna_1	1	150	40 NC.03956	40	1	60	150M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.9848	0.9671	0	0					
1986	chmerna_1	1	150	40 NC.03956	53	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	0	0				
1987	chmerna_1	1	150	40 NC.03956	65	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95774	0.95714	0	0					
1988	chmerna_1	1	150	40 NC.03956	65	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95774	0.95714	0	0					
1989	chmerna_1	1	150	40 NC.03956	67	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95774	0.95714	0	0					
1990	chmerna_1	1	150	40 NC.03956	67	1	60	1184325	1	1	2	0	4246	4246	4246	0	1	10	10	0	0	0	0	0	0	32	32	118	1	1	0	0				
1991	chmerna_1	1	150	40 NC.03956	69	1	60	150M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95774	0.95714	0	0					
1992	chmerna_1	1	150	40 NC.03956	71	1	60	1009445	1	1	2	0	4237	4237	4237	0	16	16	0	0	0	0	0	0	0	41	41	150	1	1	0	0				
1993	chmerna_1	1	150	40 NC.03956	77	1	60	150M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.94306	0.84296	0	0					
1994	chmerna_1	1	150	40 NC.03956	79	0	60	1009445	1	1	2	0	4234	4234	4234	0	1	17	17	0	0	0	0	0	0	44	44	106	1	1	0	0				
1995	chmerna_1	1	150	40 NC.03956	84	0	60	112388	1	1	2	0	5197	5197	5197	0	1	10	10	0	0	0	0	0	0	38	38	112	0.9377	0.9648	0	0				
1996	chmerna_1	1	150	40 NC.03956	85	0	60	112388	1	1	2	0	5196	5196	5196	0	1	20	20	0	0	0	0	0	0	99	39	111	0.93634	0.9647	0	0				
1997	chmerna_1	1	150	40 NC.03956	88	0	60	150M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95683	0.95652	1	1	1	0	0			
1998	chmerna_1	1	150	40 NC.03956	89	0	60	155135M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15	0	15	1	1	0	0				
1999	chmerna_1	1	150	40 NC.03956	89	0	60	308120M	0	1	2	0	1973	1973	1973	0	1	1	0	0	0	0	0	0	0	30	30	0.98197	0.96352	0	0					
2000	chmerna_1	1	150	40 NC.03956	89	0	60	415100M	1	1	2	0	1962	1962	1962	0	1	15	15	0	0	0	0	0	0	41	41	0.98411	0.96482	0	0					
2001	chmerna_1	1	150	40 NC.03956	89	1	60	969545	1	1	2	0	4234	4234	4234	0	1	48	48	0	0	0	0	0	0	54	54	1	1	1	0	0				
2002	chmerna_1	1	150	40 NC.03956	89	1	60	355117M	1	1	2	0	1970	1970	1970	0	1	1	1	0	0	0	0	0	0	30	0	30	0.98275	0.96389	0	0				
2003	chmerna_1	1	150	40 NC.03956	90	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0			
2004	chmerna_1	1	150	40 NC.03956	90	1	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0		
2005	chmerna_1	1	150	40 NC.03956	90	1	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0		
2006	chmerna_1	1	150	40 NC.03956	91	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0		
2007	chmerna_1	1	150	40 NC.03956	91	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0		
2008	chmerna_1	1	150	40 NC.03956	91	1	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95832	0.95714	1	1	1	0	0		
2009	chmerna_1	1	150	40 NC.03956	91	0	60	949568	1	1	2	0	4222	4222	4222	0	1	52	52	0	0	0	0	0	0	56	56	94	1	1	1	0	0			
2010	chmerna_1	1	150	40 NC.03956	92	0	60	295121M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	29	29	29	0.98047	0.96338	0	0				
2011	chmerna_1	1	150	40 NC.03956	92	0	60	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	3064	72	72	0.98047	0.96338	1	1	1	0	0			
2012	chmerna_1	1	150	40 NC.03956	92	0	60	66684M	1	1	2	0	3070	3070	3070	0	1	59	59	0	0	0	0	0	0	66	66	0.98611	0.96565	1	0					
2013	chmerna_1	1	150	40 NC.03956	92	0	60	115139M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	11	11	11	1	1	1	1	0	0			
2014	chmerna_1	1	150	40 NC.03956	92	0	60	385134M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	16	16	0.97424	0.96941	1	0					
2015	chmerna_1	1	150	40 NC.03956	92	0	60	352147M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	5	5	5	0.985	0.9653	0	0					
2016	chmerna_1	1	150	40 NC.03956	92	0	60	54599M	1	1	2	0	3082	3082	3082	0	1	30	30	0	0	0	0	0	0	54	54	54	0.9726	0.9653	0	0				

Figure: TSV Dataset showing Chimeric Reads

# Dataset construction and split

- Simulated feature tables:
  - Clean reads (label 0)
  - PCR-induced chimeras (label 1)
- `build_datasets.py`:
  - Concatenate tables
  - Shuffle rows (avoid file-order artefacts)
- 80/20 **stratified** train-test split
- Test set held out and used **only once** at the end

# Validation strategy

- Layer 1: 80/20 stratified train–test split
- Layer 2: 5-fold stratified cross-validation on training set
  - Train on 4 folds, validate on 1
  - Rotate so each fold is validation once
- Layer 3: Final evaluation on held-out test set
- Hyperparameter tuning:
  - RandomizedSearchCV inside CV for top models
- Goal: stable estimates and **unbiased** test performance

# Model zoo and preprocessing pipeline

- **Baseline:** dummy majority-class classifier
- **Linear models:** logistic regression, calibrated linear SVM
- **Tree ensembles:**
  - Random Forest, Extra Trees
  - Gradient Boosting, XGBoost, LightGBM, CatBoost
- **Others:** bagging trees, k-NN, Gaussian NB, shallow MLP
- Common scikit-learn pipeline:
  - Median imputation (numeric missing values)
  - Standardisation (zero mean, unit variance)
- Ensures a **fair comparison** across models

# Effect of hyperparameter tuning

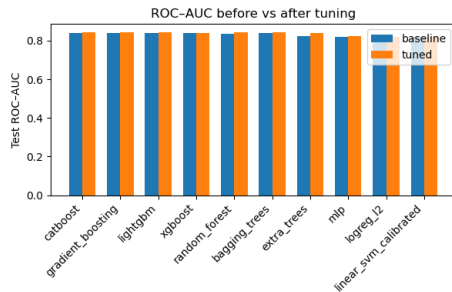
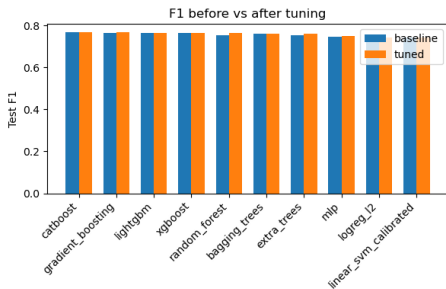
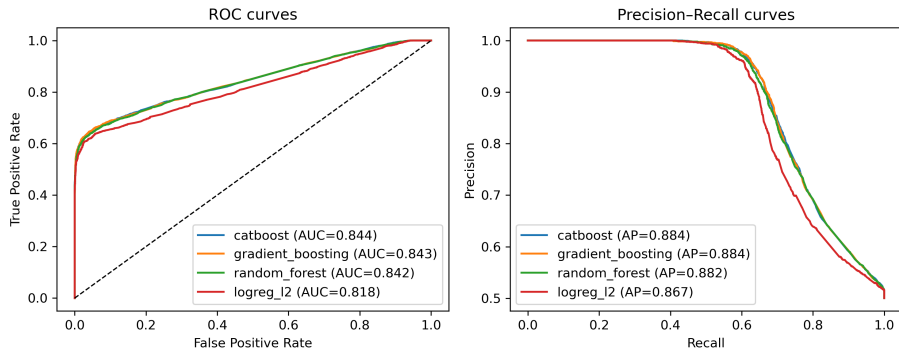


Figure: Test F1: baseline vs tuned.

Figure: Test ROC-AUC: baseline vs tuned.

- Tuning done with RandomizedSearchCV on training set
- Small but consistent gains ( $\Delta F1$ ,  $\Delta AUC \approx 0.001-0.01$ )
- Top-ranked models remain the same (CatBoost, Gradient Boosting, LightGBM)

# ROC and precision–recall curves

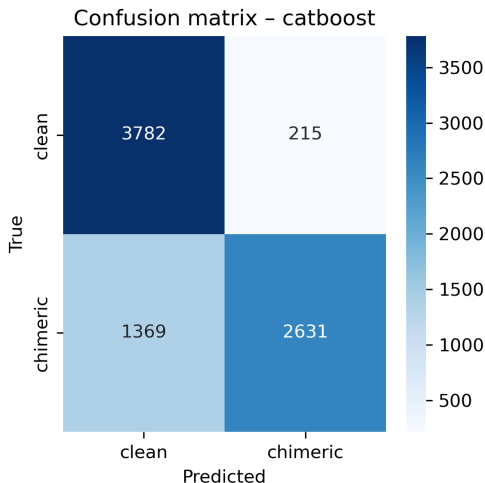


**Figure:** ROC (left) and PR (right) curves for CatBoost, Gradient Boosting, Random Forest, and logistic regression.

- Ensembles: ROC–AUC  $\approx 0.84$ ; logreg:  $\approx 0.82$
- Average precision  $\approx 0.88$  for ensembles
- Precision  $> 0.9$  up to recall  $\approx 0.5$ – $0.6$



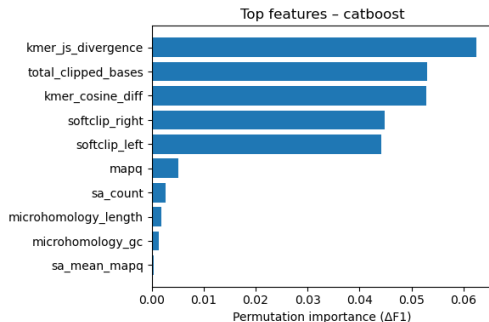
# Confusion matrix: CatBoost (test set)



- Clean reads:
  - Recall  $\approx 0.95$  (3782 / 3997)
- Chimeric reads:
  - Precision  $\approx 0.92$
  - Recall  $\approx 0.66$  (2631 / 4000)
- Behaviour at default threshold:
  - **Conservative chimera filter**
  - Protects clean reads, misses some subtle chimeras

Figure: Confusion matrix heatmap for CatBoost.

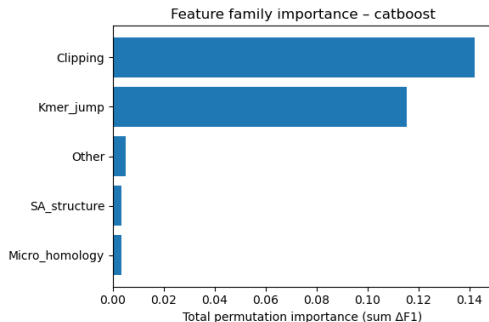
# Top features for CatBoost



**Figure:** Permutation importance ( $\Delta F1$ ) for CatBoost.

- Strongest signals:
  - kmer\_js\_divergence
  - total\_clipped\_bases
  - kmer\_cosine\_diff
- Also important:
  - Left/right soft-clipping
  - Mapping quality (MAPQ)
  - SA count (supplementary alignments)
- Consistent with PCR chimera junctions

# Feature family importance



**Figure:** Aggregated feature families for CatBoost.

- Aggregated permutation importance:
  - **Clipping** features dominate
  - **K-mer jump** features also strong
- Smaller contributions:
  - SA structure
  - Micro-homology
  - Other alignment context
- Same pattern for Gradient Boosting and Random Forest

# Take-Home Messages

- Tree-based ensembles (CatBoost, Gradient Boosting, LightGBM) consistently outperform linear baselines.
- Best models reach  $F1 \approx 0.77$  and  $ROC-AUC \approx 0.84$  on held-out reads.
- Key predictive signals match chimera biology:
  - k-mer composition jumps along the read
  - extensive soft-clipping and total clipped bases
- At the default threshold, the filter is conservative:
  - preserves most clean reads
  - removes a substantial fraction of chimeras
- Overall: our feature set and model ensemble provide a practical pre-filter before mitochondrial assembly.