

# 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 8, 2025

# Outline

1 Objectives

2 Scope and Limitations

3 Methodology

# General Objective

- Develop and evaluate a machine-learning pipeline (MitoChime) to detect PCR-induced chimeric reads in *S. lemuru* 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

# Methodology

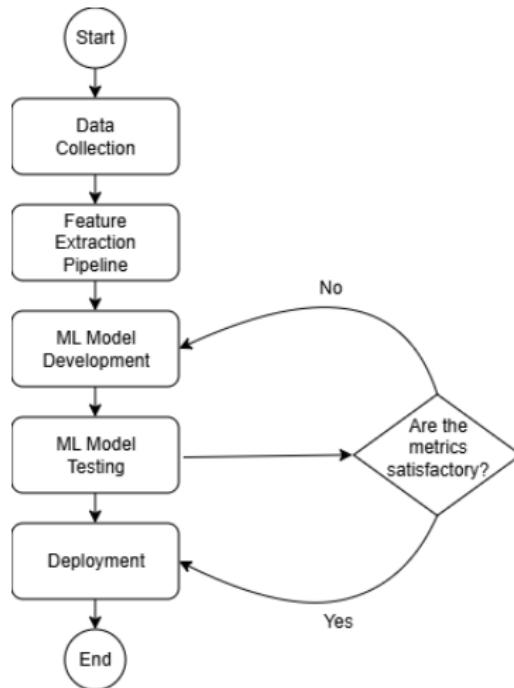


Figure: Process Diagram of the Special Project

# Data Collection

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

# Data Preprocessing

- 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.

# Data Preprocessing

- 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	
TTGTAGCTTAAACAAAGCATAACACTGAAAGATGGTCGTATAAGCCCCAACGCACTGAAAGTTGGCTCTGGCTTATTATCAGCTTACCGGAATTACACCAGCGAGGCCCTCCGCGGCCGTGAGGATGCCCTCA									NM:i:8 ms:i:220
AS:i:220	nn:i:8 tp:A:P cm:i:8 si:i:164	s2:i:0	def:f:0.0533	r1:i:0					
NC_039553_1_4_430_13:0:0 11:0:0 _243d	163	NC_039553.1	4	60	150M	=	281	427	
GGTGTAGCTTAAACAAAGCATAACACTGCAAGATGATCCGCTGGCCGTGATAAGCCGCAAGCAGGAGTGAAAGTTGGTCAGGCTTATTATCAGCTTACCCAAATTACACATGCGAGCCTCCGCGGCCGTGAGGATGCCCTAG									NM:i:13 ms:i:170
AS:i:170	nn:i:0 tp:A:P cm:i:9 si:i:135	s2:i:0	def:f:0.0867	r1:i:0					
NC_039553_1_5_495_6:0:0 11:0:0 _1d49	163	NC_039553.1	5	60	150M	=	346	491	
GTGTAGCTTACACAAAGCATAACACTGAAAGATGTTAAGATGGCCGTGATCAGCCCCAACAGCACTGAAAGTTGGCTCTGGCTTATTATCAGTTTCCCCAATTACACATGCGAGCCTCCGCGGCCGTGAGGATGCCCTAGC									NM:i:6 ms:i:170
AS:i:240	nn:i:0 tp:A:P cm:i:12 si:i:148	s2:i:0	def:f:0.04	r1:i:0					
NC_039553_1_6_523_6:0:0 9:0:0 _82c	163	NC_039553.1	6	60	150M	=	374	518	
TGTAGCTTAAACAAAGCATAACACTGAAAGATGTTAAGATGGCCGTGATAAGCCCCAACAGCACTGAAAGTTGGCTCTGGCTATTACACGTTTACCCAAATTACACATGCGCCCTCCGCGGCCGTGAGGATGCCCTAGCC									NM:i:6 ms:i:240
AS:i:240	nn:i:0 tp:A:P cm:i:10 si:i:157	s2:i:0	def:f:0.04	r1:i:0					
NC_039553_1_9_574_7:0:0 7:0:0 _181b	163	NC_039553.1	9	60	150M	=	425	566	
AGCTTAAACAAAGCATAACACTGCAAGATGTTAAGCTGGCCGTGATAAGCCCCAACAGCACTGAAAGTTGGCTCTGGCTTATTATCAGCTTACCCAAATTACACATGCGAGCCTCCGCGGCCGTGAGGCTGCCCTCCGCTCC									NM:i:7 ms:i:230
AS:i:230	nn:i:0 tp:A:P cm:i:12 si:i:176	s2:i:0	def:f:0.0467	r1:i:0					
NC_039553_1_10_391_9:0:0 8:0:0 _256b	99	NC_039553.1	10	60	150M	=	242	382	
GCTTAAACAAAGCCAACACTGAAAGATGTTAAGATGGCCGTGATAAGCCCCAACAGCACAGAACAGAAAGTTGGCTCTGGCTTATTAAACAGCTTACCCAAATTAGACATGCGAGCCTCCGCGGCCGTGATGCTGCCCTAGCCCTCC									NM:i:9 ms:i:210
AS:i:210	nn:i:0 tp:A:P cm:i:5 si:i:156	s2:i:0	def:f:0.06	r1:i:0					
NC_039553_1_11_509_6:0:0 11:0:0 _a19	99	NC_039553.1	11	60	150M	=	360	499	
CTTCACAAAGCATAACACTGAAAGATGTTAAGATGGCCGTATAAGCCCACAGCACTGAAAGTTGGCTCTGGCTTATTATGAGCTTACCCAAATTACACATGCGATCTCCGCGGCCGTGAGGATGCCCTAGCCTCCG									NM:i:6 ms:i:242
AS:i:242	nn:i:0 tp:A:P cm:i:10 si:i:150	s2:i:0	def:f:0.04	r1:i:0					
NC_039553_1_12_427_9:0:0 9:0:0 _157	163	NC_039553.1	12	60	150M	=	278	416	
TTAACAAAGCATAACACTGAAAGATTCAGATGGCCGTAAAGCCCCAACAGCACTGAAAGTTGGCTCTGGCTTATTATCAGCTTACCCAAATTACACATGCGAGCCTCGGGGGGCCGTGAGGATGCCCTCCGCTCCGT									NM:i:9 ms:i:210
AS:i:210	nn:i:0 tp:A:P cm:i:8 si:i:158	s2:i:0	def:f:0.06	r1:i:0					

Figure: SAM File of Clean Reads

# Data Preprocessing

```
AS:i:240      nn:i:0  tp:A:P  cm:1:19  s1:i:109      s2:i:0  def:f:0  SA:Z:NC_039553.1,2062,+,34M116S,1,0;  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_40985_41514_0:0:0  161  NC_039553.1  89  60  415109M = 7383  7444
CTTAATTATAGGAGGTCCCGCCCTGGCTTGACCAAAAAGTTTATTATCAGCTTACCCAAATTTCACACATGGAGGCCCTCGGGCCCCCGTGAGGATGCCCTCAGGCTCCCGTCCGGAGATGAGGAGCCGGCATCAGGCACGATGTCG
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:218
AS:i:218      nn:i:0  tp:A:P  cm:1:16  s1:i:194  s2:i:0  def:f:0  SA:Z:NC_039553.1,2051,+,45M105S,15,0;  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_105028_105471_0:0:0  81  NC_039553.1  89  60  96M54S = 13068  12885
TTTATTATCAGCTTACCCCAAATTTCACACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTAGCCCTCCGTCGGAGATGAGGAGCCGGCATCACCACTTGACAAGGCCACCGGCCTGTACAATTGCGTTACAGCTCTAGCACTCA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:192
AS:i:192      nn:i:0  tp:A:P  cm:1:15  s1:i:87  s2:i:0  def:f:0  SA:Z:NC_039553.1,4313,-,92558M,48,0;  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_40665_41142_0:0:0  2371  81  NC_039553.1  89  60  335117M = 3362  3158
-TATAGGAGGTCCGCCCTGCCCTTGACCAAAAAGTTTATTATCAGCTTACCCAAATTTCACACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTCGAGGCTCCGTCAGGAGATGAGGAGCCGGCATCAGGCACGATGTCGCCCATGA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:234
AS:i:234      nn:i:0  tp:A:P  cm:1:19  s1:i:109      s2:i:0  def:f:0  SA:Z:NC_039553.1,2059,-,37M113S,1,0;  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_41927_41581_0:0:0  ae9  97  NC_039553.1  90  60  150M = 7450  7467
TTTATATCAGCTTACCCCAAATTTCACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTCGGGAGATGAGGAGCCGGCATCAGGCACGATGTTCCGGCCATGAGCCCTGTAGCCACACCCCCAAGGGAAATTCA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:300
AS:i:1300     nn:i:0  tp:A:P  cm:1:25  s1:i:139      s2:i:0  def:f:0  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_5784_6251_0:0:0  145  NC_039553.1  90  60  150M = 6133  5895
TTTATATCAGCTTACCCCAAATTTCACACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTAGGGAGATGAGGAGCCGGCATCAGGCACGATGTTCCGGCCATGAGCCCTGTAGCCACACCCCCAAGGGAAATTCA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:300
AS:i:1300     nn:i:0  tp:A:P  cm:1:25  s1:i:139      s2:i:0  def:f:0  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_5788_6251_0:0:0  1913  81  NC_039553.1  90  60  150M = 6133  5895
TTTATATCAGCTTACCCCAAATTTCACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTAGGGAGATGAGGAGCCGGCATCAGGCACGATGTTCCGGCCATGAGCCCTGTAGCCACACCCCCAAGGGAAATTCA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:0  ms:i:300
AS:i:300      nn:i:0  tp:A:P  cm:1:25  s1:i:139      s2:i:0  def:f:0  rl:i:0
chimera_1_A9831-10051_B14983-15061_MH0_32277_32777_0:0:0  be3  161  NC_039553.1  91  60  150M = 6793  6812
NTTATCAGCTTACCCCAAATTTCACACATGGAGGCCCTCGGGCCCCGGTAGGGATGCCCTAGGGAGATGAGGAGCCGGCATCAGGCACGATGTTCCGGCCATGAGCCCTGTAGCCACACCCCCAAGGGAAATTCA
||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||||| NM:i:2  ms:i:296
AS:i:296      nn:i:2  tp:A:P  cm:1:25  s1:i:139      s2:i:0  def:f:0:0.0133  rl:i:0
```

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_basenr	ref_start	3'strand	mappq	cigar	has_sa	sa_count	num_segs_in_sa	diff	co_sa_min	dk_sa_max	d_sa_mean	sa_mean	sa_same	sa_opp	st_sa_max	m_sa_mean	sa_min	nr_sa_mean	softclip_le	softclip_ie	n_total	clippbreakpoint	kmer	cosi	kmer_js_d	microhom	microhomole
NC_09950	0	150	13 NC_09950	3	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.9720	0.97143	1	0		
NC_09950	0	150	13 NC_09950	4	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98591	0.98571	1	0		
NC_09950	0	150	13 NC_09950	5	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95887	0.95714	0	0		
NC_09950	0	150	13 NC_09950	6	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97183	0.97143	1	1		
NC_09950	0	150	13 NC_09950	7	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98644	0.98571	1	0		
NC_09950	0	150	13 NC_09950	10	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97296	0.97143	0	0		
NC_09950	0	150	13 NC_09950	11	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	1	1		
NC_09950	0	150	13 NC_09950	12	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98648	0.98571	1	1		
NC_09950	0	150	13 NC_09950	12	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.95869	0.95714	1	1		
NC_09950	0	150	13 NC_09950	14	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	0	0		
NC_09950	0	150	13 NC_09950	15	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98649	0.98571	1	0		
NC_09950	0	150	13 NC_09950	17	0	60	146M4S	0	0	1	0	0	0	0	0	0	0	0	0	0	4	4	146	0	0.5	0	0			
NC_09950	0	150	13 NC_09950	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98649	0.98571	3	0		
NC_09950	0	150	13 NC_09950	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97221	0.97143	3	0		
NC_09950	0	150	13 NC_09950	18	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98629	0.98571	3	0		
NC_09950	0	150	13 NC_09950	19	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97221	0.97143	3	0		
NC_09950	0	150	13 NC_09950	20	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97221	0.97143	0	0		
NC_09950	0	150	13 NC_09950	21	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0		
NC_09950	0	150	13 NC_09950	23	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98667	0.98571	0	0		
NC_09950	0	150	13 NC_09950	25	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98629	0.98571	0	0		
NC_09950	0	150	13 NC_09950	26	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0		
NC_09950	0	150	13 NC_09950	32	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97258	0.97143	2	1		
NC_09950	0	150	13 NC_09950	34	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	1	1	0	0		
NC_09950	0	150	13 NC_09950	34	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	0	0		
NC_09950	0	150	13 NC_09950	35	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0		
NC_09950	0	150	13 NC_09950	36	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98648	0.98571	0	0		
NC_09950	0	150	13 NC_09950	38	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98611	0.98571	1	0		
NC_09950	0	150	13 NC_09950	39	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.98684	0.98571	0	0		
NC_09950	0	150	13 NC_09950	41	0	60	150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	75	0.97296	0.97143	2	0.5		
NC_09950	0	150	13 NC_09950	43	0	60	150M	0	0	1	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

1	read_id	label	-3	read_id	mean	ref_nuc	ref_stal	strand	w	mpg	cigar	has_sai	sai_csq	num_sai	sai_dif	w	sai_min	w	sai_max	w	sai_pct	w	sai_max	w	sai_min	w	sai_max	w	sai_pct	w	sai_pct	w	total_c	w	breakup	w	kmer_c	w	kmer_j	w	microf	w	microf	w	biology_gc
19985	chimeric_3	1	150	40 NC_09565	40	1	60 150M	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0					
19986	chimeric_3	1	150	40 NC_09565	53	0	0 00000H	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0							
19987	chimeric_3	1	150	40 NC_09565	65	0	0 00000H	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0								
19988	chimeric_3	1	150	40 NC_09565	66	0	0 00000H	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0								
19989	chimeric_3	1	150	40 NC_09565	67	0	0 00000H	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
19990	chimeric_3	1	150	40 NC_09565	67	1	60 110M32S	1	1	2	0	0	4246	4246	0	0	1	10	10	0	0	0	0	0	0	32	32	32	118	1	1	1	0	0	0	0	0	0							
19991	chimeric_3	1	150	40 NC_09565	69	1	60 150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0								
19992	chimeric_3	1	150	40 NC_09565	76	0	0 60 100M41S	1	1	2	0	0	4237	4237	4237	0	0	1	16	16	0	0	0	0	0	0	41	41	41	109	1	1	1	0	0	0	0	0							
19993	chimeric_3	1	150	40 NC_09565	77	1	60 150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
19994	chimeric_3	1	150	40 NC_09565	79	0	0 60 100M44S	1	1	2	0	0	4234	4234	4234	0	0	1	17	17	0	0	0	0	0	0	44	44	44	106	1	1	1	0	0	0	0	0							
19995	chimeric_3	1	150	40 NC_09565	84	0	0 60 120M38S	1	1	2	0	0	5197	5197	0	0	1	10	10	0	0	0	0	0	0	38	38	38	101	0	0	0	0	0	0	0	0								
19996	chimeric_3	1	150	40 NC_09565	85	0	0 60 110M39S	1	1	2	0	0	5180	5180	5180	0	0	1	20	20	0	0	0	0	0	0	39	39	39	111	0	0	0	0	0	0	0	0							
19997	chimeric_3	1	150	40 NC_09565	88	0	0 60 150M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
19998	chimeric_3	1	150	40 NC_09565	89	0	0 60 150M39M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	15	0	15	15	1	1	1	0	0	0	0	0									
19999	chimeric_3	1	150	40 NC_09565	89	0	0 60 305120M	1	1	2	0	1973	1973	1973	0	0	1	1	1	0	0	0	0	0	0	30	0	30	30	0	0	0	0	0	0	0	0								
20000	chimeric_3	1	150	40 NC_09565	89	0	0 60 415109M	1	1	2	0	1962	1962	1962	0	0	1	15	15	0	0	0	0	0	0	41	0	41	41	0	0	0	0	0	0	0	0								
20001	chimeric_3	1	150	40 NC_09565	89	1	60 969M54S	1	1	2	0	4224	4224	4224	0	0	1	48	48	0	0	0	0	0	0	54	96	1	1	1	0	0	0	0	0	0	0								
20002	chimeric_3	1	150	40 NC_09565	89	1	60 335117M	1	1	2	0	1970	1970	1970	0	0	1	1	1	0	0	0	0	0	0	33	0	33	33	0	0	0	0	0	0	0	0								
20003	chimeric_3	1	150	40 NC_09565	90	0	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	0	0	0	0	0	0	0	0	0						
20004	chimeric_3	1	150	40 NC_09565	90	1	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	0	0	0	0	0	0	0	0	0						
20005	chimeric_3	1	150	40 NC_09565	90	1	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	0	0	0	0	0	0	0	0	0						
20006	chimeric_3	1	150	40 NC_09565	91	0	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	0	0	0	0	0	0	0	0	0						
20007	chimeric_3	1	150	40 NC_09565	91	0	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	0	0	0	0	0	0	0	0	0						
20008	chimeric_3	1	150	40 NC_09565	91	1	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	0	0	0	0	0	0	0	0	0						
20009	chimeric_3	1	150	40 NC_09565	91	1	60 949M65S	1	1	2	0	4222	4222	4222	0	0	1	52	52	0	0	0	0	0	0	56	56	94	1	1	1	0	0	0	0	0	0	0							
20010	chimeric_3	1	150	40 NC_09565	92	0	0 295121M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	29	0	29	0	0	0	0	0	0	0	0									
20011	chimeric_3	1	150	40 NC_09565	92	0	0 275278M	1	1	2	0	3064	3064	3064	0	0	1	60	60	0	0	0	0	0	0	72	0	72	72	0	0	0	0	0	0	0	0								
20012	chimeric_3	1	150	40 NC_09565	92	0	0 665848M	1	1	2	0	3070	3070	3070	0	0	1	59	59	0	0	0	0	0	0	66	0	66	66	0	0	0	0	0	0	0	0								
20013	chimeric_3	1	150	40 NC_09565	92	0	0 665849M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	11	0	11	11	0	0	0	0	0	0	0									
20014	chimeric_3	1	150	40 NC_09565	92	0	0 665120M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
20015	chimeric_3	1	150	40 NC_09565	92	0	0 651547M	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0									
20016	chimeric_3	1	150	40 NC_09565	92	0	0 54599M	1	1	2	0	3082	3082	3082	0	0	1	30	30	0	0	0	0	0	0	54	0	54	54	0	0	0	0	0	0	0	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

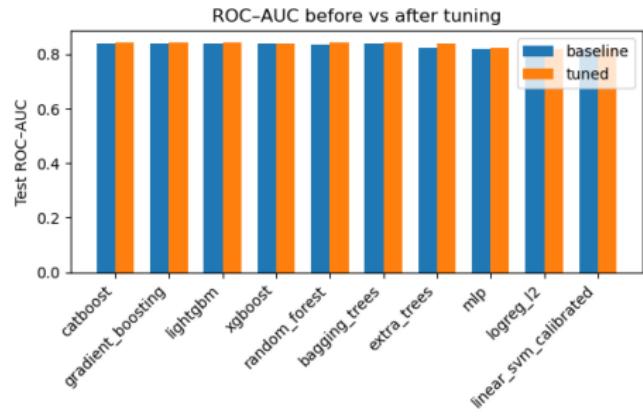
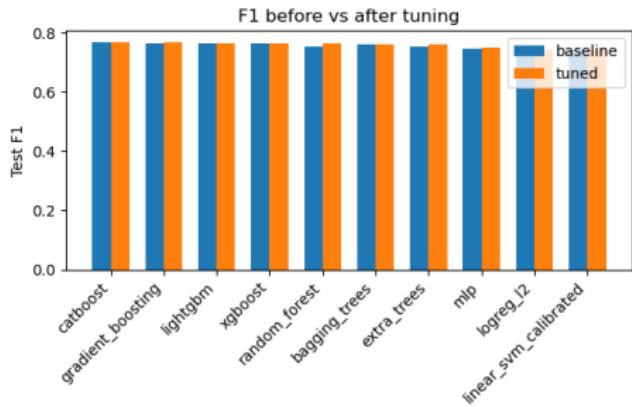


Figure: Test F1: baseline vs tuned.

Figure: Test ROC-AUC: baseline vs tuned.

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

# ROC and precision–recall curves

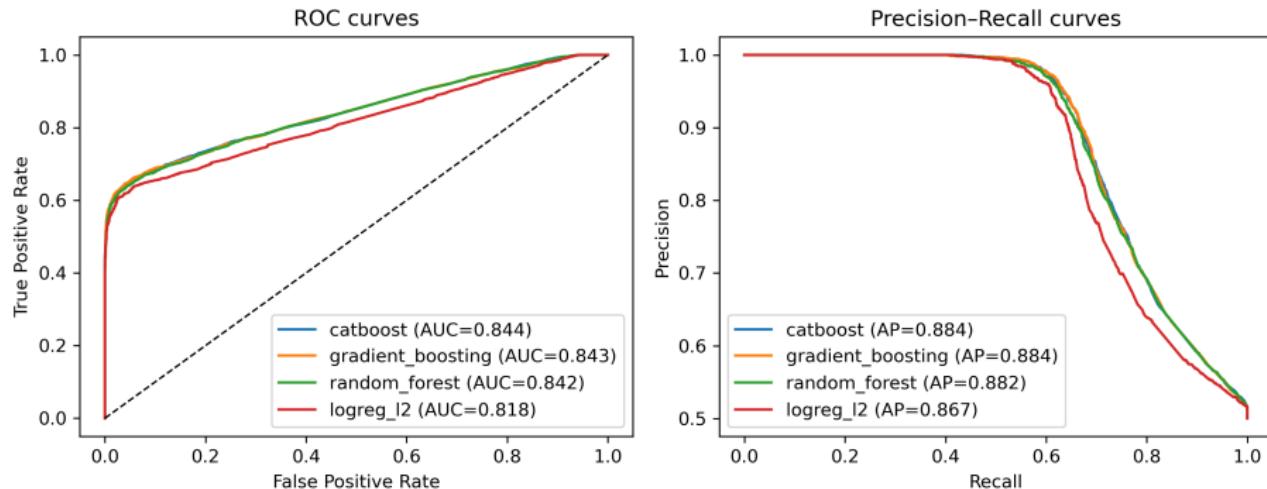
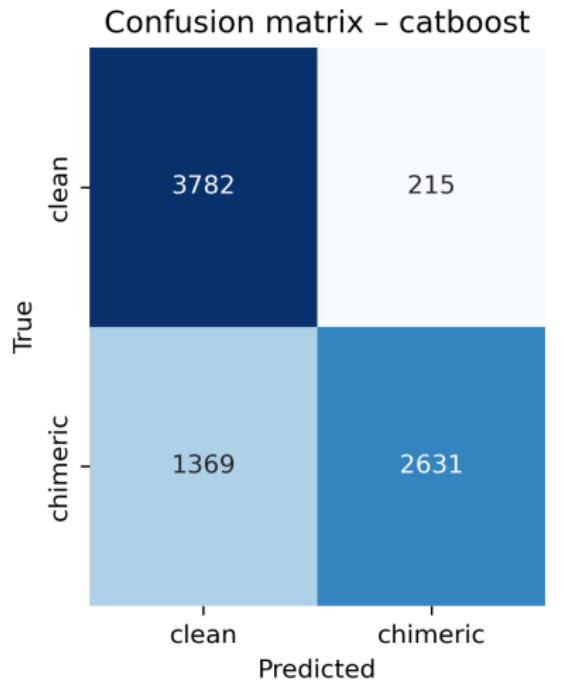


Figure: ROC (left) and PR (right) curves for CatBoost, Gradient Boosting, Random Forest, and logistic regression.

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

# 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

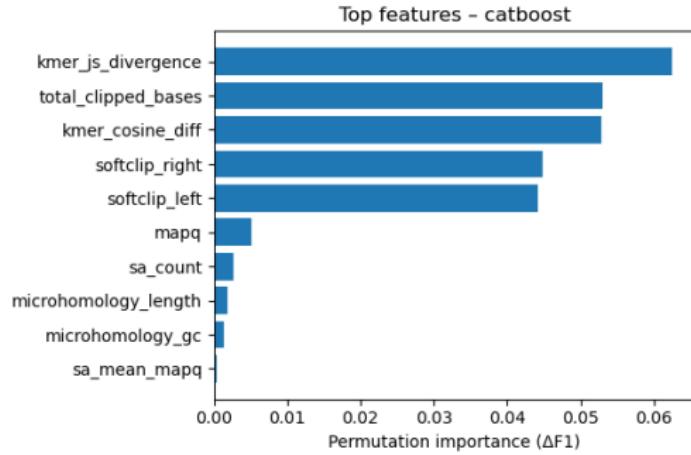


Figure: Permutation importance ( $F_1$ ) for CatBoost.

- 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

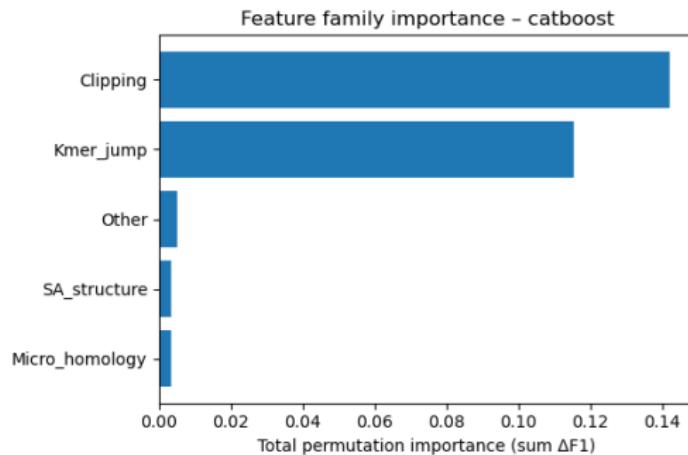


Figure: Aggregated feature families for CatBoost.

- 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

# ML component: summary and implications

- Per-read classifier for clean vs PCR chimeric reads
- Evaluation:
  - Stratified 80/20 split, 5-fold CV, held-out test set
- Best models: tree-based ensembles
  - CatBoost, Gradient Boosting, LightGBM
  - Test F1  $\approx 0.77$ , ROC-AUC  $\approx 0.84$
- Default threshold:
  - Conservative chimera filter (high clean recall, high chimera precision)
  - Removes  $\sim 2/3$  of chimeras
- Features match PCR chimera biology
- Practical, interpretable pre-filter before mitochondrial assembly