

# 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
- 4 Train–Test Split and Validation
- 5 Model Zoo and Training
- 6 Metrics and Interpretation
- 7 Results
- 8 Discussion and Conclusion

# 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, junction complexity, 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 (Nanopore, PacBio) 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



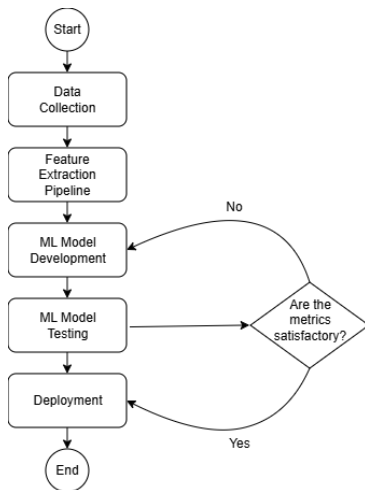


Figure: Process Diagram of the Special Project

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
TGGTGTAGCTTAAACAAGCATAAAGCTGAAGATGTTACGATGGGCGGTGAAGAGCCACAGCAGCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCGGAATTTACACACCGAGAGCTCCGCGGCGCGGTGAGGATGGCTCA
..... 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
GGTGTAGCTTAAACAAGCATAAAGCTGAGATGATCCGCTGGGCGGTGAAGAGCCAGCAGGAGTGAAAGTTTGGTCCAGGCTTTATTATCAGCTTTACCCCAATTTACACATGCGAGCTCCGCGGCGCGGTGAGGATGGCTCAG
..... 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
GTGTAGCTTACACAAGCATAAAGCTGAAGATGTTAAGATGGGCGGTGATCAGCCCAAGCAGCTGAAAGGTTAGGTCCTGGCTTTATTATCAGGTTTCCCCCAATTTACACATGCGAGCTCCGCGGCGCGGTGAGGATGGCTCAGC
..... 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
TGTAGCTTAAACAAGCATAAAGCTGAAGATGTTAAGATGGGCGGTGATCAGCCCAAGCAGCTGCAAGGTTTGGTCTGGCTATATTACAGCTTTACCCCAATTTACACATGCGAGCTCCGCGGCGCGGTGAGGATGGCTCAGCC
..... 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
AGCTTAAACAAGCATAAAGCTGAGATGTTAAGCTGGGCGGTGATAAGCCCAAGCAGCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCGCAATTTACACATGCGAGCTCCGCGGCGCGGTGAGGCTGCCCTCCGCTCC
..... 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
GCTTAAACAAGCATAAAGCTGAGATGTTAAGATGGGCGGTGATAAGCCCAAGCAGCTGAAAGGTTTGGTCTGGCTTTATTATCAGCTTTACCCCAATTTAGACATGCGAGCTCCGCGGCGCGGTGATGCTGCCCTCAGCTCCC
..... 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
CTTCAACAAGCATAAAGCTGAAGATGTTAAGATGGGCGGTGATAAGCCCAAGCAGCTGAAAGGTTAGGTCCTGGCTTTATTATGAGCTTTACCCCAATTTACACATGCGATCTCCGCGGCGCGGTGAGGATGCCCTCAGCTCCCG
..... 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
TTAAACAAGCATAAAGCTGAAGATTTAGATGGGCGGTGATAAGCCCAAGCAGCTGAAAGTTTGGTCTGGCTTTATTATCAGCTTTACCCCAATTTACACATGCGAGCTCCGCGGCGCGGTGAGGATGCCCTCCGCTCCCGT
..... 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

chimer1_A1981-10051	B14983-15061	M40	40985	41514	0:0:0	0:0:0	0:0:0	1047	161	NC	039553.1	89	60	415109M =	7383	7444		
CTCAATATATAGGAGGTCGCCGCTGCCCTGTGACCAAAAGTTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCG																		
chimer1_A1981-10051	B14983-15061	M40	105028	105471	0:0:0	0:0:0	0:0:0	e2e	81	NC	039553.1	89	60	96W5AS =	13068	12885	NM:i:0	ms:i:218
TTTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCACCACCTTGACAGGCCCAACGCCCTGACAAATTCGCGTTACAGCTTAGCACTCA																		
chimer1_A1981-10051	B14983-15061	M40	40665	41142	0:0:0	0:0:0	0:0:0	2371	81	NC	039553.1	89	60	335117M =	3362	3158	NM:i:0	ms:i:192
TATAGGAGGTCGCCGCTGCCCTGTGACCAAAAGTTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCGCCGCCATGA																		
chimer1_A1981-10051	B14983-15061	M40	41027	41581	0:0:0	0:0:0	0:0:0	aer	97	NC	039553.1	90	60	150M =	7450	7467	NM:i:0	ms:i:234
TTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCGCCGCCATGACGCCCTGTTAGGCACACCCCCAAGGGAATTCAG																		
chimer1_A1981-10051	B14983-15061	M40	5784	6251	0:0:0	0:0:0	0:0:0	1330	145	NC	039553.1	90	60	150M =	6133	5895	NM:i:0	ms:i:300
TTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCGCCGCCATGACGCCCTGTTAGGCACACCCCCAAGGGAATTCAG																		
chimer1_A1981-10051	B14983-15061	M40	5788	6251	0:0:0	0:0:0	0:0:0	1913	81	NC	039553.1	89	60	150M =	6133	5895	NM:i:0	ms:i:300
TTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCGCCGCCATGACGCCCTGTTAGGCACACCCCCAAGGGAATTCAG																		
chimer1_A1981-10051	B14983-15061	M40	32227	32777	0:0:0	0:0:0	0:0:0	be3	161	NC	039553.1	91	60	150M =	6793	6812	NM:i:0	ms:i:300
TTTATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCGCCGCCATGACGCCCTGTTAGGCACACCCCCAAGGGAATTCAG																		
chimer1_A1981-10051	B14983-15061	M40	40985	41514	0:0:0	0:0:0	0:0:0	1047	161	NC	039553.1	89	60	415109M =	7383	7444	NM:i:2	ms:i:296
CTCAATATATAGGAGGTCGCCGCTGCCCTGTGACCAAAAGTTATATACAGCTTACCCCAATTTACACATGCGAGCTCCGCGGGCCCGTGAGGATGCCTTCAGCTCCCGTCGGAGATGAGGAGCGGGGATCAGGCACAGATGTCG																		

### Figure: SAM File of Chimeric Reads

# Stratified Train–Test Split

- First step: **create a held-out test set** for final evaluation.
- Use `build_datasets.py`:
  - 1 Combine clean and chimeric feature tables.
  - 2 Attach labels (0 = clean, 1 = chimeric) if missing.
  - 3 Shuffle and perform **stratified** split:

Train : Test = 80% : 20%

with the same class proportions in each split.

- Output:
  - `train.tsv` (used for model selection and cross-validation).
  - `test.tsv` (kept untouched until the very end).

# 5-Fold Stratified Cross-Validation

- On the **training set only**, we perform:

5-fold stratified cross-validation

- Procedure:

- ① Split training data into 5 folds with balanced 0/1 labels.
- ② For each fold:
  - Train the model on 4 folds.
  - Evaluate on the remaining fold.
- ③ Average metrics across the 5 folds:

mean F1  $\pm$  std,    mean accuracy  $\pm$  std

- This tells us:

- **Typical performance** on unseen data.
- **Stability** of each model (via standard deviation).
- Helps guide which algorithms are promising before going to the test set.



# Model Zoo: Algorithms Compared

- We implemented a panel of 13 classifiers using scikit-learn and gradient boosting libraries:
  - **Baseline:** Dummy (always predicts most frequent class).
  - **Linear models:** Logistic regression (logreg\_12), linear SVM with calibration.
  - **Tree ensembles:**
    - Random Forest, Extra Trees.
    - Gradient Boosting (sklearn).
    - XGBoost, LightGBM, CatBoost.
    - Bagging with decision trees.
  - **Others:** k-NN, Gaussian Naive Bayes, shallow MLP.
- All models use the same preprocessing pipeline:

Imputer (median) → StandardScaler → Classifier

# Hyperparameter Tuning for Top Models

- For the 10 strongest families, we perform **RandomizedSearchCV** with 5-fold CV:
  - Logistic regression, linear SVM (calibrated).
  - Random Forest, Extra Trees, Gradient Boosting.
  - XGBoost, LightGBM, CatBoost.
  - Bagging (trees), MLP.
- Each search explores combinations of:
  - Tree depth, number of estimators, learning rate, subsample ratios, etc.
  - For MLP: hidden layer sizes, regularization ( $\alpha$ ), learning rate.
- Selection criterion:
  - Choose the hyperparameters with the best **cross-validated F1-score**.
  - Re-fit the best model on the **full training set**, then evaluate on the held-out test set.

# Classification Metrics (Per-Read)

- For each model, on the test set we compute:

- Accuracy:**

$$\frac{\# \text{ correct predictions}}{\# \text{ all predictions}}$$

- Precision** (for chimeras):

$$\frac{TP}{TP + FP}$$

Of the reads we call “chimeric”, how many are truly chimeric?

- Recall** (for chimeras):

$$\frac{TP}{TP + FN}$$

Of all true chimeric reads, how many did we detect?

- F1-score** (for chimeras):

$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Harmonic mean: high only if both precision and recall are high.

# Threshold-Free Metrics: ROC–AUC and PR Curves

- Our models output a **score** per read (probability of being chimeric).
- By sweeping a threshold on this score, we can draw:
  - **ROC curve**:
    - x-axis: False Positive Rate (FPR).
    - y-axis: True Positive Rate (TPR = recall).
    - **ROC–AUC** = area under the curve.
  - **Precision–Recall (PR) curve**:
    - x-axis: Recall.
    - y-axis: Precision.
    - **Average Precision (AP)** = area under PR curve.
- Intuition for ROC–AUC:

$\text{AUC} \approx 0.84 \Rightarrow 84\%$  chance a random chimera is scored higher than a random non-chimera

# Overall Performance Across Models (Test Set)

Model	CV Acc	CV F1	Test Acc	Test F1	ROC-AUC
Dummy baseline	0.50	0.67	0.50	0.67	0.50
Logistic regression	0.79	0.75	0.79	0.74	0.82
Linear SVM (cal.)	0.79	0.75	0.79	0.74	0.82
Random Forest	0.80	0.77	0.79	0.75	0.83
Extra Trees	0.80	0.77	0.79	0.75	0.82
Gradient Boosting	0.81	0.78	0.80	0.77	0.84
XGBoost	0.81	0.77	0.80	0.76	0.84
LightGBM	0.81	0.77	0.80	0.76	0.84
CatBoost	0.81	0.78	0.80	0.77	0.84
k-NN	0.78	0.75	0.78	0.75	0.81
Gaussian NB	0.75	0.66	0.74	0.65	0.82
Bagging (trees)	0.80	0.77	0.79	0.76	0.84
MLP	0.79	0.75	0.79	0.75	0.82

**Table:** Summary of cross-validation and test performance (chimeric class F1).

# ROC and PR Curves (Placeholder)

ROC curves (CatBoost, GBM, RF,  
logreg)

Precision–Recall curves

- ROC–AUC for top models:  $\approx 0.84$ .
- Curves pushed towards the top-left / top-right illustrate strong separation between clean and chimeric reads across thresholds.

# Confusion Matrix and Class-Wise Behaviour (CatBoost)

Placeholder: confusion matrix for CatBoost on test set

## CatBoost (test set, illustrative)

clean: precision  $\approx 0.73$ , recall  $\approx 0.95$   
chimeric: precision  $\approx 0.92$ , recall  $\approx 0.66$   
overall accuracy  $\approx 0.80$

- **Clean reads:**
  - Very high recall: most true clean reads are correctly kept.
- **Chimeric reads:**
  - High precision: when we call a read chimeric, it is usually correct.
  - Moderate recall: we detect about two-thirds of all chimeras.
- Practical behaviour: **conservative chimera filter** that prioritizes not discarding clean reads.

# Effect of Hyperparameter Tuning (F1 and ROC–AUC)

Model	F1 (base)	AUC (base)	F1 (tuned)	AUC (tuned)
CatBoost	0.767	0.839	0.769	0.844
Gradient Boosting	0.766	0.840	0.767	0.843
LightGBM	0.764	0.838	0.766	0.842
XGBoost	0.765	0.839	0.765	0.839
Random Forest	0.755	0.834	0.763	0.842
Bagging (trees)	0.760	0.837	0.763	0.842
Extra Trees	0.753	0.824	0.760	0.837
MLP	0.748	0.819	0.749	0.821
Logistic reg.	0.744	0.821	0.743	0.818
Linear SVM (cal.)	0.744	0.820	0.743	0.818

**Table:** Test F1 and ROC–AUC before vs after hyperparameter tuning.

- Tuning yields **modest but consistent gains** in F1 and ROC–AUC.
- Confirms that the initial defaults were already reasonable, but performance can be further refined.



# Permutation Feature Importance (Placeholder)

Placeholder: permutation importance for CatBoost

- Top features across CatBoost, GBM, RF:
  - total\_clipped\_bases
  - kmer\_js\_divergence, kmer\_cosine\_diff
  - softclip\_left, softclip\_right
  - mapq
- Interpretation:
  - Chimeras are characterized by **large clipped segments** and **abrupt k-mer composition shifts**.
  - Aligners are already “seeing” the breakpoint signal; the ML model learns to combine these signals into a chimera score.

# Summary of Findings

- We built a per-read feature table capturing:
  - Alignment and clipping patterns.
  - Supplementary alignments and breakpoint distances.
  - Sequence-level k-mer divergence and microhomology.
- A broad panel of ML models was evaluated:
  - Tree-based ensembles (CatBoost, Gradient Boosting, Random Forest, LightGBM, XGBoost) achieved the **best performance**.
  - Test F1 for chimeras  $\approx 0.76$ – $0.77$ , ROC–AUC  $\approx 0.84$ .
- Model behaviour:
  - Conservative on clean reads (high recall).
  - High precision on chimeric reads, moderate recall.

# Implications for Mitochondrial Assembly

- The ML classifier can be used as a **pre-filter** before assembling mitochondrial genomes:
  - Remove high-confidence chimeric reads to reduce false junctions.
  - Retain the majority of clean reads to preserve coverage.
- Especially useful for:
  - Small, circular, and repetitive organellar genomes where chimeras are particularly harmful.
  - Scenarios without high-quality reference genomes or abundance information.
- The feature importance analysis provides biological insight:
  - Confirms the role of soft-clipping, supplementary alignments, and k-mer jumps as core signals of chimeric structure.

# Limitations and Future Work

- Current study uses **simulated** chimeras and a single species:
  - Need to validate on real experimental datasets.
  - Extend to other organellar genomes and library preparations.
- Classifier currently treats each read independently:
  - Future work: incorporate read-pair information, local read depth, or graph features.
- Integration into practical pipelines:
  - Wrap as a command-line tool interfacing with standard BAM/FASTQ workflows.
  - Benchmark impact on final assembly quality (contiguity, misassemblies).

# Conclusion

- We developed a **machine learning pipeline** that:
  - Learns from alignment- and sequence-based features.
  - Achieves strong separation between clean and chimeric reads.
- Tree-based gradient boosting models (CatBoost, GBM, RF) provide:
  - High test F1 and ROC-AUC.
  - Interpretable feature importance aligned with known chimera mechanisms.
- This framework is a step towards **reference-free chimera detection** tailored for organellar genomes and low-resource settings.

# Thank You

Questions?