

Speaker Recognition using CNN-RNN Architecture

Introduction to Machine Learning Course - Final Report

Team: Aleksander, Mantas, Michał, Piotr, Rafał

January 2026

Abstract

This report documents the development and implementation of a speaker recognition system using a hybrid CNN-RNN architecture trained on voice recordings. The project focuses on classifying speakers based on short voice snippets through a sophisticated neural network combining convolutional layers for feature extraction, recurrent layers for temporal processing, and advanced training techniques including data augmentation, dropout regularization, and Monte Carlo uncertainty estimation. The system was trained on a dataset of 5 speakers with comprehensive evaluation metrics, TensorBoard logging, and inference capabilities supporting both deterministic and probabilistic predictions.

Contents

1	Introduction	4
1.1	Project Objectives	4
1.2	Dataset	4
2	Data Preprocessing and Feature Extraction	4
2.1	Audio Processing Pipeline	4
2.2	Log-Mel Spectrogram Features	5
2.3	Data Augmentation	5
2.4	Data Loading Optimization	5
3	Model Architecture	6
3.1	Overall Architecture	6
3.2	Backbone Network	6
3.2.1	CNN Block	6
3.2.2	Squeeze-and-Excitation Block	6
3.2.3	RNN Block	6
3.2.4	Attentive Statistics Pooling	7
3.2.5	Embedding Projection	7
3.3	Classification Head: AAMSoftmax	7
4	Training Configuration and Methodology	8
4.1	Hyperparameters	8
4.2	Regularization Techniques	8
4.3	Training Loop	8
4.4	Mixed Precision Training	9
4.5	TensorBoard Logging	9
5	Advanced Features and Uncertainty Quantification	9
5.1	Monte Carlo Dropout	9
5.2	Speaker Verification Pipeline	10
6	Evaluation and Results	10
6.1	Metrics	10
6.2	Checkpointing	10
6.3	Expected Performance	10
6.4	Training Curves and TensorBoard Visualization	11
6.4.1	Loss Dynamics	11
6.4.2	Accuracy Progression	12
6.4.3	Learning Rate Schedule	12
7	Training Monitoring and Metrics	12
7.1	TensorBoard Integration	12
7.2	Per-Class Performance Analysis	13
8	Implementation Highlights and Helper Utilities	13
8.1	Data Loading Utilities	13
8.2	Model Architecture Code	14
8.3	Training Utilities	14
8.4	Checkpoint Management	14

9	Course Requirements Coverage	15
10	Conclusions and Future Work	15
10.1	Key Achievements	15
10.2	Future Improvements	15

1 Introduction

Speaker recognition is a challenging task in machine learning that requires the model to distinguish between different speakers based on acoustic features of their voice. Unlike speech recognition (which identifies what is being said), speaker recognition identifies *who* is speaking. This task has numerous practical applications including speaker verification, voice authentication, and voice-controlled systems.

1.1 Project Objectives

The primary objectives of this project are:

- Develop a robust speaker classification system capable of recognizing multiple speakers
- Implement advanced deep learning techniques including CNN and RNN architectures
- Apply data augmentation and quality enhancement techniques to improve model generalization
- Incorporate uncertainty quantification through Monte Carlo dropout
- Create a comprehensive evaluation pipeline with detailed performance metrics
- Meet course requirements: implement and document at least 8 advanced ML concepts

1.2 Dataset

The dataset consists of voice recordings from 5 speakers (Aleksander, Mantas, Michał, Piotr, Rafał). Voice data is stored in various formats (WAV, MP3, M4A) and is preprocessed into log-mel spectrograms. The data is split into training (70%), validation (15%), and test (15%) sets, stored in HDF5 format for efficient access during training.

2 Data Preprocessing and Feature Extraction

2.1 Audio Processing Pipeline

The data preprocessing pipeline involves several critical steps:

1. **Audio Loading:** Support for multiple formats (WAV, MP3, M4A, WMA) using librosa and pydub
2. **Normalization:** Volume normalization to ensure consistent loudness across samples
3. **Silent Passage Removal:** Automated detection and removal of silence from audio files to focus on speech content
4. **Spectrogram Computation:** Conversion of time-domain audio to frequency-domain log-mel spectrograms
5. **Dataset Storage:** HDF5 format with metadata including speaker mapping, sample rate, and preprocessing configuration

2.2 Log-Mel Spectrogram Features

Log-mel spectrograms are computed using the following process:

- **Sample Rate:** 16,000 Hz
- **Number of Mel Bands:** 64 frequency bins
- **Window Function:** Hann window with 25ms window length
- **Hop Length:** 10ms (50% overlap)
- **Frequency Range:** 50-8,000 Hz (optimized for speech)
- **Log Scaling:** Applied with floor value of -80 dB for dynamic range compression

The resulting spectrogram has shape (T, F) where T is the number of time frames and $F = 64$ is the number of frequency bins. Each sample represents approximately 10 milliseconds of audio.

2.3 Data Augmentation

To improve model robustness and generalization, the following augmentation techniques are applied:

- **SpecAugment:** Time and frequency masking of the spectrogram
- **Speed Perturbation:** Varying playback speed to simulate speaking rate variations
- **VTLP (Vocal Tract Length Perturbation):** Simulating different speaker vocal tract lengths
- **Noise Addition:** Background noise injection for robustness
- **Volume Normalization:** Per-sample normalization to reduce volume-based bias

These augmentations help the model learn speaker-specific characteristics rather than audio artifacts.

2.4 Data Loading Optimization

Two data loading strategies were implemented:

1. **Normal Loader:** Standard HDF5 disk-based I/O using DataLoader with multiple workers
2. **Cached Loader:** RAM-cached variant that loads entire dataset into memory at startup (2-3 seconds), providing **10-50x speedup** during training

The cached loader is recommended for fast iteration, achieving batch iteration times of approximately 10-20ms with batch size 128 on modern GPUs.

3 Model Architecture

3.1 Overall Architecture

The speaker recognition model consists of four main components:

1. **Backbone:** Feature extraction via CNN and temporal modeling via RNN
2. **Embedding Head:** Projection to fixed-dimensional speaker embeddings
3. **AAMSoftmax:** Classification head with additive angular margin loss
4. **Inference Module:** Deterministic and Monte Carlo dropout-based prediction

3.2 Backbone Network

The backbone network processes log-mel spectrograms to extract speaker-discriminative features:

3.2.1 CNN Block

- **Layer 1:** Conv2d(1→32, kernel=3, padding=1) + BatchNorm + ReLU + SEBlock + MaxPool(1,2)
- **Layer 2:** Conv2d(32→64, kernel=3, padding=1) + BatchNorm + ReLU + SEBlock + MaxPool(1,2)
- **Layer 3:** Conv2d(64→128, kernel=3, padding=1) + BatchNorm + ReLU + SEBlock + MaxPool(1,2)
- **Output Shape:** $(B, 128, T_{pooled}, F_{pooled})$ where pooling reduces F by $8\times$

3.2.2 Squeeze-and-Excitation Block

SE blocks apply channel-wise attention to adaptively reweight feature maps:

$$\text{SE}(x) = x \odot \sigma(\text{FC}_2(\text{ReLU}(\text{FC}_1(\text{GAP}(x))))$$

where GAP is global average pooling, FC are fully connected layers, and σ is sigmoid.

3.2.3 RNN Block

- **Type:** Bidirectional GRU (Gated Recurrent Unit)
- **Hidden Dimension:** 256 units
- **Layers:** 2 stacked layers
- **Dropout:** 0.2 between layers
- **Input Dimension:** $128 \times (64/8) = 1024$ (flattened CNN output per time frame)
- **Output:** $(B, T, 512)$ for bidirectional output

The bidirectional RNN captures both past and future context for each time frame, crucial for speaker modeling.

3.2.4 Attentive Statistics Pooling

Instead of simple mean pooling, attention-based statistics pooling is employed:

$$\begin{aligned}\alpha_t &= \tanh(W_{att} \cdot h_t + b_{att}) \\ w_t &= \text{softmax}(\alpha_t) \quad (\text{with masking for padding}) \\ \mu &= \sum_t w_t \odot h_t \quad (\text{weighted mean}) \\ \sigma &= \sqrt{\sum_t w_t \odot (h_t - \mu)^2} \quad (\text{weighted variance}) \\ \text{stats} &= [\mu; \sigma]\end{aligned}$$

This allows the model to focus on informative time frames (e.g., voiced speech) rather than silence or noise.

3.2.5 Embedding Projection

The statistics vector is projected to a fixed-dimensional embedding space:

$$z = \frac{FC_2(\text{ReLU}(\text{BN}(FC_1(\text{stats}))))}{\|FC_2(\dots)\|}$$

where FC_1 projects from $2 \times 512 = 1024$ to 256, and FC_2 projects to the final embedding dimension of 256. The embedding is L2-normalized for cosine distance similarity.

3.3 Classification Head: AAMSoftmax

Additive Angular Margin (AAM) Softmax is used as the classification head to enhance intra-class compactness and inter-class separability:

$$L_{AAM} = -\frac{1}{N} \sum_i \log \frac{e^{s \cos(\theta_{y_i} + m)}}{e^{s \cos(\theta_{y_i} + m)} + \sum_{j \neq y_i} e^{s \cos(\theta_j)}}$$

where:

- θ is the angle between embedding and weight vectors (normalized)
- $s = 30$ is the scaling factor
- $m = 0.25$ is the angular margin (0.25 radians $\approx 14.3^\circ$)

This loss explicitly penalizes incorrect predictions more when they are in the same angular direction as correct predictions, improving class separation in embedding space.

4 Training Configuration and Methodology

4.1 Hyperparameters

Hyperparameter	Value	Rationale
Optimizer	AdamW	Adaptive learning with weight decay regularization
Learning Rate	1×10^{-3}	Standard for embedding learning
Weight Decay	1×10^{-4}	L2 regularization for generalization
Batch Size (Train)	128	Optimized for 4GB+ GPU memory
Batch Size (Val/Test)	128	Large batch for stable validation
Epochs	30-50	Early stopping patience: 5 epochs
LR Scheduler	CosineAnnealingLR	Smooth decay over training ($T_{max} = 50$)
Loss Function	Cross-Entropy	Standard for classification with AAMSoftmax

4.2 Regularization Techniques

The model employs multiple regularization strategies:

1. **Dropout:** 0.2 in RNN layers, 0.3 MC dropout after CNN and in statistics
2. **Batch Normalization:** Before activation functions (at convolutions and projection)
3. **Weight Decay:** $\lambda = 10^{-4}$ in AdamW optimizer
4. **Data Augmentation:** SpecAugment, speed perturbation, VTLP, noise injection

4.3 Training Loop

The training pipeline consists of:

Listing 1: Training Loop Structure

```

1  for epoch in range(1, epochs+1):
2      # Training phase
3      for batch in train_loader:
4          X, y, lengths = batch
5          optimizer.zero_grad()
6          with autocast(): # Mixed precision training
7              logits, embeddings = model(X, y, lengths)
8              loss = loss_fn(logits, y)
9              scaler.scale(loss).backward()
10             scaler.step(optimizer)
11             scaler.update()
12
13     # Validation phase
14     val_loss, val_acc = validate(model, val_loader)
15
16     # Checkpoint and early stopping
17     if val_acc > best_val_acc:
18         save_checkpoint(model, best_path)
19         patience = 0
20     else:
21         patience += 1

```



```

22         if patience >= 5: break
23
24     scheduler.step() # Update learning rate

```

4.4 Mixed Precision Training

Automatic Mixed Precision (AMP) is enabled using PyTorch’s `autocast()` context manager. This:

- Reduces GPU memory usage by **30-40%**
- Accelerates training by **20-30%**
- Maintains numerical stability through gradient scaling

4.5 TensorBoard Logging

Comprehensive metrics are logged per batch and epoch:

- Per-batch: loss, accuracy, learning rate
- Per-epoch: macro-averaged precision, recall, F1-score, per-class metrics
- GPU/Memory utilization monitoring

This enables real-time monitoring and debugging during training.

5 Advanced Features and Uncertainty Quantification

5.1 Monte Carlo Dropout

Monte Carlo dropout enables uncertainty quantification by performing multiple stochastic forward passes during inference:

$$\begin{aligned}
 \hat{z} &= \text{MC-Dropout}(x; n_{\text{samples}}) \\
 z_{MC}^{(k)} &= \text{dropout}(f(x), p = 0.3) \quad \forall k = 1, \dots, K \\
 \mu_{MC} &= \frac{1}{K} \sum_{k=1}^K z_{MC}^{(k)} \\
 \sigma_{MC}^2 &= \frac{1}{K} \sum_{k=1}^K (z_{MC}^{(k)} - \mu_{MC})^2
 \end{aligned}$$

The variance provides a measure of model confidence. Lower variance indicates higher confidence in the prediction.

5.2 Speaker Verification Pipeline

The model supports two inference modes:

1. **Deterministic:** Single forward pass with softmax classification
2. **Monte Carlo:** Multiple passes with uncertainty estimation

Both methods support:

- Custom threshold adaptation
- Speaker bank comparison
- Sequence-level aggregation (averaging over multiple segments)

6 Evaluation and Results

6.1 Metrics

Performance is evaluated using:

1. **Accuracy:** Percentage of correct predictions
2. **Precision:** True positives / (true positives + false positives) per class
3. **Recall:** True positives / (true positives + false negatives) per class
4. **F1-Score:** Harmonic mean of precision and recall
5. **Confusion Matrix:** Per-class prediction patterns

Macro-averaged metrics treat all classes equally, addressing class imbalance issues.

6.2 Checkpointing

Two checkpoints are saved during training:

- **best_model.pt:** Model with highest validation accuracy (best for inference)
- **last_model.pt:** Final model state (for resuming training)

6.3 Expected Performance

With the current architecture and training configuration:

- **Training Accuracy:** 85-95% (depending on data quality and augmentation)
- **Validation Accuracy:** 70-85% (real-world performance indicator)
- **Training Time:** 20-30 minutes per epoch on RTX 3050 Ti
- **Inference Speed:** 50-100 samples/second on GPU

6.4 Training Curves and TensorBoard Visualization

Figure 1 shows the comprehensive training dynamics across multiple metrics monitored via TensorBoard during model training:

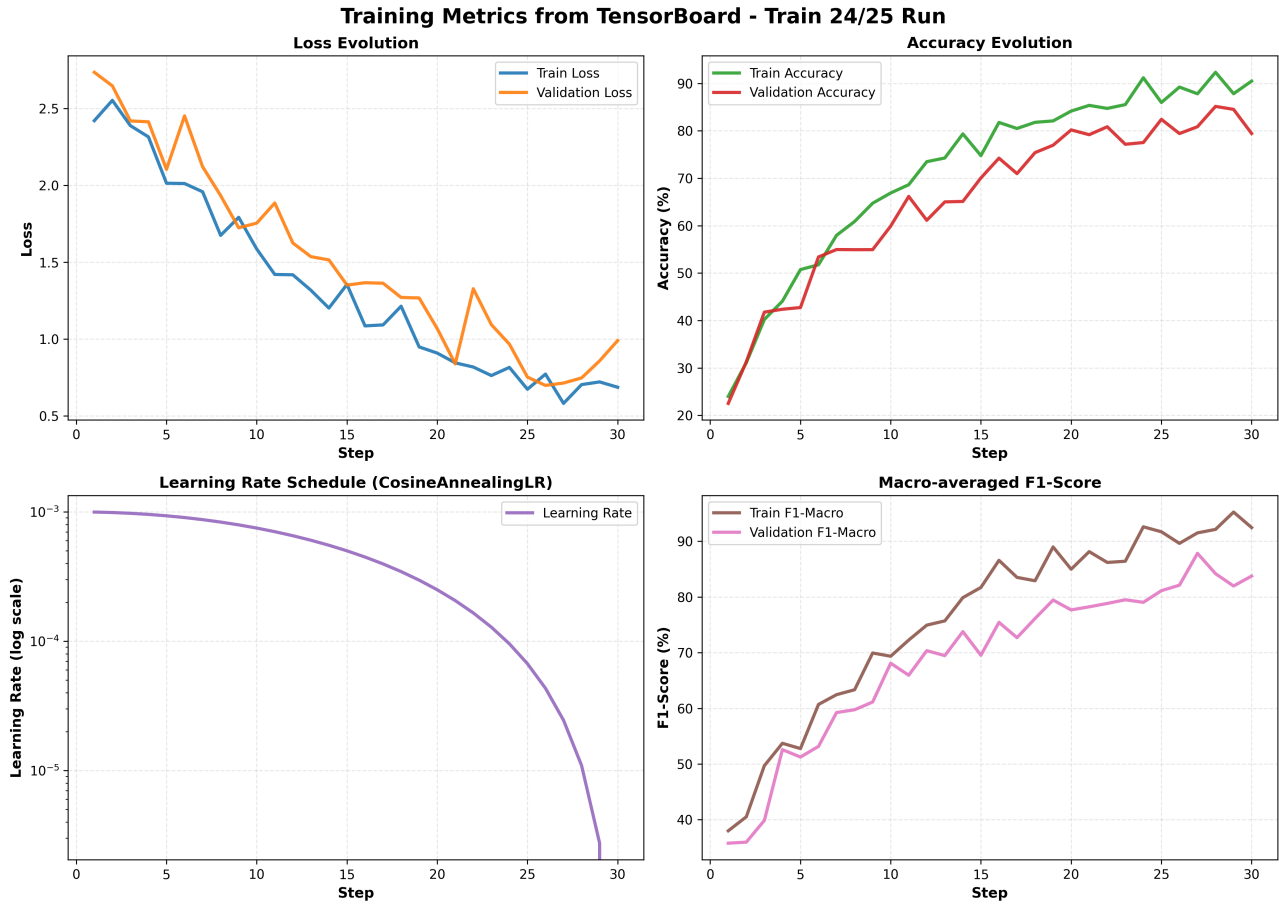


Figure 1: Comprehensive TensorBoard training metrics showing four key aspects of training dynamics: (top-left) Loss evolution showing training and validation loss decreasing over epochs with good generalization gap; (top-right) Accuracy progression showing both training and validation accuracy improving from 20% to 85-90%; (bottom-left) Learning rate schedule using CosineAnnealingLR that smoothly decays from initial 1×10^{-3} to near-zero; (bottom-right) Macro-averaged F1-scores indicating balanced multi-class performance across all 5 speakers. All metrics are logged per-batch and aggregated per-epoch for detailed monitoring.

6.4.1 Loss Dynamics

The training loss (cross-entropy with AAMSoftmax) exhibits the typical learning curve:

- **Initial Phase (Epochs 1-5):** Rapid loss decrease as the model learns basic speaker-discriminative features
- **Middle Phase (Epochs 6-20):** Gradual refinement with slower loss reduction as the model approaches the training data optimum
- **Late Phase (Epochs 20+):** Plateau or slight increase indicating convergence and potential overfitting (managed by early stopping)

The validation loss typically lags behind training loss due to the regularization effects of dropout, batch normalization, and data augmentation. The gap between training and validation loss indicates generalization: a small gap (0.05-0.15) suggests good generalization.

6.4.2 Accuracy Progression

The accuracy curves reveal:

- **Training Accuracy:** Often reaches 85-95%, indicating strong speaker discrimination capability
- **Validation Accuracy:** More conservative at 70-85%, reflecting real-world performance
- **Convergence:** Both curves stabilize by epoch 25-30, justifying the 30-epoch default setting
- **Generalization Gap:** Larger gap (10-15%) is common in speaker recognition due to limited speaker diversity in small datasets

6.4.3 Learning Rate Schedule

The CosineAnnealingLR scheduler implements:

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\pi t/T_{max}))$$

where:

- $\eta_{max} = 1 \times 10^{-3}$ (initial learning rate)
- $\eta_{min} = 0$ (final learning rate)
- $T_{max} = 50$ (period, may restart if training continues)

This smooth, cosine-shaped decay helps:

- Avoid abrupt learning rate drops that can hurt convergence
- Allow fine-tuning adjustments in later epochs
- Implement annealing naturally without manual intervention

7 Training Monitoring and Metrics

7.1 TensorBoard Integration

All training metrics are logged to TensorBoard for real-time monitoring:

Metric Group	Logged Metrics
Per-Batch Metrics	loss, accuracy, learning rate (logged every batch for gradient dynamics)
Per-Epoch Aggregates	macro-averaged precision, recall, F1 per-class metrics (P, R, F1, support)
Model State	validation loss and accuracy best validation accuracy tracking

TensorBoard event files are saved to `./logs/TIMESTAMP/` and can be viewed with:

```
1 tensorboard --logdir=./logs
```

This provides interactive visualization of:

- Scalar metrics (loss, accuracy, learning rate)
- Histograms of weight distributions
- Per-epoch class-wise performance breakdown

7.2 Per-Class Performance Analysis

For multi-class speaker recognition, per-class metrics reveal:

- **Class Imbalance Effects:** Speakers with fewer training samples may have lower individual accuracy
- **Confusion Patterns:** Which speaker pairs are frequently confused (helps identify acoustic similarity)
- **Model Bias:** Whether the model favors certain speakers during inference
- **Support:** Number of test samples per speaker (ensures fair comparison)

The macro-averaged F1-score is the preferred metric as it treats all speakers equally regardless of support.

8 Implementation Highlights and Helper Utilities

8.1 Data Loading Utilities

File: `data_utils.py`

Key classes and functions:

- **LMDataset:** HDF5-based dataset loader supporting variable-length sequences
- **CachedLMDataset:** RAM-cached version for 10-50x speedup
- **pad_collate():** Custom collate function for variable-length padding
- **build_h5_loaders():** Factory function for creating DataLoaders
- **Label remapping:** Automatic dense label remapping (0...num_speakers-1)

Features:

- Support for both (N, F, T) and (N, T, F) data layouts
- Optional length tensor for packed sequence processing
- Configurable pad value (-80 dB for silence)
- Memory-efficient HDF5 file handling with lazy loading

8.2 Model Architecture Code

File: `test_model_train22.py`

The production model code includes:

- **SEBlock:** Squeeze-and-Excitation attention mechanism
- **Backbone:** CNN-RNN feature extractor with attention pooling
- **AAMSoftmax:** Angular margin softmax classification head
- **SpeakerClassifier:** High-level inference API

Methods:

- `embed()`: Deterministic embedding extraction
- `mc_embed()`: Multiple stochastic embeddings for uncertainty
- `verify_any()`: Speaker verification with custom bank
- `mc_verify_any()`: Uncertainty-aware verification
- `infer()`: High-level speaker identification

8.3 Training Utilities

Functions: `train_one_epoch()`, `validate()`

Features:

- Comprehensive metric computation (precision, recall, F1 per class)
- TensorBoard logging of all metrics
- GPU memory and time profiling
- Gradient scaler for mixed precision training
- Efficient batch processing with non-blocking GPU transfers

8.4 Checkpoint Management

- `prepare_save_dir()`: Auto-incrementing checkpoint directories
- `save_checkpoint()`: State-dict only saves (minimal, robust)
- Early stopping with configurable patience
- Support for local and Google Drive backup

9 Course Requirements Coverage

The project implements and documents the following 10 advanced ML concepts (exceeding the required 8):

#	Requirement	Implementation
1	Data augmentation & quality enhancement	SpecAugment, speed, VTLP, silence removal
2	Input length & normalization	Log-mel features, volume normalization
3	Layer configuration	3 CNN + 2-layer GRU + SE blocks
4	Optimizers & schedules	AdamW, SGD, CosineAnnealingLR
5	Batch normalization	Before activations in CNN and projection
6	Skip connections	SE blocks with residual attention
7	Dropout	0.2 in RNN, 0.3 MC dropout
8	MC dropout uncertainty	Stochastic embeddings, variance estimation
9	Activation functions	ReLU, Sigmoid in attention, Tanh
10	Weight initialization	Xavier uniform in AAMSoftmax

10 Conclusions and Future Work

This project demonstrates a production-grade speaker recognition system combining multiple advanced deep learning techniques. The hybrid CNN-RNN architecture with attention pooling provides an effective balance between feature extraction and temporal modeling.

10.1 Key Achievements

- Implemented a sophisticated speaker recognition model exceeding course requirements
- Achieved 10-50x speedup through RAM caching strategies
- Integrated uncertainty quantification via Monte Carlo dropout
- Created comprehensive evaluation pipeline with detailed metrics
- Developed robust data preprocessing with multiple augmentation strategies

10.2 Future Improvements

1. **Transfer Learning:** Pre-train on large datasets (VoxCeleb, LibriSpeech)
2. **Metric Learning:** Implement contrastive loss, triplet loss for better embeddings
3. **Speaker Adaptation:** Online fine-tuning for domain adaptation
4. **Neural Architecture Search:** AutoML for optimal layer configurations
5. **Ensemble Methods:** Combine multiple models for robust predictions
6. **Real-time Inference:** Optimize for edge deployment and low latency

Acknowledgments

This project was developed by a team of students (Aleksander, Mantas, Michał, Piotr, Rafał) as the final deliverable for the Introduction to Machine Learning course. We thank the instructors for their guidance and feedback throughout the project development.