

# Comprehensive Report on Speech Command Classification

Eugene Romanov

March 20, 2025

## 1 Introduction

This report details the implementation and analysis of a neural network-based system for speech command classification using mel-filterbank features. The project consists of two main parts:

1. Development of a customized `LogMelFilterBanks` class
2. Training and evaluating CNN models with various configurations

The goal was to explore how different parameters affect model performance, computational efficiency, and training time.

## 2 Part 1: Implementation of `LogMelFilterBanks`

### 2.1 Overview

We implemented a PyTorch module `LogMelFilterBanks` that transforms audio signals into log mel-filterbank features. This module follows a pipeline of:

1. Calculating Short-Time Fourier Transform (STFT)
2. Converting to power spectrum
3. Applying mel-filterbanks
4. Taking the logarithm of the result

### 2.2 Implementation Details

The implementation includes these key methods:

- `__init__`: Initializes parameters for STFT and mel-filterbanks
- `_init_melscale_fbanks`: Creates the mel-filterbank matrix
- `spectrogram`: Calculates the STFT of the input signal
- `forward`: Processes the audio signal to produce log mel-filterbank features

We used PyTorch's built-in functions like `torch.stft` for the STFT calculation and `torchaudio.functional.melscale_fbanks` for creating the mel-filterbank matrix. The power spectrum was calculated by taking the squared magnitude of the complex STFT output.

## 2.3 Validation

To verify the correctness of our implementation, we compared it with the native `torchaudio.transforms.MelSpec` implementation:

```
1 our_output = our_melspec(waveform)
2 torchaudio_output = torchaudio_melspec(waveform)
3 torchaudio_log_output = torch.log(torchaudio_output + 1e-6)
```

Visual inspection of the spectrograms and numerical comparison confirmed that our implementation produces identical results to the native implementation:

- Mean absolute difference: 0.000000
- Maximum absolute difference: 0.000000

## 3 Part 2: CNN Model for Binary Speech Command Classification

### 3.1 Dataset Preparation

We used the Google Speech Commands dataset from `torchaudio` and converted it into a binary classification problem by selecting only "yes" and "no" samples:

```
1 class BinarySpeechCommands(Dataset):
2     def __init__(self, root_dir="./data", subset="training", n_mels=80):
3         self.dataset = SPEECHCOMMANDS(root=root_dir, subset=subset, download=True)
4
5         # Filter only "yes" and "no" samples
6         self.indices = []
7         for i in range(len(self.dataset)):
8             waveform, sample_rate, label, *_ = self.dataset[i]
9             if label.lower() in ["yes", "no"]:
10                self.indices.append(i)
```

This resulted in:

- 6,358 training examples
- 803 validation examples
- 824 testing examples

### 3.2 Model Architecture

We designed a CNN architecture using PyTorch Lightning:

```
1 class SpeechCNN(pl.LightningModule):
2     def __init__(self, n_mels=80, groups=1):
3         super(SpeechCNN, self).__init__()
4
5         # CNN layers
6         self.conv1 = nn.Conv2d(1, 16, kernel_size=3, stride=1, padding=1, groups=1)
7         self.bn1 = nn.BatchNorm2d(16)
8         self.pool1 = nn.MaxPool2d(kernel_size=2)
9
10        self.conv2 = nn.Conv2d(16, 32, kernel_size=3, stride=1, padding=1, groups=
groups)
```

```

11     self.bn2 = nn.BatchNorm2d(32)
12     self.pool2 = nn.MaxPool2d(kernel_size=2)
13
14     self.conv3 = nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1, groups=
groups)
15     self.bn3 = nn.BatchNorm2d(64)
16     self.pool3 = nn.MaxPool2d(kernel_size=2)
17
18     # Adaptive pooling and fully connected layers
19     self.adaptive_pool = nn.AdaptiveAvgPool2d((4, 4))
20     self.fc1 = nn.Linear(64 * 4 * 4, 128)
21     self.dropout = nn.Dropout(0.5)
22     self.fc2 = nn.Linear(128, 2)

```

The model includes:

- Three convolutional blocks with batch normalization and max pooling
- Adaptive average pooling to handle variable input sizes
- Two fully connected layers with dropout for regularization

### 3.3 Training and Evaluation

We implemented a comprehensive training and evaluation pipeline that:

1. Tracks training loss, validation accuracy, and epoch training time
2. Calculates model parameters and FLOPs
3. Tests the model on a separate test set
4. Handles variable-length audio inputs through custom collation

## 4 Experiments

### 4.1 Experiment 1: Varying the Number of Mel-Filterbanks

We trained models with different numbers of mel-filterbanks ( $n\_mels \in \{20, 40, 80\}$ ) and analyzed their performance:

<b>n_mels</b>	<b>Test Accuracy</b>	<b>Parameters</b>	<b>FLOPs</b>	<b>Avg Epoch Time</b>
20	99.39%	154,978	5,253,888	56.93s
40	98.91%	154,978	10,376,192	70.32s
80	98.18%	154,978	20,620,032	99.67s

Table 1: Performance metrics for different numbers of mel-filterbanks

Key observations:

- Lower n\_mels (20) achieved the highest accuracy
- FLOPs increased proportionally with n\_mels
- Training time increased with n\_mels
- Parameter count remained constant across different n\_mels values

## 4.2 Experiment 2: Varying the Groups Parameter

Using the best `n_mels` value (20), we experimented with different group convolution settings ( $groups \in \{1, 2, 4, 8\}$ ):

groups	Test Accuracy	Parameters	FLOPs	Avg Epoch Time
1	98.67%	154,978	5,253,888	57.97s
2	97.69%	143,458	2,949,888	54.79s
4	98.54%	137,698	1,797,888	60.04s
8	98.06%	134,818	1,221,888	59.02s

Table 2: Performance metrics for different groups values

Key observations:

- Highest accuracy was achieved with `groups=1`
- The relationship between groups and accuracy was non-monotonic
- Parameter count and FLOPs decreased as groups increased
- Training time showed no clear trend with changing groups

## 5 Analysis and Conclusions

### 5.1 Impact of Mel-Filterbanks

The experiments demonstrate that for this binary speech command classification task, a smaller number of mel-filterbanks (`n_mels=20`) is not only more computationally efficient but also yields better accuracy. This suggests that:

1. Lower dimensional features are sufficient to capture the relevant information for distinguishing between "yes" and "no" commands
2. Higher dimensional features may introduce noise or lead to overfitting
3. The computational cost increases substantially with higher `n_mels` values (4x increase from `n_mels=20` to `n_mels=80`) without providing any benefit in accuracy

### 5.2 Impact of Grouped Convolutions

Grouped convolutions offer a compelling trade-off between computational efficiency and accuracy:

1. Standard convolutions (`groups=1`) achieved the highest accuracy, but at the highest computational cost
2. Using `groups=4` reduced FLOPs by approximately 3x with only a 0.13% drop in accuracy
3. `Groups=8` offered the greatest reduction in FLOPs (4.3x) with a moderate 0.61% accuracy drop
4. The non-monotonic relationship between groups and accuracy (`groups=4` outperforming `groups=2`) suggests that certain grouped architectures may be more effective for this task

### 5.3 Optimal Configuration

For this binary speech command classification task:

1. **Best Performance:** n\_mels=20 with groups=1 provides the highest accuracy (99.39%)
2. **Balanced Efficiency:** n\_mels=20 with groups=4 offers excellent accuracy (98.54%) with substantially reduced computational requirements
3. **Maximum Efficiency:** n\_mels=20 with groups=8 provides good accuracy (98.06%) with minimal computational cost

## 6 Implementation Challenges and Solutions

During implementation, we faced several challenges:

1. **Variable-length audio inputs:** Solved by implementing a custom collation function that pads spectrograms to a uniform size within each batch
2. **Dimension mismatch in matrix operations:** Resolved by carefully managing tensor dimensions in the `forward` method of `LogMelFilterBanks`
3. **Efficient batch processing:** Addressed using PyTorch’s broadcasting capabilities and batch matrix multiplication

## 7 Future Work

Potential directions for further research:

1. Extending the model to multi-class classification with more speech commands
2. Exploring different CNN architectures, such as residual networks or attention mechanisms
3. Implementing data augmentation techniques to improve robustness
4. Deploying the model on resource-constrained devices using the efficient configurations discovered

## 8 Summary

This project successfully implemented and evaluated a speech command classification system using custom mel-filterbank features and CNN models. Through systematic experimentation with different parameters, we identified configurations that balance accuracy and computational efficiency. The findings demonstrate that carefully tuned smaller models can achieve excellent performance on this task, which has important implications for deploying such systems in resource-constrained environments.