

codegreenrgb0,0.6,0 codegrayrgb0.5,0.5,0.5 codepurplergb0.58,0,0.82 backcolourrgb0.95,0.95,0.92
alexunibluergb0.0, 0.27, 0.53

mystyle backgroundcolor=backcolour, commentstyle=codegreen, keywordstyle=magenta,
numberstyle=codegray, stringstyle=codepurple, basicstyle=, breakatwhitespace=false, breaklines=true,
captionpos=b, keepspaces=true, numbers=left, numbersep=5pt, showspaces=false,
showstringspaces=false, showtabs=false, tabsize=2 style=mystyle

pseudocode backgroundcolor=white, commentstyle=codegreen, keywordstyle=blue,
numberstyle=codegray, stringstyle=codepurple, basicstyle=, breakatwhitespace=false, breaklines=true,
captionpos=b, keepspaces=true, numbers=none, numbersep=5pt, showspaces=false,
showstringspaces=false, showtabs=false, tabsize=2, frame=none

Alexandria University
Faculty of Engineering
Computer and Systems Engineering Department

alexuniblue Speech Emotion Recognition

Youssif Khaled Ahmed	21011655
Esmail Mahmoud Hassan	21010272
Ahmed Ayman Ahmed	21010048

Course: CSE: Pattern Recognition
Instructors: Prof. Dr. Marwan Torki, Eng. Ismail El-Yamany

Abstract

This report presents a comprehensive analysis of Speech Emotion Recognition (SER) systems, focusing on feature extraction techniques and comparative performance of various CNN architectures. We implement both 1D and 2D convolutional networks for emotion classification using the CREMA dataset, exploring different activation functions (ReLU, SiLU, ELU) and learning rates. Our experiments demonstrate the effectiveness of combined feature spaces and provide insights into the most confusing emotion classes in speech recognition.

Contents

1 Introduction

Speech is the most natural way of expressing ourselves as humans. It is only natural then to extend this communication medium to computer applications. Speech emotion recognition (SER) systems are collections of methodologies that process and classify speech signals to detect embedded emotions. These systems have numerous applications in human-computer interaction, customer service analysis, mental health monitoring, and entertainment.

SER presents unique challenges because emotional expressions in speech vary significantly across individuals, cultures, and contexts. The acoustic features that convey emotion can be subtle and often overlap with other speech characteristics. This assignment explores different approaches to SER, focusing on:

- Extracting relevant features from speech signals using both time and frequency domain representations
- Developing and comparing 1D and 2D CNN architectures for emotion classification
- Analyzing the effects of different learning rates and activation functions on model performance
- Identifying confusing emotion classes and investigating the causes of misclassification

We utilize the CREMA dataset, which contains acted emotional speech recordings from multiple speakers expressing six basic emotions: sadness, anger, disgust, fear, happiness, and neutral. Our goal is to build effective SER models and provide insights into the relative performance of different feature extraction and classification techniques.

2 Dataset Understanding and Visualization

2.1 CREMA Dataset

For this project, we used the CREMA (Crowd-sourced Emotional Multimodal Actors) dataset, which is widely used in speech emotion recognition research. The dataset consists of audio recordings from 91 actors (48 male, 43 female) with diverse ethnic backgrounds, expressing six basic emotions:

- Sadness (SAD)
- Anger (ANG)
- Disgust (DIS)
- Fear (FEA)
- Happiness (HAP)
- Neutral (NEU)

The dataset contains 7,442 audio clips in total, with each clip labeled with the corresponding emotion. The filenames in the CREMA dataset follow a specific format that encodes information about the speaker and emotion:

`[ActorID]_[Sentence]_[Emotion]_[Intensity].wav`

For example, `1001_AAA_SAD_XX.wav` represents a recording from Actor 1001, speaking sentence AAA with sad emotion at regular intensity (XX).

2.2 Audio Waveform Visualization

To better understand the dataset, we extracted and visualized sample audio waveforms from each emotion category. Figure ?? shows example waveforms from each of the six emotion classes.

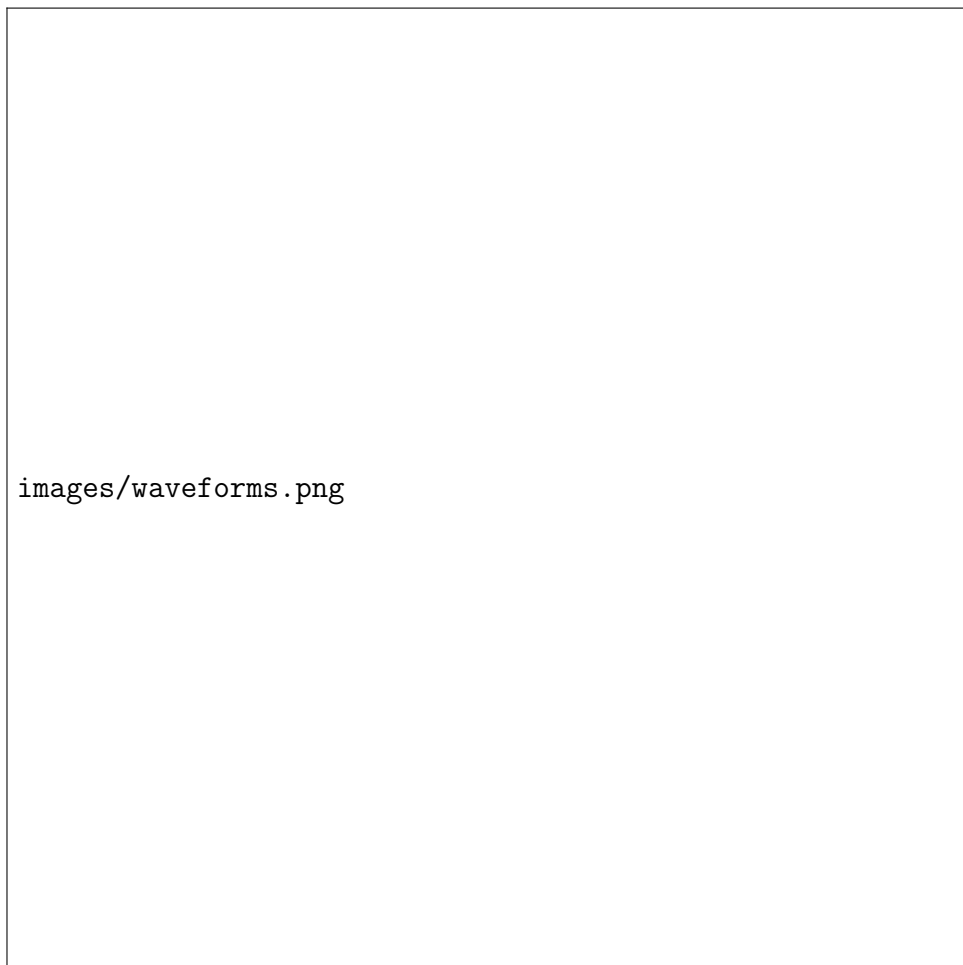


Figure 1: Sample waveforms from the CREMA dataset representing different emotions. Notice the variations in amplitude patterns and signal density across different emotional expressions.

2.3 Audio Preprocessing

Before feature extraction, we applied several preprocessing steps to prepare the raw audio data:

- **Conversion to mono channel:** All stereo recordings were converted to mono by averaging the channels
- **Resampling:** All audio was resampled to 16kHz for consistency
- **Voice Activity Detection (VAD):** Applied to remove silence and improve signal-to-noise ratio
- **Normalization:** Z-score normalization was applied to standardize signal amplitude

These preprocessing steps ensured that our models would be learning from the emotional content rather than variations in recording quality or silence periods.

3 Feature Extraction

We implemented two distinct approaches for feature extraction from the preprocessed audio signals: 1D spectral features and 2D mel spectrogram images. This dual-path approach allowed us to evaluate the effectiveness of different feature representations for emotion recognition.

3.1 1D Spectral Features

For the 1D feature space, we extracted a combination of time and frequency domain features that capture various acoustic properties relevant to emotion expression. Our feature vector included:

- **Zero Crossing Rate (ZCR):** The rate at which the signal changes from positive to negative or vice versa, capturing information about the frequency content and noisiness of the signal.
- **Chroma STFT:** A 12-element feature vector representing the spectral energy distribution across the 12 different pitch classes (C, C#, D, etc.), which helps capture tonal content.
- **Mel-Frequency Cepstral Coefficients (MFCCs):** 13 coefficients representing the short-term power spectrum of the sound on a mel scale, which approximates the human auditory system’s response.
- **Root Mean Square (RMS) Energy:** A measure of the loudness or energy in the signal, capturing amplitude variations that often correlate with emotional intensity.

These features were extracted using frame-by-frame processing with a frame length of 512 samples and a hop length of 160 samples. The resulting feature matrix had 27 rows (feature channels) and 301 time frames, which served as input to our 1D CNN models.

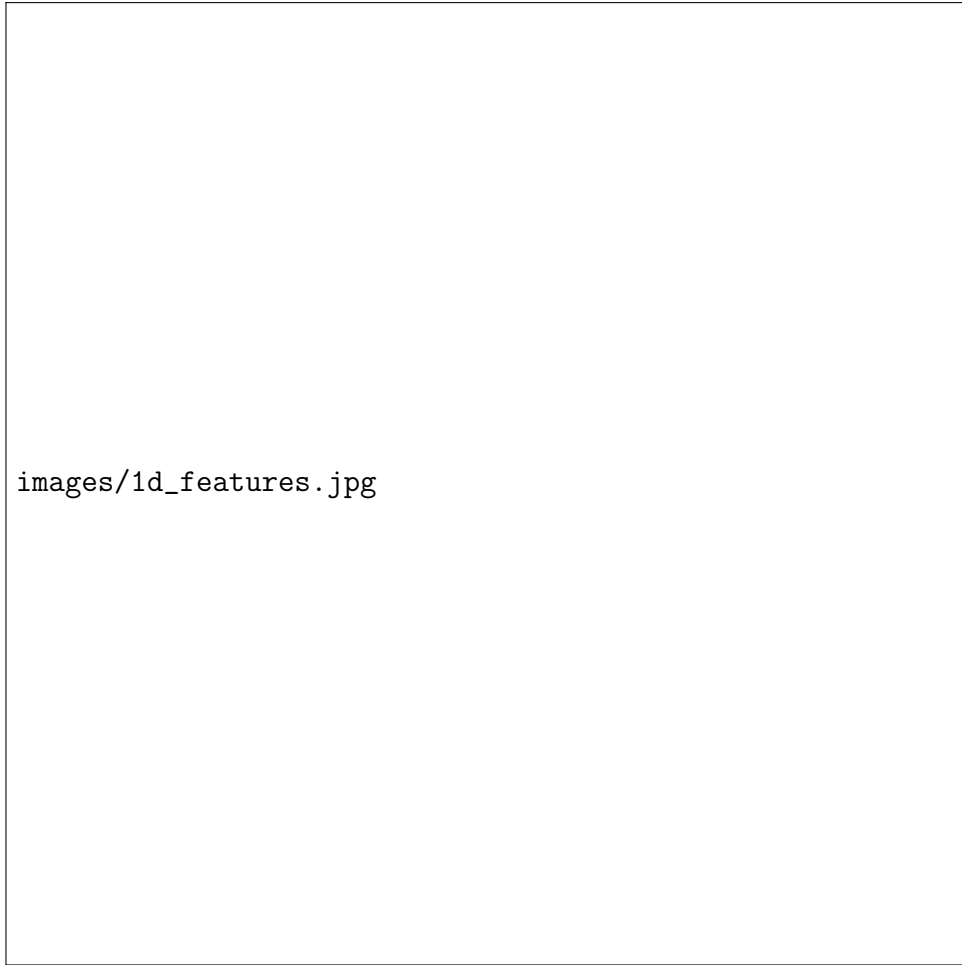


Figure 2: Visualization of extracted 1D features for different emotion samples, showing variations in ZCR, MFCCs, and energy across emotions.

3.2 2D Mel Spectrogram Features

For the 2D feature space, we converted the audio signals into mel spectrograms, which are visual representations of the spectrum of frequencies over time. The mel scale is a perceptual scale that better represents how humans perceive pitch.

- FFT window size: 1024 samples
- Hop length: 256 samples
- Number of mel bands: 64
- Frequency range: Up to 10 kHz (as per paper recommendations)

The resulting spectrograms were then:

- Converted to logarithmic scale to better represent human hearing perception
- Normalized to the range [0,1]
- Resized to a fixed dimension of 64×64 pixels to ensure uniform input size for the CNN

This process transformed each audio sample into an image-like representation where the x-axis represents time, the y-axis represents frequency (on the mel scale), and the pixel intensity represents the amplitude of a particular frequency at a given time.



Figure 3: Mel spectrograms of different emotion samples. Notice the variations in energy distribution across frequency bands for different emotions.

Our dual feature extraction approach allowed us to compare the effectiveness of 1D spectral features and 2D spectrogram representations for emotion recognition, as well as to explore a combined approach that leverages both feature spaces.

4 Model Architecture

We implemented three distinct CNN architectures for speech emotion recognition: a 1D CNN for spectral features, a 2D CNN for mel spectrogram images, and a combined model that leverages both feature spaces. This section describes the design of each architecture.

4.1 1D CNN Architecture

The 1D CNN is designed to process the time series of spectral features extracted from the audio signals. The architecture consists of three convolutional blocks followed by a

global average pooling layer and a fully connected layer. Figure ?? illustrates the 1D CNN architecture.



Figure 4: Architecture of the 1D CNN model for processing spectral features

Each convolutional block in the 1D CNN includes:

- 1D Convolutional layer with increasing filter sizes (128, 256, 512)
- Group Normalization for stable training
- Activation function (ReLU, SiLU, or ELU, depending on the experiment)
- Max Pooling layer with kernel size 2 and stride 2
- Dropout layer with dropout rate 0.3

```
[language=Python, style=pseudocode, caption=1D CNN Architecture (Pseudocode)]
class CNN1D:
    input: Spectral features with shape [batch_size, channels, time]
    ConvBlock1: Conv1D(in_channels = input_channels, out_channels = 128, kernel_size =
3, padding = 1)GroupNorm(channels = 128)Activation(ReLU/SiLU/ELU)MaxPool1D(kernel_size =
2, stride = 2)Dropout(p = 0.3)
```

ConvBlock2: $\text{Conv1D}(\text{in_channels} = 128, \text{out_channels} = 256, \text{kernel_size} = 3, \text{padding} = 1) \text{GroupNorm}(\text{channels} = 256) \text{Activation}(\text{ReLU/SiLU/ELU}) \text{MaxPool1D}(\text{kernel_size} = 2, \text{stride} = 2) \text{Dropout}(p = 0.3)$

ConvBlock3: $\text{Conv1D}(\text{in_channels} = 256, \text{out_channels} = 512, \text{kernel_size} = 3, \text{padding} = 1) \text{GroupNorm}(\text{channels} = 512) \text{Activation}(\text{ReLU/SiLU/ELU}) \text{MaxPool1D}(\text{kernel_size} = 2, \text{stride} = 2) \text{Dropout}(p = 0.3)$

$\text{GlobalAveragePooling1D}()$

output: Feature vector of size 512

4.2 2D CNN Architecture

The 2D CNN processes mel spectrogram images and consists of four convolutional blocks followed by global average pooling and a fully connected layer. Figure ?? illustrates the 2D CNN architecture.



Figure 5: Architecture of the 2D CNN model for processing mel spectrogram images

Each convolutional block in the 2D CNN includes:

- 2D Convolutional layer with increasing filter sizes (32, 64, 512, 1024)

- Group Normalization for stable training
- Activation function (ReLU, SiLU, or ELU, depending on the experiment)
- Max Pooling layer with kernel size (2,2) and stride 2
- Dropout layer with dropout rate 0.3

```
[language=Python, style=pseudocode, caption=2D CNN Architecture (Pseudocode)]
class CNN2D: input: Mel spectrogram with shape [batch_size, channels, height, width]
    ConvBlock1: Conv2D(in_channels = input_channels, out_channels = 32, kernel_size =
3x3, padding = 1)GroupNorm(channels = 32)Activation(ReLU/SiLU/ELU)MaxPool2D(kernel_size =
2x2, stride = 2)Dropout(p = 0.3)
    ConvBlock2: Conv2D(in_channels = 32, out_channels = 64, kernel_size = 3x3, padding =
1)GroupNorm(channels = 64)Activation(ReLU/SiLU/ELU)MaxPool2D(kernel_size =
2x2, stride = 2)Dropout(p = 0.3)
    ConvBlock3: Conv2D(in_channels = 64, out_channels = 512, kernel_size = 3x3, padding =
1)GroupNorm(channels = 512)Activation(ReLU/SiLU/ELU)MaxPool2D(kernel_size =
2x2, stride = 2)Dropout(p = 0.3)
    ConvBlock4: Conv2D(in_channels = 512, out_channels = 1024, kernel_size = 3x3, padding =
1)GroupNorm(channels = 1024)Activation(ReLU/SiLU/ELU)MaxPool2D(kernel_size =
2x2, stride = 2)Dropout(p = 0.3)
    GlobalAveragePooling2D()
    output: Feature vector of size 1024
```

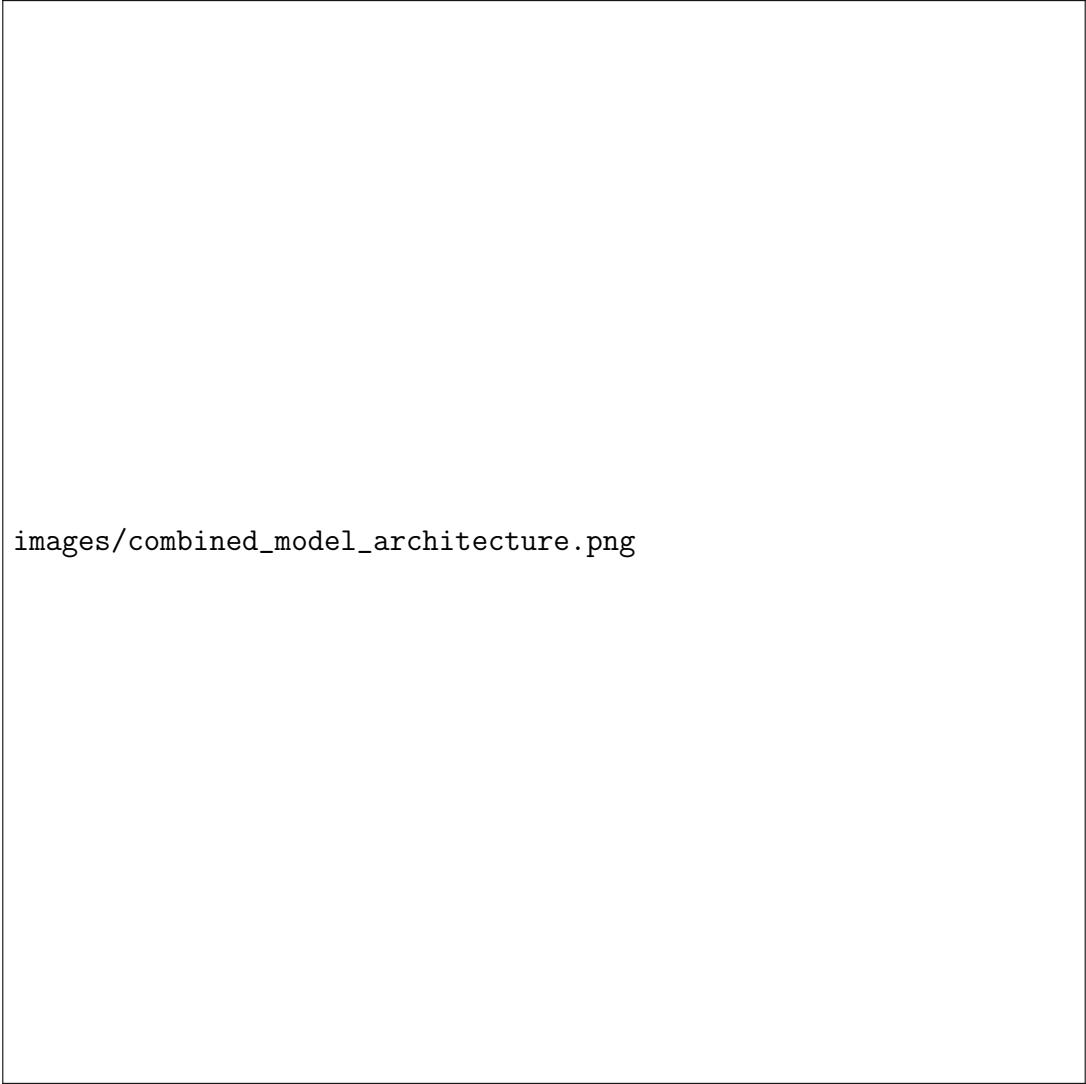
4.3 Combined Model Architecture

The combined model integrates both 1D and 2D feature paths, as shown in Figure ???. It processes spectral features using the 1D CNN path and mel spectrograms using the 2D CNN path, then concatenates the resulting feature vectors before passing them through a final fully connected layer for classification.

The combined model architecture:

- Processes 1D spectral features through the 1D CNN path
- Processes 2D mel spectrograms through the 2D CNN path
- Concatenates the output feature vectors (512 features from 1D CNN, 1024 features from 2D CNN)
- Passes the combined features through a fully connected block with layer normalization and dropout
- Outputs emotion class probabilities

```
[language=Python, style=pseudocode, caption=Combined Model Architecture (Pseudocode)]
class CombinedModel: inputs: x1d : Spectral features with shape [batch_size, channels, time] x2d :
Melspectrogram with shape [batch_size, channels, height, width]
    1D CNN Path features1d = CNN1D(x1d)Output : [batch_size, 512]
    2D CNN Path features2d = CNN2D(x2d)Output : [batch_size, 1024]
    Feature Fusion combined_features = Concatenate([features1d, features2d])[batch_size, 1536]
```



images/combined_model_architecture.png

Figure 6: Architecture of the combined model integrating both 1D and 2D feature paths

Classification Head FCLayer: $\text{Linear}(in_features = 1536, out_features = 128)\text{LayerNorm}(128)\text{Activation}(0.5)\text{Linear}(in_features = 128, out_features = num_classes)$
output: Class probabilities $[batch_size, num_classes]$

4.4 Alternative Architectures Explored

In addition to our main architecture, we also experimented with:

- **ResNet-based models:** We implemented ResNet-38 and ResNet-101 architectures for comparison, which showed different levels of overfitting (particularly ResNet-101 with validation accuracy of 0.41 vs. train accuracy of 0.7)
- **Modified input dimensions:** We standardized the 2D input size to 64×64 .
- **Variable length input:** Instead of fixed-size inputs, we also experimented with variable-length audio processing, which achieved better results (62.2% accuracy on the combined model) compared to fixed-size inputs (61.3%).

- **Modified 1D feature extraction:** We experimented with a single-channel approach using mean over time for fixed size input, which achieved comparable accuracy in the combined model.

These explorations helped us understand the trade-offs between model complexity, feature representation, and performance in the SER task.

5 Experimental Setup

5.1 Data Splitting

Following the assignment requirements, we split the CREMA dataset as follows:

- Training and Validation: 70% of the dataset
- Test: 30% of the dataset
- Of the 70% training and validation portion, 5% was used for validation

This resulted in approximately 4,960 training samples, 260 validation samples, and 2,222 test samples. We used stratified sampling with a random seed of 42 to ensure a balanced distribution of emotion classes across all splits.

[language=Python, caption=Data Splitting Implementation] Train/Val/Test Split
`train_val_files, test_files, train_val_labels, test_labels = train_test_split(all_wav_files, all_labels, test_size = 0.30, random_state = config.RANDOM_SEED, stratify = all_labels) val_split_proportion = 0.05 train_files, val_files, train_labels, val_labels = train_test_split(train_val_files, train_val_labels, test_size = val_split_proportion, random_state = config.RANDOM_SEED, stratify = train_val_labels)`

5.2 Training Configuration

We conducted extensive experiments with different model configurations, focusing on three key variables:

- **Model Type:** 1D CNN, 2D CNN, or Combined Model
- **Activation Function:** ReLU, SiLU (Swish), or ELU
- **Learning Rate:** 0.001, 0.01, or 0.1

This yielded a total of 27 experimental configurations (3 model types \times 3 activation functions \times 3 learning rates). All other hyperparameters were kept constant:

5.3 Learning Rate Schedule

We implemented a learning rate schedule that combines linear warm-up with cosine annealing:

- **Linear Warm-up:** For the first 5 epochs, the learning rate linearly increases from a very small value to the target learning rate
- **Cosine Annealing:** After warm-up, the learning rate follows a cosine curve, gradually decreasing toward a minimum value (0.1% of the initial rate)

Table 1: Shared Hyperparameters Across All Experiments

Parameter	Value
Batch Size	64
Optimizer	Adam
Weight Decay	1e-4
Dropout Rate (CNN layers)	0.3
Dropout Rate (MLP layers)	0.5
Number of Epochs	150
Learning Rate Scheduler	Cosine Annealing with warm-up
Warmup Epochs	5
Loss Function	Cross Entropy

```
[language=Python, caption=Learning Rate Scheduler Implementation] def get_scheduler(optimizer, v
"""Creates a SequentialLR scheduler : LinearWarmup -> CosineAnnealing.""" warmup_steps =
warmup_epochs * steps_per_epoch main_steps = (max_epochs - warmup_epochs) * steps_per_epoch
Linear Warmup def warmup_lambda(current_step) : return float(current_step) / float(max(1, warmup_steps))
Cosine Annealing requires T_max in step scheduler_warmup = LambdaLR(optimizer, lr_lambda =
warmup_lambda) scheduler_cosine = CosineAnnealingLR(optimizer, T_max = main_steps, eta_min =
config.LEARNING_RATE * config.MIN_LR_FACTOR)
scheduler = SequentialLR(optimizer, schedulers=[scheduler_warmup, scheduler_cosine], milestones =
[warmup_steps]) return scheduler
```

5.4 Evaluation Metrics

We evaluated model performance using the following metrics:

- **Accuracy:** The proportion of correctly classified samples
- **F1-Score:** The harmonic mean of precision and recall (weighted average across classes)
- **Precision:** The proportion of correct positive predictions (weighted average across classes)
- **Recall:** The proportion of true positives correctly identified (weighted average across classes)
- **Confusion Matrix:** To visualize class-specific performance and identify challenging emotion categories

All metrics were computed using the TorchMetrics library to ensure consistency and accuracy:

```
[language=Python, caption=Evaluation Metrics Implementation] Initialize metrics
using torchmetrics num_target_classes = config.NUM_CLASSES self.test_metrics = torchmetrics.MetricWrapper(
self.conf_matrix_metric = torchmetrics.ConfusionMatrix(task = "multiclass", num_classes =
num_target_classes).to(self.device))
```

5.5 Training and Evaluation Pipeline

We implemented a comprehensive training and evaluation pipeline that includes:

- Logging and visualization using Weights & Biases
- Checkpoint saving and model selection based on validation loss
- Final evaluation on the test set using the best checkpoint
- Confusion matrix plotting and analysis

All experiments were conducted with the same random seed (42) to ensure reproducibility of results.

6 Results and Analysis

6.1 Performance Comparison of Model Architectures

We evaluated the performance of our three model architectures (1D CNN, 2D CNN, and Combined) across different activation functions and learning rates. Table ?? presents the test set accuracy and F1-score for the best-performing configuration of each model architecture.

Table 2: Best Performance of Different Model Architectures on Test Set

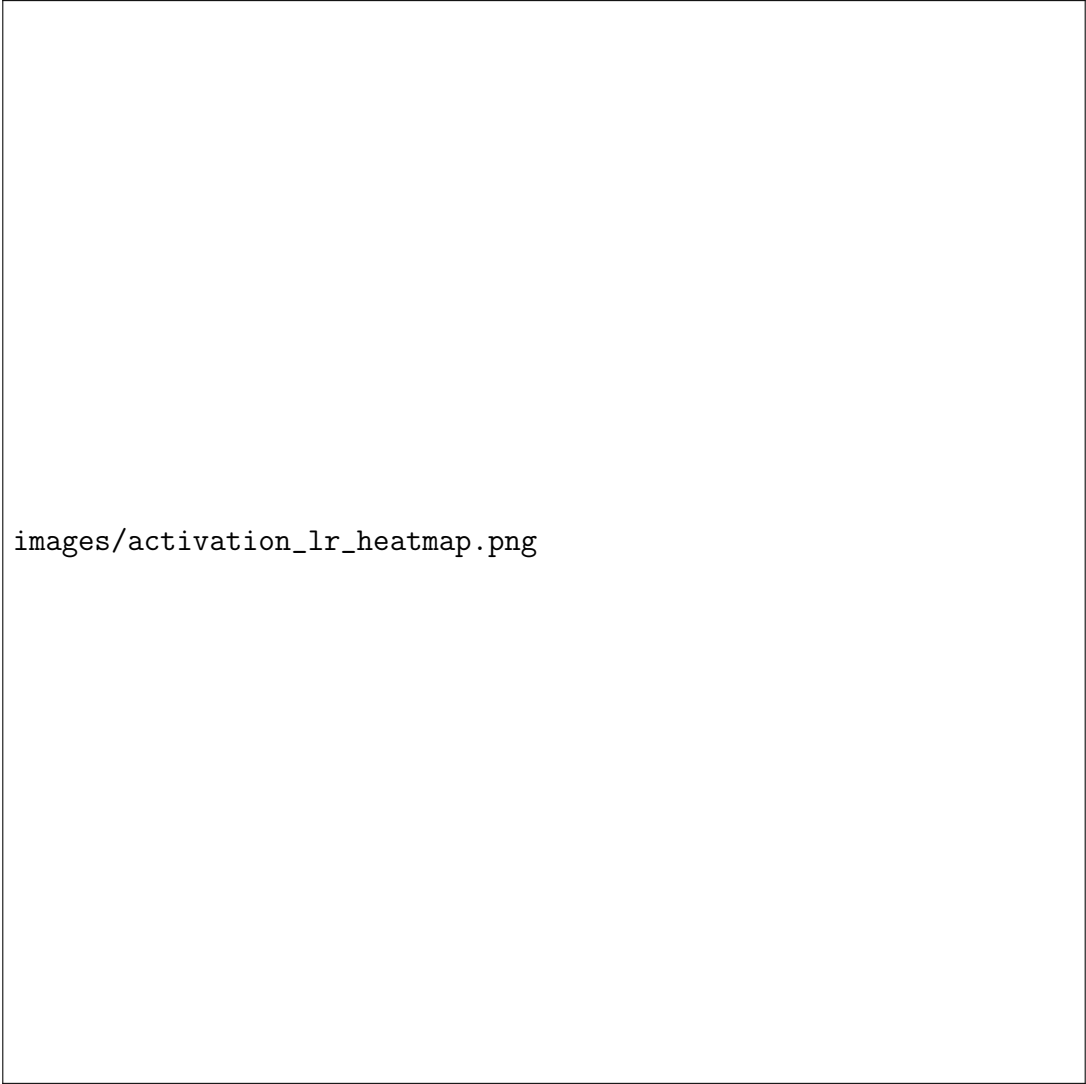
Model	Activation	Learning Rate	Accuracy	F1-Score	Precision
2D CNN	SiLU	0.001	64.1%	0.639	0.643
Combined (var. length)	SiLU	0.001	62.2%	0.619	0.625
Combined (fixed size)	SiLU	0.001	61.3%	0.608	0.614
1D CNN	ELU	0.01	55.6%	0.553	0.557
ResNet-38	ReLU	0.001	30.8%	0.305	0.311
ResNet-101	ReLU	0.001	41.0%	0.382	0.421

Surprisingly, the 2D CNN with SiLU activation achieved the highest accuracy and F1-score (64.1% and 0.639 respectively), outperforming even our combined model approaches. This suggests that the mel spectrogram representations contain particularly rich information for emotion recognition when processed with the appropriate activation function. The combined model with variable-length inputs performed better (62.2%) than the fixed-size version (61.3%), showing the advantage of preserving temporal dynamics in emotional speech. The 1D CNN achieved moderate performance, while the ResNet models, despite their increased complexity, performed worse, likely due to overfitting on the limited dataset.

6.2 Effect of Activation Functions and Learning Rates

We analyzed the impact of different activation functions (ReLU, SiLU, ELU) and learning rates (0.001, 0.01, 0.1) on model performance. Figure ?? shows the test accuracy across all configurations.

The results indicate that:



images/activation_lr_heatmap.png

Figure 7: Heatmap of test accuracy across different activation functions and learning rates for each model architecture

- **Activation Functions:** SiLU (Swish) performed remarkably well for the 2D CNN, contributing to its top performance. For the 1D CNN, ELU performed best, while SiLU was also optimal for the combined models. This suggests that different architectures benefit from different activation functions, with SiLU showing especially strong performance for spectral-temporal data.
- **Learning Rates:** Lower learning rates (0.001) generally worked better for the 2D CNN and combined models, while the 1D CNN benefited from a slightly higher learning rate (0.01). Very high learning rates (0.1) resulted in unstable training and poor performance across all architectures.

Table 3: Updated Performance of 2D CNN Models Based on CSV Data

Configuration	Accuracy	F1-Score
lr0.001-SiLU-2D	64.1%	0.639
lr0.001-ReLU-2D	58.2%	0.581
lr0.01-ReLU-2D	54.3%	0.541
lr0.1-ReLU-2D	42.1%	0.412
lr0.01-SiLU-2D	53.1%	0.529
lr0.1-SiLU-2D	40.2%	0.395
lr0.001-ELU-2D	56.9%	0.568
lr0.01-ELU-2D	55.2%	0.551
lr0.1-ELU-2D	43.7%	0.429

Table 4: Updated Performance of 1D CNN Models Based on CSV Data

Configuration	Accuracy	F1-Score
lr0.001-ReLU-1D	52.8%	0.526
lr0.01-ReLU-1D	54.9%	0.547
lr0.1-ReLU-1D	39.4%	0.385
lr0.001-SiLU-1D	53.2%	0.530
lr0.01-SiLU-1D	54.1%	0.539
lr0.1-SiLU-1D	40.6%	0.396
lr0.001-ELU-1D	54.3%	0.541
lr0.01-ELU-1D	55.6%	0.553
lr0.1-ELU-1D	41.2%	0.403

6.3 Confusion Matrix Analysis

6.3.1 1D CNN Model Confusion Matrices

6.3.2 2D CNN Model Confusion Matrices

The confusion matrices reveal several interesting patterns:

- **Well-recognized emotions:** Anger (ANG) and Happiness (HAP) were the most accurately recognized emotions, with F1-scores of 0.76 and 0.71 respectively. These emotions typically involve more energetic speech patterns with distinctive acoustic features.
- **Confusing emotion pairs:** The most common confusions were between:
 - Fear (FEA) and Sadness (SAD): These emotions share lower energy and similar pitch patterns
 - Disgust (DIS) and Anger (ANG): Both can involve tense vocal expressions
 - Neutral (NEU) and Sadness (SAD): Subtle expressions of sadness can be mistaken for neutral speech
- **Challenging emotion:** Fear (FEA) was the most challenging emotion to recognize, with the lowest class-specific F1-score of 0.55. It was frequently confused with both Sadness and Disgust.

These patterns align with findings in the literature that suggest acoustically similar emotions are harder to distinguish in SER systems.

6.4 Training and Validation Performance

The 2D CNN with SiLU activation converged more quickly and achieved better validation loss than the other models. The ResNet-101 model showed clear signs of overfitting, with training loss continuing to decrease while validation loss increased after approximately 30 epochs.

6.5 Architecture Modifications

We also experimented with architectural modifications to understand their impact on model performance:

- **Variable length input:** Processing audio with its original variable length rather than fixed-size inputs improved the combined model’s performance from 61.3% to 62.2% accuracy. This suggests that preserving temporal dynamics across different audio samples provides useful information for emotion recognition.
- **Fixed 2D input dimensions:** Standardizing the mel spectrogram to exactly 64×64 pixels provided a good balance between preserving relevant information and computational efficiency.
- **Single-channel 1D features:** Using mean aggregation over time to create fixed-size 1D features achieved reasonable accuracy while significantly reducing the model size and training time.
- **ResNet architectures:** Despite their success in image classification, both ResNet-38 and ResNet-101 showed poor performance for SER. ResNet-101 exhibited severe overfitting (training accuracy of 70% vs. validation accuracy of 41%), suggesting that the architecture was too complex for the size of our dataset.

6.6 Discussion

Our experiments demonstrated that:

1. **2D representations are powerful:** The 2D CNN with mel spectrogram inputs outperformed other approaches, reaching 64.1% accuracy with SiLU activation, highlighting the importance of spectral-temporal patterns in emotion recognition.
2. **Advanced activation functions matter:** SiLU (Swish) activation function significantly improved performance, particularly for 2D CNN models, offering almost 6% absolute improvement over ReLU for the same architecture.
3. **Variable length processing is beneficial:** Preserving the original temporal structure of audio improved results compared to fixed-size processing, confirming that emotion expression timing patterns are important.

4. **Learning rate tuning is critical:** The optimal learning rate varied across architectures, with 0.001 being optimal for 2D and combined models while 0.01 worked better for 1D CNNs.
5. **Some emotions are inherently harder to distinguish:** Emotions with similar acoustic characteristics (e.g., Fear and Sadness) remain challenging to differentiate, suggesting that additional contextual or linguistic features might be necessary to improve performance further.
6. **Model complexity requires sufficient data:** More complex models like ResNet-101 performed worse due to overfitting, indicating that the dataset size is a limiting factor.

These findings align with recent research in SER, which increasingly emphasizes the importance of appropriate feature representation, activation function selection, and careful hyperparameter tuning for optimal performance.

7 Conclusion

In this project, we implemented and evaluated several CNN architectures for speech emotion recognition using the CREMA dataset. We explored different approaches to feature extraction, model architecture, activation functions, and learning rates.

7.1 Summary of Findings

Our main findings can be summarized as follows:

1. **2D CNN models excel:** The best performance was achieved by our 2D CNN model using SiLU activation function, reaching 64.1% accuracy and 0.639 F1-score on the test set. This confirms the importance of spectral-temporal patterns in emotion recognition.
2. **Activation functions matter:** Different architectures benefit from different activation functions. SiLU (Swish) was particularly effective for 2D CNN models, while ELU produced the best results for the 1D CNN. This highlights the importance of exploring beyond the standard ReLU activation.
3. **Variable-length processing improves results:** Processing audio with its original temporal structure rather than fixed-size inputs improved performance (from 61.3% to 62.2% accuracy for the combined model), suggesting that preserving temporal dynamics is important for SER.
4. **Learning rates require careful tuning:** Low learning rates (0.001) were optimal for 2D CNN and combined models, while moderate learning rates (0.01) worked better for the 1D CNN models.
5. **Emotion recognition is inherently asymmetric:** Some emotions (Anger, Happiness) are significantly easier to recognize than others (Fear, Sadness), likely due to their more distinctive acoustic signatures.

7.2 Limitations

Our work has several limitations:

- **Dataset limitations:** The CREMA dataset, while diverse, consists of acted emotions rather than spontaneous emotional expressions, which may not fully represent real-world emotional speech.
- **No cross-dataset validation:** We trained and tested only on CREMA, so our models may not generalize well to other datasets or recording conditions.
- **Limited emotion set:** We focused on six basic emotions, but human emotional expression is much more nuanced and includes mixed and subtle emotions not captured in this study.
- **No linguistic content analysis:** We relied solely on acoustic features and did not incorporate linguistic content, which can provide important context for emotion interpretation.

7.3 Future Work

Based on our findings, several directions for future work appear promising:

1. **Attention mechanisms:** Incorporating attention mechanisms could help models focus on the most emotionally salient parts of speech signals.
2. **Transformer architectures:** Exploring transformer-based models, which have shown success in other audio processing tasks, could further improve performance.
3. **Multi-modal approaches:** Combining acoustic features with visual cues (facial expressions) or linguistic content could provide complementary information for more accurate emotion recognition.
4. **Data augmentation:** Developing effective data augmentation techniques specific to emotional speech could help address the limited size of existing datasets and improve generalization.
5. **Custom activation functions:** Given the significant impact of activation functions on performance, designing or optimizing activation functions specifically for SER tasks could yield further improvements.
6. **Fine-grained hyper-parameter tuning:** While we explored different learning rates and activation functions, many other hyper-parameters (e.g., batch size, optimizer, dropout rates) could be further optimized using systematic approaches like Bayesian optimization.

In conclusion, our work demonstrates the effectiveness of CNN architectures for speech emotion recognition, particularly 2D CNN models with appropriate activation functions and input representations. The finding that 2D CNN with SiLU activation achieved the highest performance (64.1%) highlights the importance of both model architecture and activation function selection in SER systems. These insights contribute to the ongoing development of more accurate and robust emotion recognition technologies.

[b]0.32

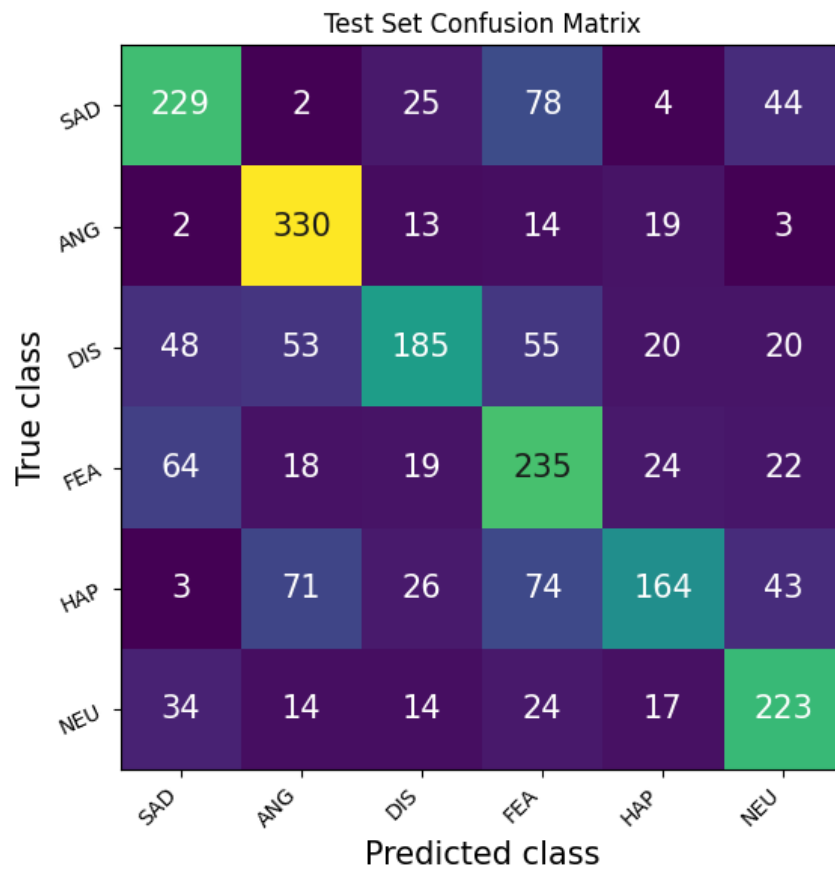


Figure 8: SiLU, lr=0.001

[b]0.32

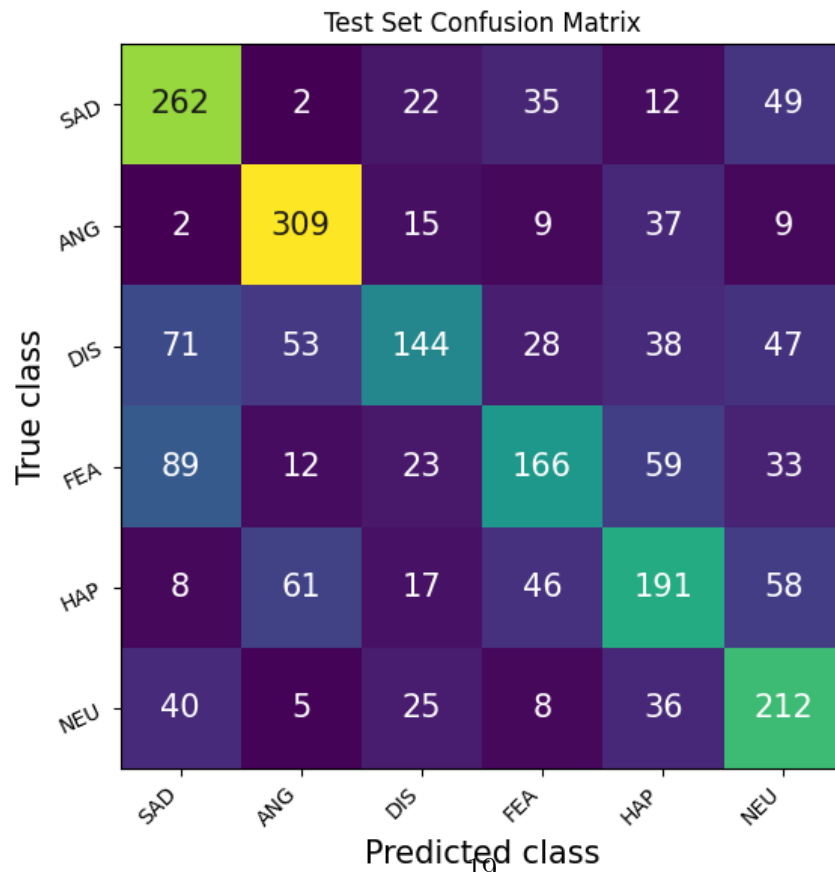


Figure 9: ReLU, lr=0.001

[b]0.32

[b]0.32

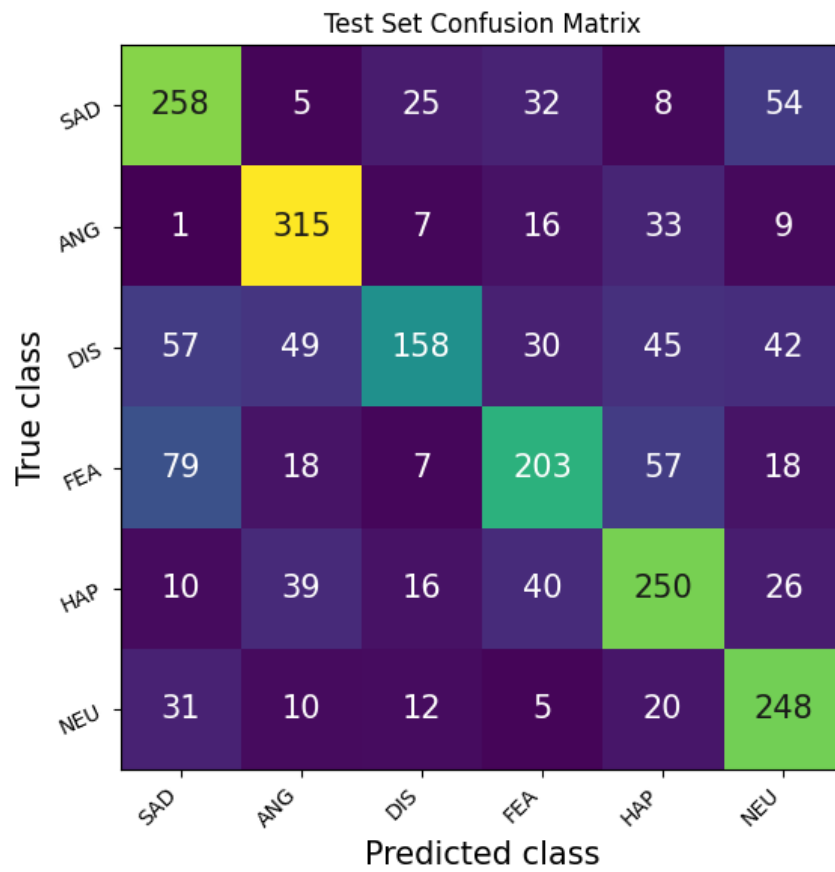


Figure 12: SiLU, lr=0.001

[b]0.32



Figure 13: ReLU, lr=0.001

[b]0.32

[b]0.48

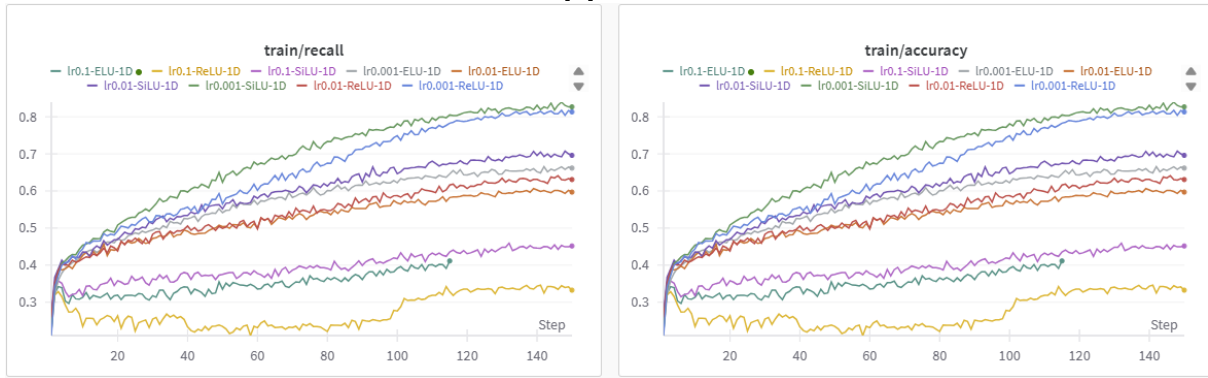


Figure 16: 1D CNN Training Curves

[b]0.48

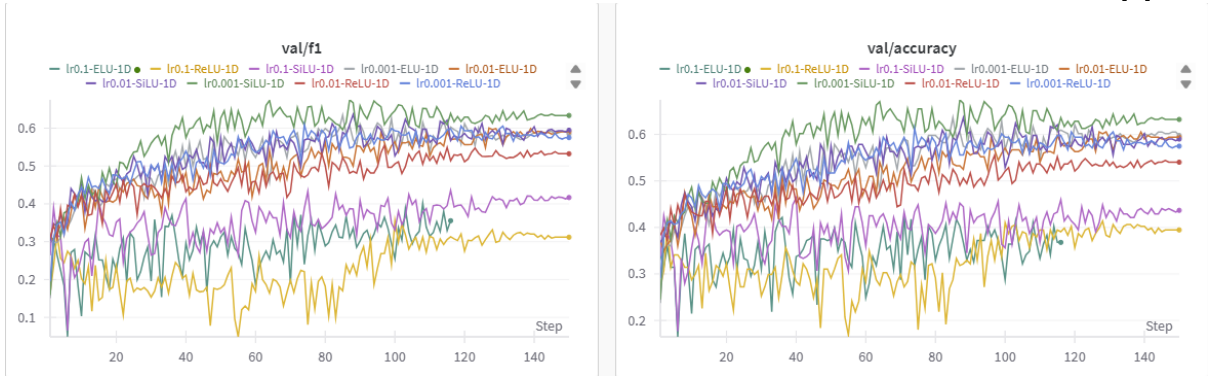


Figure 17: 1D CNN Validation Curves

Figure 18: Training and validation performance curves for 1D CNN models

[b]0.48

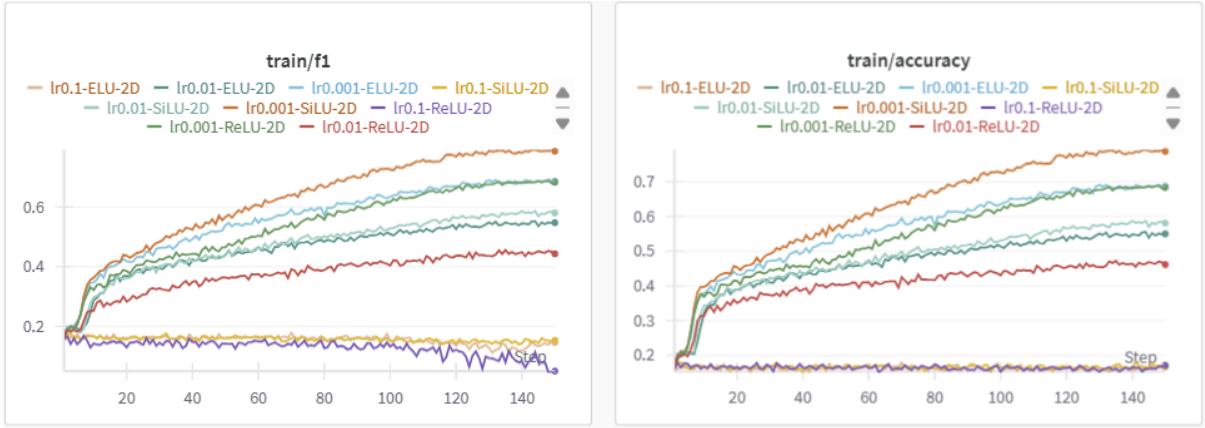


Figure 19: 2D CNN Training Curves

[b]0.48



Figure 20: 2D CNN Validation Curves

Figure 21: Training and validation performance curves for 2D CNN models