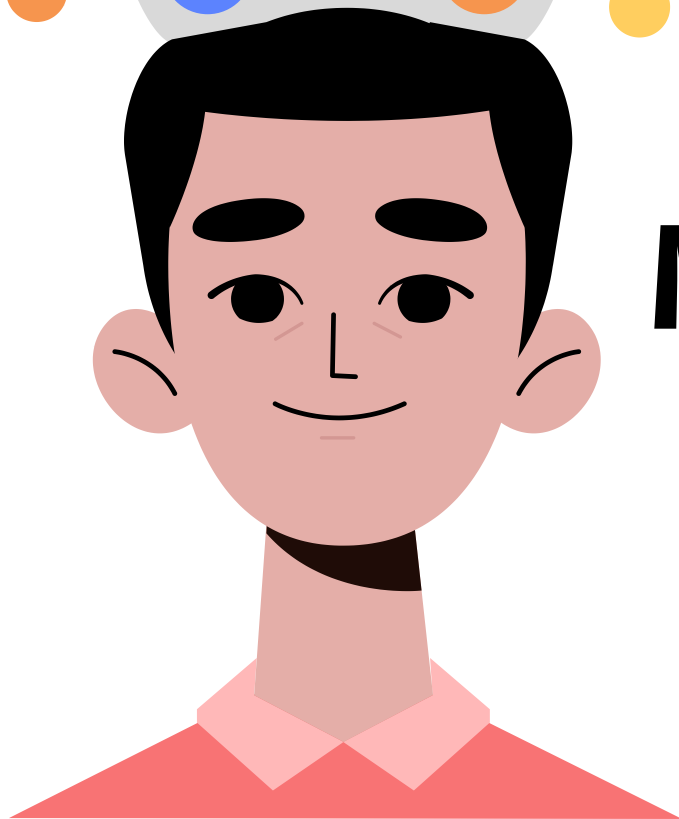# Classifying Emotions with Machine Learning

**Team 26**
Cooper Bosch, Conrad Hock,
Vincent Hock, Tomas Arevalo, Jake Pappo

# Problem Statement

to develop a sentiment analysis classifier capable of accurately identifying the emotional state of a speaker based on an audio clip of their speech
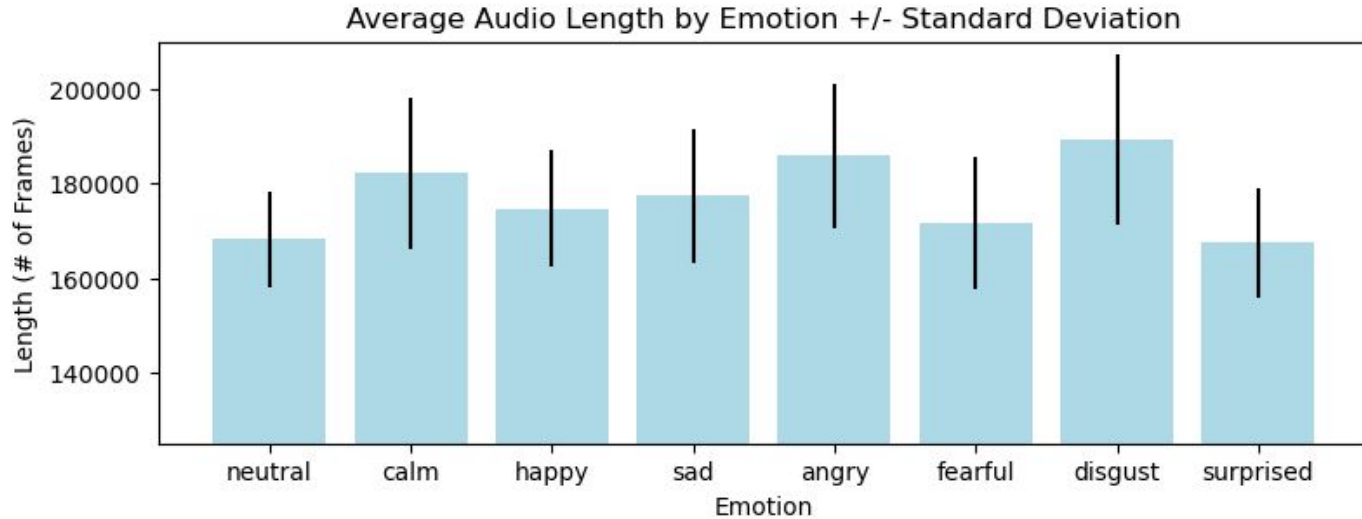
# Dataset Description

- Sourced from Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS)
- 24 voice actors tasked with recording phrases with a specific emotional mood
- About 1500 recordings total containing neutral voicings, and 7 different emotions:

  - Calm
  - Happy
  - Sad
  - Angry
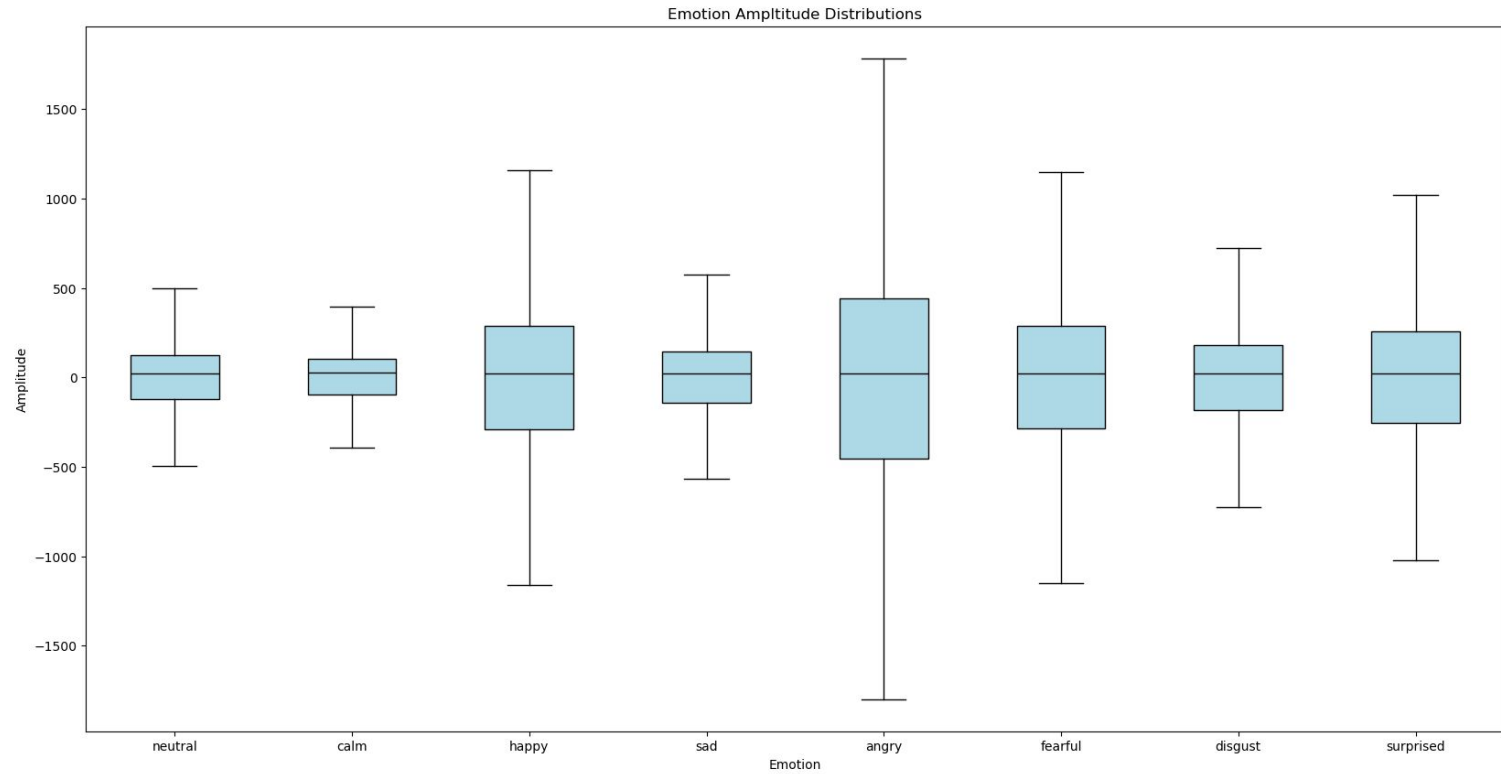  - Fearful
  - Disgusted
  - Surprised

# EDA

| | Emotion_Number | Emotion | Intensity | Statement_Number | Statement | Repetition | Actor | Gender_Number | Gender | Frame_Rate | Num_Frames | Data |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 4 | angry | 1 | 2 | Dogs | 1 | 16 | 0 | Female | 48000 | 187387 | [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... |
| 1 | 5 | fearful | 1 | 2 | Dogs | 2 | 16 | 0 | Female | 48000 | 171371 | [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... |
| 2 | 5 | fearful | 2 | 1 | Kids | 2 | 16 | 0 | Female | 48000 | 179379 | [0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ... |



Average Audio Length by Emotion +/- Standard Deviation
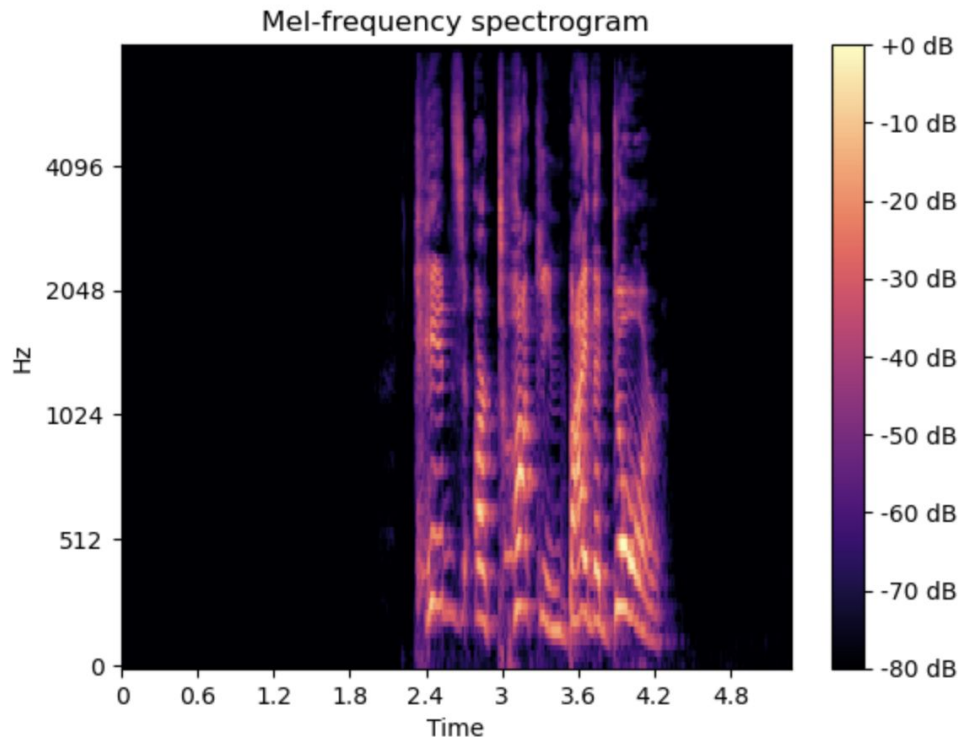
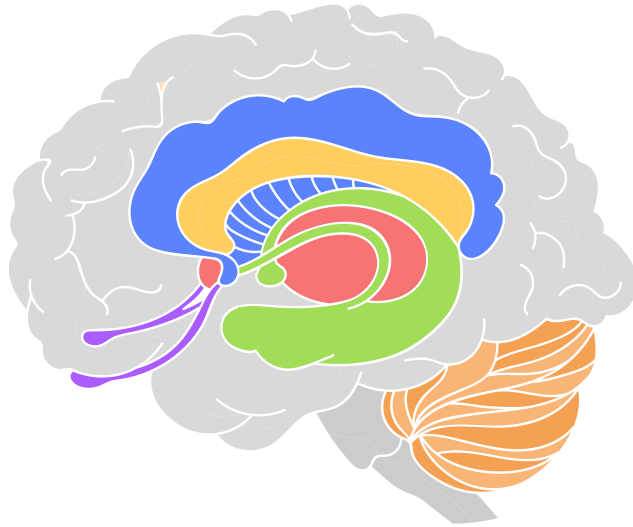# EDA



Emotion Ampltitude Distributions

# Mel Spectrogram

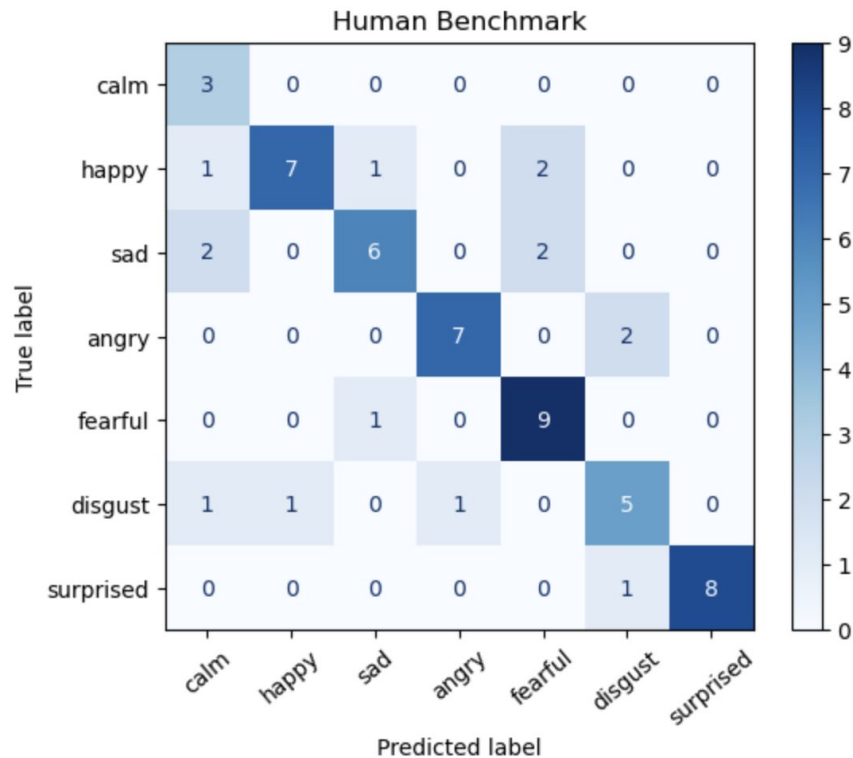Important preprocessing step with 2 main components:

1. Run FFT on small windows of signal to convert to pitch data

2. Transform pitches using mel scale, which is specifically designed to sound linear to humans



Mel-frequency spectrogram
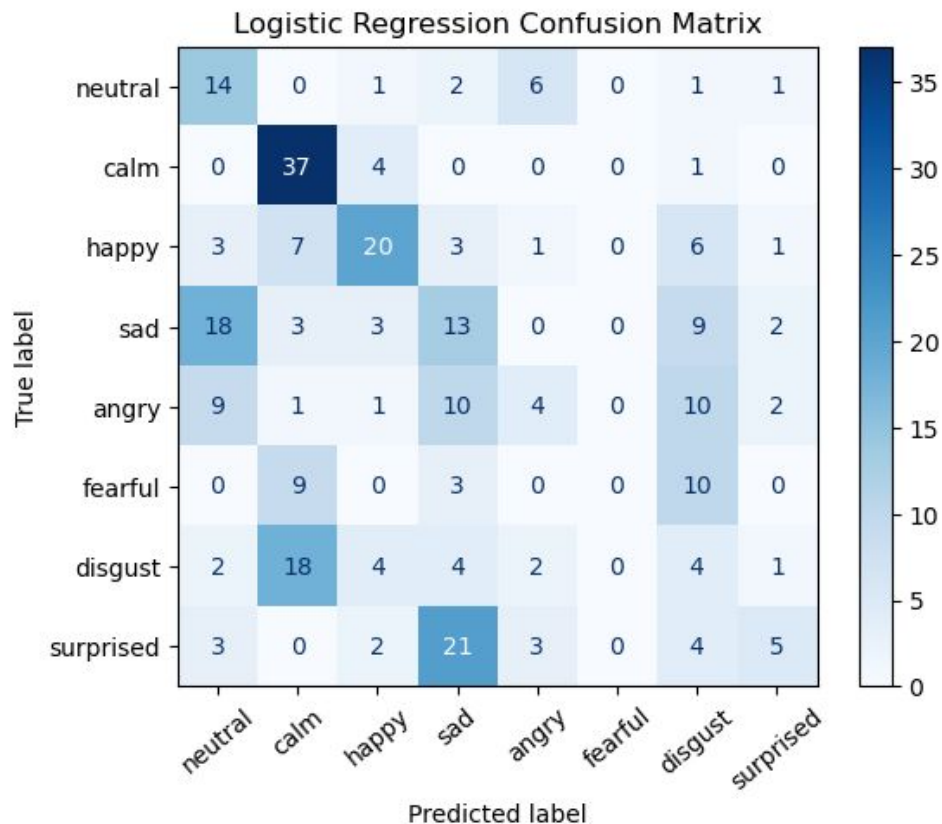
# Model Implementation and Analysis

# Human Benchmark



- 75% accuracy
- Worst at classifying sadness

# Logistic Regression - Baseline Model



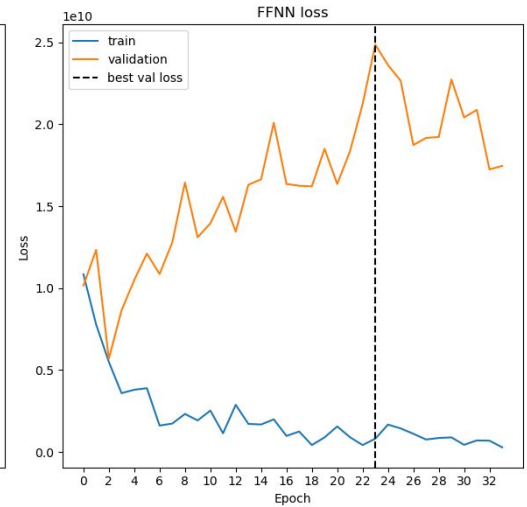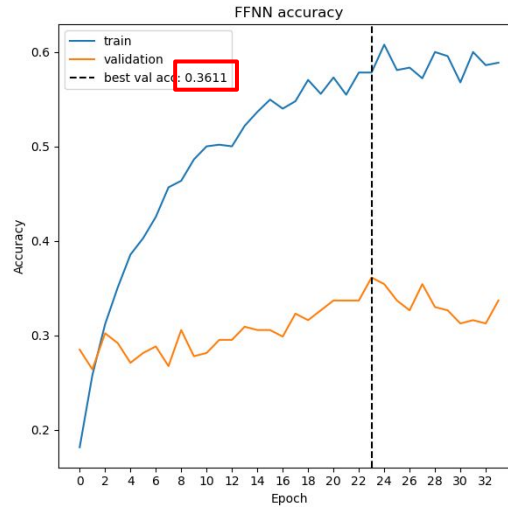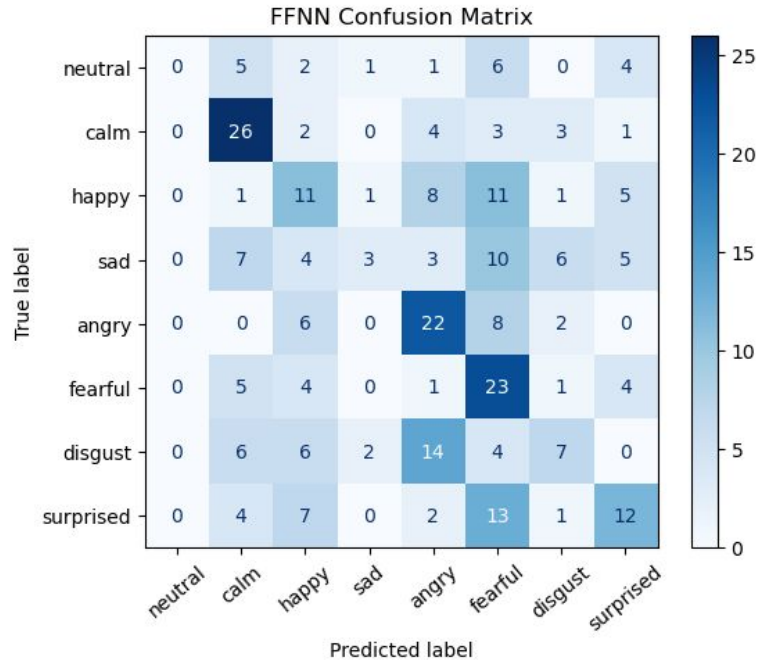Logistic Regression Confusion Matrix

- 34% Accuracy
- Takes as input:
  - Standard deviation of amplitudes
  - Length of clip
- Best at classifying calm emotion, which has lowest STD

# FFNN

**Parameters:** 31,935,008
**Train time:** ~4 min
**Epochs:** 34
**Val Acc:** 36.1%

Inputs

↓

Dense (500)

↓

Dropout (50%)

↓

Dense (500)

↓

Dropout (50%)

↓

Output

# FFNN



FFNN Confusion Matrix

FFNN accuracy

FFNN loss

# LSTM



## SIMPLE is better

**Parameters:** 3,445,320
**Train time:** ~13.5 min
**Epochs:** 100
**Val Acc:** 55.9%

# LSTM



Model accuracy

Model loss



Confusion Matrix

| True label \ Predicted label | neutral | calm | happy | sad | angry | fearful | disgust | surprised |
|---|---|---|---|---|---|---|---|---|
| neutral | 6 | 1 | 6 | 5 | 1 | 1 | 1 | 0 |
| calm | 5 | 24 | 2 | 9 | 0 | 2 | 2 | 0 |
| happy | 4 | 1 | 13 | 0 | 1 | 2 | 6 | 9 |
| sad | 2 | 5 | 3 | 15 | 0 | 4 | 0 | 5 |
| angry | 0 | 0 | 2 | 1 | 20 | 0 | 2 | 1 |
| fearful | 1 | 1 | 8 | 9 | 1 | 18 | 5 | 5 |
| disgust | 1 | 0 | 2 | 0 | 5 | 3 | 24 | 6 |
| surprised | 0 | 0 | 4 | 4 | 5 | 1 | 2 | 22 |

Ranking accuracy

1. Angry: 0.769
2. Disgust: 0.585
3. Surprised: 0.579
4. Calm: 0.545
5. Sad: 0.441
6. Happy: 0.361
7. Fearful: 0.375
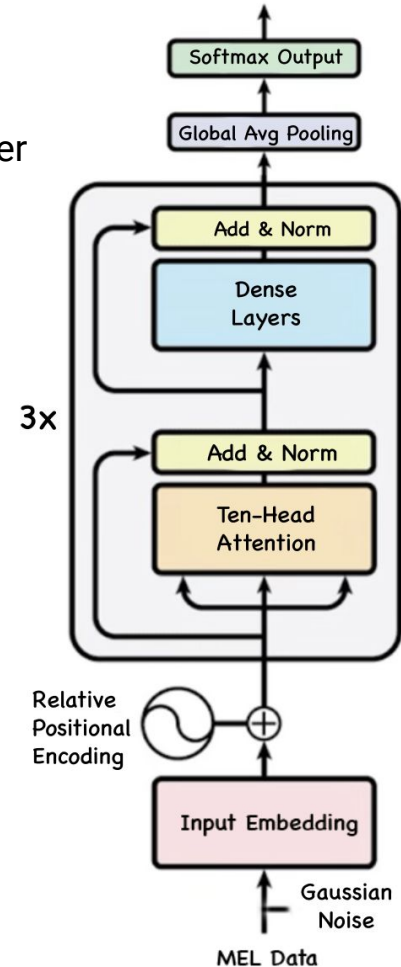8. Neutral: 0.286

Ranking precision:

1. Calm: 0.889
2. Angry: 0.606
3. Fearful: 0.581
4. Disgust: 0.632
5. Sad: 0.349
6. Neutral: 0.286
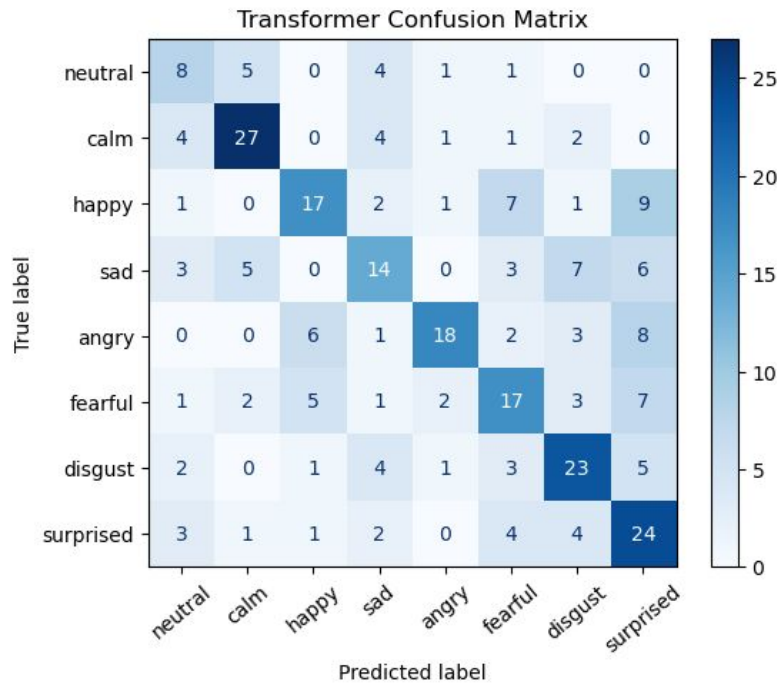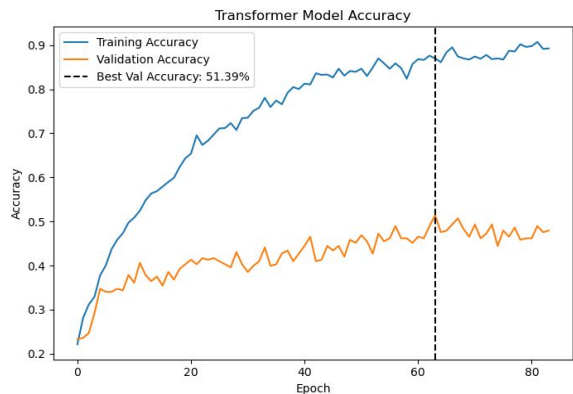7. Surprised: 0.458
8. Happy: 0.325

# Transformer

Specific Architecture: Bidirectional encoder



Let's take a look under the hood!

**Parameters:** 376,834
**Train time:** ~7 minutes
**Epochs:** 83
**Embed Dimen:** 50



Softmax Output

Global Avg Pooling

3x

Add & Norm

Dense Layers

Add & Norm

Ten-Head Attention

Relative Positional Encoding

Input Embedding

Gaussian Noise

MEL Data

# Transformer



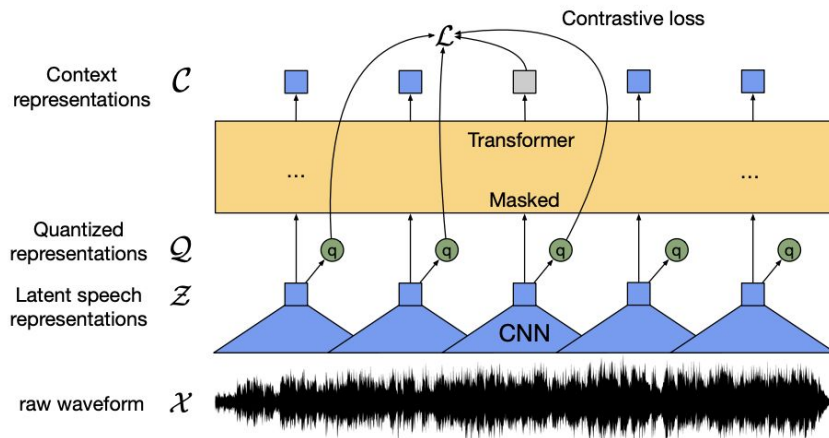Transformer Model Accuracy

Transformer Model Loss



Transformer Confusion Matrix

- **Validation Accuracy:** 51.4%
- "Calm" accuracy: 69.2%
- "Sad" accuracy: 36.8%
- "Surprised" precision: 40.7% (35 FPs)
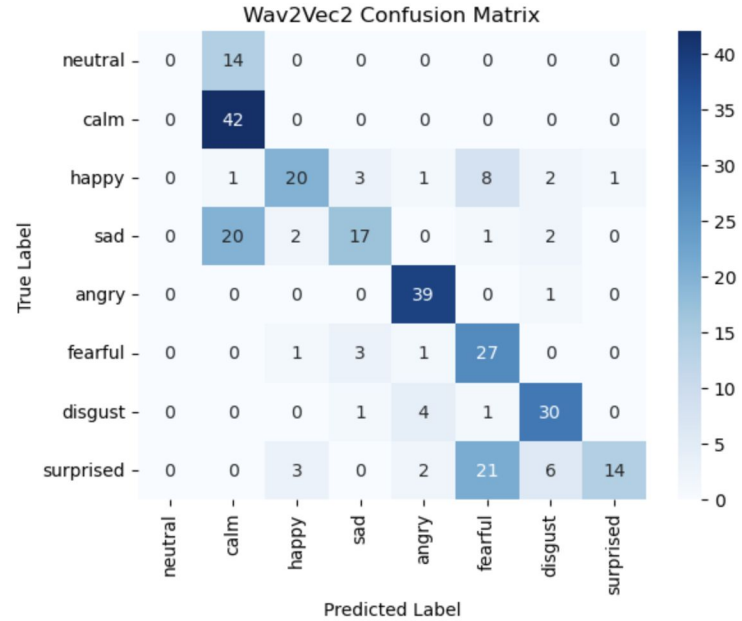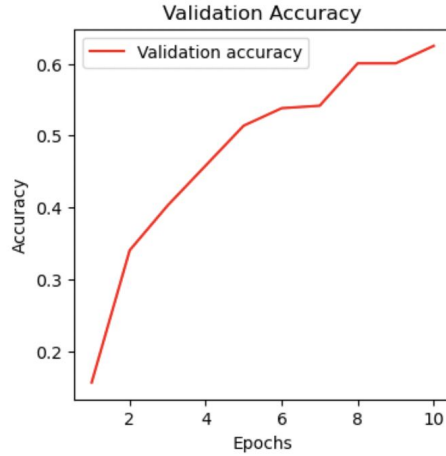- Relatively distributed error

# SOTA Model: Wav2Vec2

**Parameters:** 94,570,632
**Train time:** ~32 minutes
**Epochs:** 10
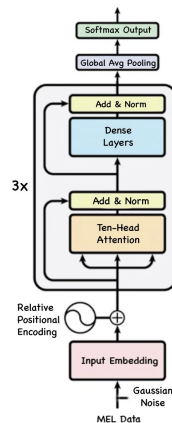**Val Acc:** 62.5%



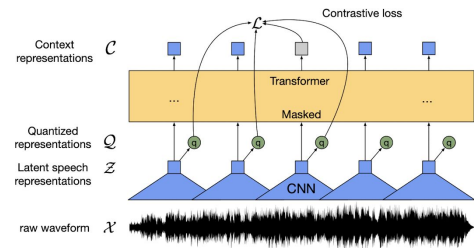Baevski et al., 2020

16

## SOTA Model: Wav2Vec2

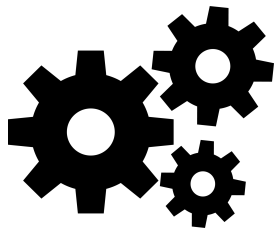# Conclusion



Fundamentally difficult problem

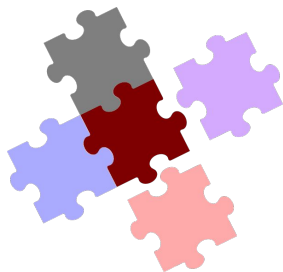Don't need a ton of parameters for decent accuracy with transformer
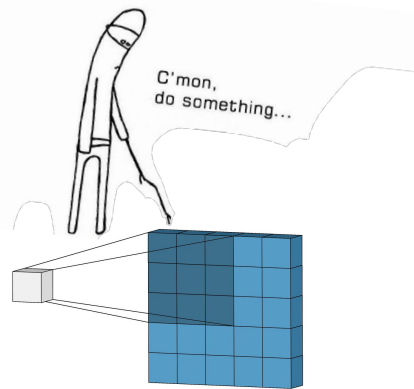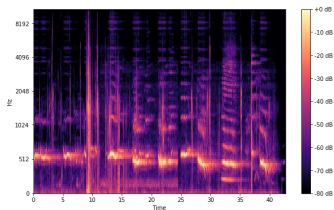
Transfer learning is usually the way to go

# Future Steps



Tweak architectures &
tune hyperparameters



Combining
emotions or
removing neutrals



C'mon,
do something...

Further exploration
of CNNs



Mel spectrogram
hyperparameters



Throw more compute
at Wave2Vec2



MORE DATA!

# Thanks!