# Transfer Learning for Emotion Recognition

Jessica Huynh, Amith Ananthram, Kailash Karthik
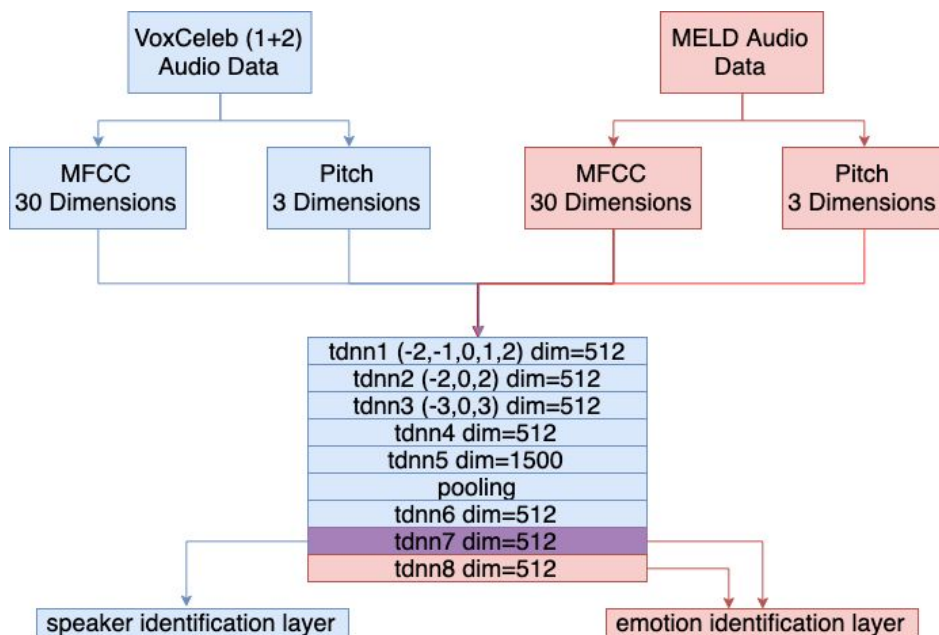jyh2127, aa4461, ks3740

# Problem Description

- **problem:** given an utterance, classify its emotion into one of {happiness, sadness, fear/surprise, anger/disgust, neutral}

- **approach:** transfer learning from larger speaker identification corpus to generate MFCC + pitch based speech embeddings; combine with BERT-based text embeddings; neural classification, LDA/PLDA

- **training:** VoxCeleb1 and 2, MELD
- **test:** IEMOCAP

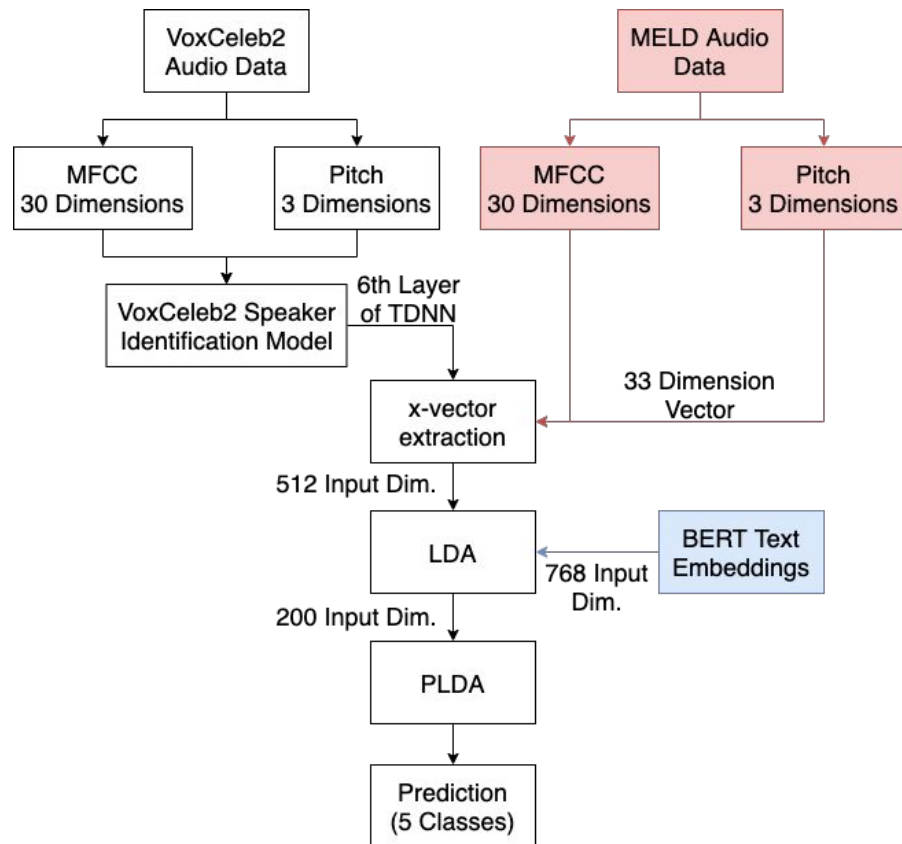| Mapped Emotion | IEMOCAP | MELD |
|---|---|---|
| Happiness | Happiness Excitement | Joy |
| Sadness | Sadness | Sadness |
| Fear/Surprise | Fear Surprise | Fear Surprise |
| Anger/Disgust | Anger Disgust Frustration | Anger Disgust |
| Neutral | Neutral | Neutral |

**Table 1**. Emotion Label Mapping

# Approach

**Network Training**

**LDA/PLDA Training**

VoxCeleb (1+2) Audio Data

MELD Audio Data

MFCC 30 Dimensions

Pitch 3 Dimensions

MFCC 30 Dimensions

Pitch 3 Dimensions

tdnn1 (-2,-1,0,1,2) dim=512
tdnn2 (-2,0,2) dim=512
tdnn3 (-3,0,3) dim=512
tdnn4 dim=512
tdnn5 dim=1500
pooling
tdnn6 dim=512
tdnn7 dim=512
tdnn8 dim=512

speaker identification layer

emotion identification layer

VoxCeleb2 Audio Data

MELD Audio Data

MFCC 30 Dimensions

Pitch 3 Dimensions

MFCC 30 Dimensions

Pitch 3 Dimensions

VoxCeleb2 Speaker Identification Model

6th Layer of TDNN

33 Dimension Vector

x-vector extraction

512 Input Dim.

LDA

BERT Text Embeddings

768 Input Dim.

200 Input Dim.

PLDA

Prediction (5 Classes)

# Neural Network Results



Confusion matrix

# LDA/PLDA Results

| LDA Input | EER |
|---|---|
| MELD Speech xvectors | 47.21% |
| BERT Text Embeddings | 46.23% |
| Speech & Text Embeddings | 43.05% |

**Table 3**. EER On All IEMOCAP

| Test Sess./LDA Input | Speech | Text | Both |
|---|---|---|---|
| Session 1 | 40.91% | 41.48% | 35.49% |
| Session 2 | 41.04% | 40.95% | 34.78% |
| Session 3 | 39.86% | 41.91% | 34.87% |
| Session 4 | 39.21% | 40.82% | 33.71% |
| Session 5 | 40.11% | 41.46% | 34.70% |
| Weighted Average | 40.19% | 41.33% | 34.69% |

**Table 4**. EER On Fifths Of IEMOCAP

# Challenges

1. cross-domain evaluation (train on VOX/MELD, evaluate on IEMOCAP)
2. partial signal: human annotators used visual cues + audio cues
3. imbalanced corpora: not all emotions are represented equally
4. context independent combination of text + speech vectors
   a. we generate context dependent dependings of text and speech separately and combine them

# Future Work

1. neural network
   a. reference model
      i. train speaker identification model with additional layers
      ii. train speaker identification model with silence frames
   b. domain adaptation
      i. include portions of IEMOCAP in re-training (in-progress)
   c. emotion conversion
      i. using the emotion detector, train an auto-encoder (in-progress)
2. target emotions (experiment with other clusterings)
3. model entire conversations for latent emotional state
4. LDA/PLDA
   a. context dependent combinations of speech + text embeddings
   b. Use embeddings from 8th layer of neural network for LDA/PLDA (in-progress)