# A Study on Elderly Speech Emotion Recognition for Low-Resource Languages

Lo Shun Hang Lincoln, 1155176490
Wong Ka Long Myron

Supervisor: Prof. Wu, Xixin

The Chinese University of Hong Kong
SEEM4999 Final Year Project II
Final Report

April 2025

# 1 Abstract

# 2 Introduction

## 2.1 Background

Hong Kong is a city facing major population aging problems[8]. Alongside the aging population, Hong Kong does not have enough infrastructure to support the elderly care industry, which is critical to taking good care of the elderly citizens.

Artificial Intelligence can help with a lot of things, including elderly care. Studies have found that for the elderly, loneliness is a significant issue that can lead to various cognitive or mental problems[13]. To maintain the quality of life for the elderly, we can create Artificial Intelligence systems that can analyze whether an older adult is mentally unstable and provide them with the right support.

## 2.2 Problems

In the modern world where Artificial Intelligence is blooming, human-computer interaction is becoming more and more important[5]. As a result, the popularity of speech emotion recognition is also rising, unfortunately, the majority of studies tend to focus on young adults that speak English[6]. Despite the popularity in speech emotion recognition, there is little research on elderly speech

emotion recognition, making it a hard problem to solve[14]. It is even harder for low-resource languages (e.g., Cantonese) to train a powerful enough elderly speech emotion recognition model that can support intensive use in real life with the lack of data.

To train powerful models under such harsh conditions, researchers have derived techniques that can fully utilize the small datasets to train a better model, these methods include automatic feature extraction and cross-domain learning[12, 4].

## 2.3 Objective

In this paper, two different approaches to low-resource training will be explored: cross-domain learning and concatenating low-level features. For cross-domain learning, we study the possibility of transferring emotion recognition abilities between different age groups (Adult - Elderly). By utilizing pre-trained speech models (e.g., XLSR-Wav2Vec 2.0), a representation of the speech can be acquired, a classifier will then be used to get predictions from the representations. Comparing the results, a conclusion of whether cross-domain learning can be beneficial for elderly speech emotion recognition can then be made. Text-based features will also be added to test the usefulness of text-based features in speech emotion recognition. Texts will be transcribed with a Cantonese fine-tuned Whisper model[1], and then it will be fed into a Cantonese fine-tuned fastText model[18] to extract its text representations.

By testing these approaches, we hope to answer the following questions:

- Whether cross-domain learning is viable in elderly speech emotion recognition

- Whether text-based features contribute to the training of elderly speech emotion recognition models

# 3 Related Works

## 3.1 Feature Concatenation

Feature concatenation involves combining two or more distinct feature vectors each representing different aspects of the same speech signal into a single, unified feature vector at the frame level. The goal is to leverage complementary information captured by these features to enhance the robustness and accuracy of the system.

Emotional content in speech is conveyed through variations in pitch, energy, speaking rate, and spectral characteristics all of which can be captured by different feature types. For example:

- **Mel-frequency cepstral coefficients (MFCCs)**: Widely used in SER, MFCCs capture the spectral envelope of speech, which reflects timbre and vocal tract characteristics influenced by emotion.

2

- **Gamma-tone filterbank-based features (e.g., GTF-CC)**: These model the human auditory system more closely and are sensitive to pitch and harmonic structures, which are critical for detecting prosodic cues like intonation patterns tied to emotions.

- **Prosodic features**: Features like fundamental frequency (F0), energy, and duration can directly indicate emotional arousal or valence.

** feature concatenation (1)

## 3.2 Cross-Domain Learning

Cross-domain learning refers to the process of training a model to recognize emotions from speech across different datasets (or corpora) where the feature distributions vary due to factors like speaker characteristics, recording conditions, languages, or corpora-specific attributes. This is a critical challenge in real-world SER applications because models trained on one dataset often perform poorly when applied to another due to these domain differences, a phenomenon known as domain divergence or domain shift.

### 3.2.1 Domain Divergence Problem

In SER, datasets (e.g., IEMOCAP, MSP-Improv, SAVEE, Emo-DB) differ in aspects such as language (English vs. German), speaker demographics (gender, number of speakers), recording type (acted vs. hybrid), and emotional annotations. These differences lead to variations in the feature distributions of the source (training) and target (testing) datasets. Traditional SER models, trained and tested on the same dataset, struggle to generalize to unseen datasets because they inadvertently learn domain-specific, non-affective information (e.g., speaker identity, recording environment) alongside emotional cues.

### 3.2.2 Domain Adversarial Neural Networks (DANN)

Objective Function:

1. Feature Extractor: A deep convolutional neural network (CNN) combined with a bidirectional LSTM (BLSTM) extracts features from speech spectrograms.

2. Emotion Classifier: Predicts emotions (e.g., arousal and valence) using these features.

3. Domain Classifier: Attempts to identify domain-specific attributes (e.g., corpus, language, gender).

4. Gradient Reversal Layer (GRL): Inserted between the feature extractor and domain classifier, this layer reverses the gradient during backpropagation. This forces the feature extractor to maximize the domain classifier' s loss (making domain features indistinguishable) while minimizing

the emotion classifier' s loss, thus learning domain-invariant emotional features.

### 3.2.3   Center Loss Integration

Objective Function:

1. Softmax Loss: Separates different emotion classes by finding decision boundaries

2. Center Loss: Minimizes the Euclidean distance between feature representations and their corresponding emotion class centers, reducing intra-class variation (e.g., ensuring features of "happy" from different datasets cluster together)

 ** cross domain learning (3)

## 3.3   Low-Level Features

Low-level features in SER are typically time- or frequency-domain descriptors of the audio signal that reflect acoustic properties such as pitch, energy, and spectral characteristics. These features are often hand crafted or extracted using well established signal processing techniques, making them "knowledge-based" as they rely on prior understanding of how sound relates to human perception. In the paper, MFCC is highlighted as a key low-level feature due to its widespread use in speech processing and its ability to mimic the auditory characteristics of the human ear.

1. MFCC (Mel-frequency Cepstral Coefficients): MFCCs are derived from the audio signal by applying a series of transformations, including the Fourier transform, Mel-scale filtering, and discrete cosine transform. They represent the short-term power spectrum of sound on a perceptually relevant scale, capturing how humans perceive frequencies.

2. Alternative: ComParE feature set and eGeMAPs, which include additional low-level descriptors such as pitch, jitter, shimmer, and formants. However, MFCC is chosen as the target low-level feature for its complementary nature to high-level features.

 ** low level features (1)

## 3.4   Multi-Label Classification

## 3.5   Semi-Supervised Models

## 3.6   Speech Emotion Recognition

Fill in here...

# 4 Benchmark and Dataset

## 4.1 Benchmark

Our main objective is to test the usefulness of various techniques in low-resource learning, and the benchmarking will be done by comparing each model's:

- Accuracy

- F1 Score

For the classification of elderly and adult, the elderly are $\geq 55$ years old, and adults are $\leq 54$ years old. A vanilla version of the data will be used to train a baseline model, which means that a baseline model will be trained only with data from the elderly, and others will adopt the techniques that we have mentioned. Then, we will compare the accuracy and the F1 score of each model to the baseline model and see if there is a positive impact on the performance of a model after applying the techniques.

## 4.2 Dataset

Our dataset comprises 1 publicly available dataset (YueMotion[7]) and 1 dataset that we have gathered ourselves. YueMotion is a Cantonese speech emotion recognition database that has 18 speakers (11 adults, 7 elderly) with 6 emotion labels, consisting of 420 records for the elderly and 660 for adults. Our dataset is composed of audios clipped from TV shows, podcasts, documentaries, etc., it has 176 records and 5 emotion labels. As for this study, we are only using 5 labels for simplicity: anger, happy, sad, neutral, fear.

An important thing to note is that the sentimental and linguistic structure of our dataset and YueMotion is vastly different. The linguistic content in YueMotion is repetitive. A sentence like " 你宜家幾多歲" will be repeated with different emotions and with different speakers. For our dataset, the sentence in the audio will be random and will not be repeated. The emotions of the audio will also be related to the speech being said.

# 5 Model Explanations

Our model consists of various parts:

- Feature Extraction Models

  - XLSR-Wav2Vec 2.0 (Speech Representations)
  - Whisper (Transcription)
  - fastText (Text Representations)

- Classification Model

The details of the models will be explained in the following subsections.

## 5.1 Wav2Vec 2.0

Wav2Vec 2.0 is a powerful semi-supervised model developed by Facebook AI that outputs representations from speech audio[3]. The architecture of the model can be seen in the following figure:
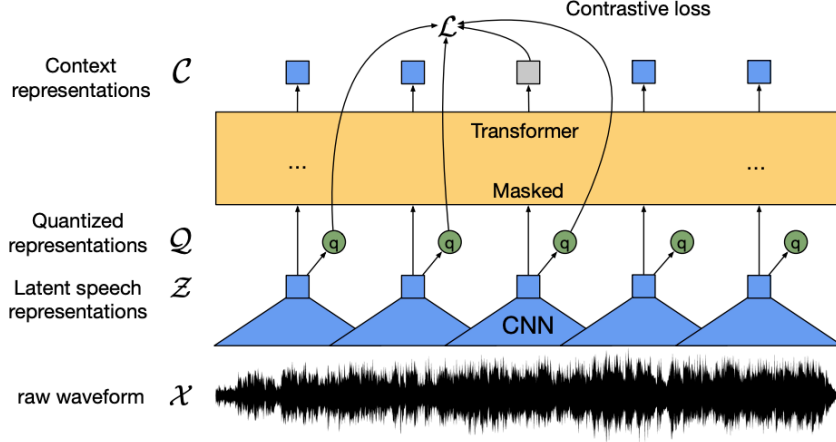


Figure 1: Architecture of Wav2Vec 2.0[3]

The model consists of a few parts: a feature encoder, a transformer, and a quantization module. A raw audio $X$ is inputted into the feature encoder. Then, a latent speech representation is fed into both the transformer and the quantization model. The transformer tries to capture information in the speech in the form of representation vectors. The quantization module discretizes the latent speech representations to represent the targets in the self-supervised objective[3]. After training on relevant speech data, the model can then output contextualized representations of an audio that contains various low-level features that are vital to speech emotion recognition tasks. A representation $L$ is formed by a true quantized latent speech representation by solving contrastive task $L_m$ and a codebook diversity loss $L_d$.

$$L = L_m + \alpha L_m$$

Where $\alpha$ is a hyperparameter that is tuned.

The contrastive loss is calculated as:

$$L_m = -\log \frac{\exp(sim(c_t, q_t)/k)}{\sum_{\bar{q} \sim Q_t} \exp(sim(c_t, \tilde{q})/k)}$$

and the cosine similarity as $sim(a, b) = a^T b / \|a\| \|b\|$.

Where $c_t$ is a output from the transformer centered over masked time step t, and $Q_t$ is a set of $K + 1$ quantized representations $\tilde{q} \in Q_t$, which includes

a true quantized representation and $K$ other distraction representations. The model needs to identify $q_t$ for a masked time step within the set $Q_t$

The Diversity loss is calculated as:

$$L_d = \frac{1}{GV} \sum_{g=1}^{G} -H(\bar{p}g) = \frac{1}{GV} \sum_{g=1}^{G} \sum_{v=1}^{V} \bar{p}_{g,v} \log \bar{p}_{g,v}$$

Where there are V entries in each of the G codebooks. And it wants to maximize the entropy of averaged softmax distribution over the codebook entries for each codebook $\bar{p}_g$.

In our previous study, Wav2Vec 2.0 was used, but as Wav2Vec 2.0 is trained on monolingual speech corpora (English), it is not suitable for speech emotion detection for other languages. Thus, we have chosen XLSR-Wav2Vec 2.0, a variation of Wav2Vec 2.0, as our model. It is trained on multilingual speech corpora, and its ability to learn multilingual representations will be crucial to our task.

For convenience, we have sourced a pre-trained model of XLSR-Wav2Vec 2.0 that is fine-tuned on Cantonese[10].

## 5.2 Whisper

Whisper is an OpenAI developed multilingual speech model that is capable of many tasks, including translation and transcription. Whisper's approach to training a multilingual speech recognition model can be an example of our cross-domain learning techniques, as it is also trained on multilingual speech corpora. The approach of Whisper can be seen in the following figure:
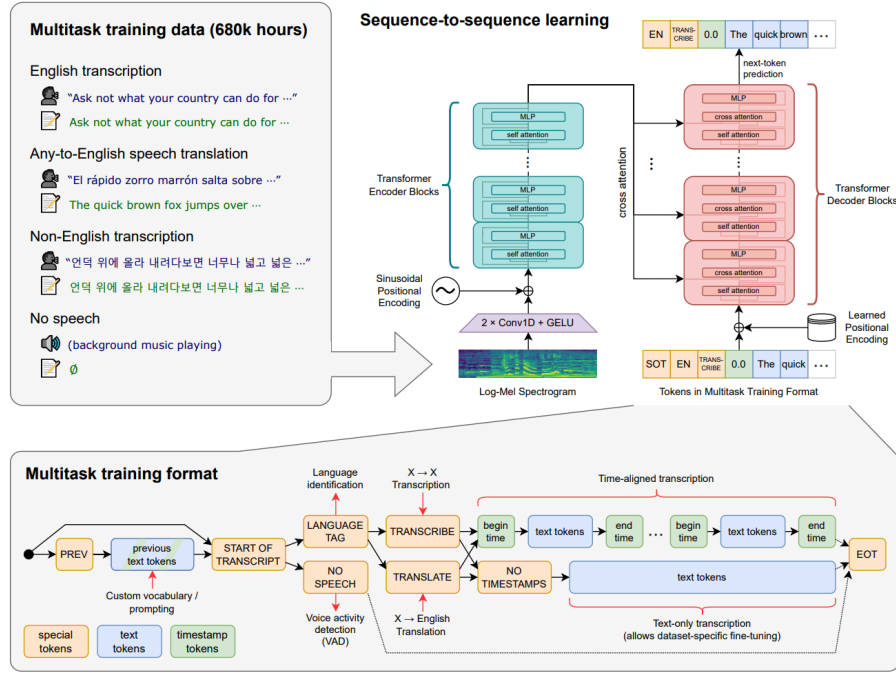


Figure 2: Overview of Whisper's approach[17]

It is a model adopting an encoder-decoder Transformer architecture, 2 convolutional layers with a filter and a GELU activation function. For the tokenizer, Whisper refits multilingual vocabulary for multilingual models to maintain accuracy for all languages.

For the training of Whisper, although it considers the history of the text being predicted, the tokens of the previous text will be masked, allowing it to resolve more ambiguous audio. Then, the model identifies if speech is present in the audio. If speech is present, transcription and translation are trained, where transcription is just time-aligned transcription in different languages.

The difference between transcription and translation in Whisper is that translation is limited to X to English instead of X to X. Transcription can be done in most languages and can be further fine-tuned for a chosen set of languages.

For convenience, we have also used a Whisper model fine-tuned on Cantonese that is available online[1].

## 5.3 fastText

aslkdjsaljd

## 5.4 Classifier

The classifier model contains 3 layers: 2 linear activated layers, a ReLU layer, and a softmax layer that outputs the probability density for the emotions. During training, the model will be trained on the cross-entropy loss between the features and the labels.

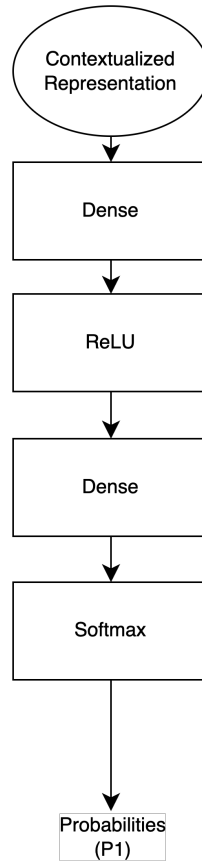The architecture of the classifier can be seen in the following figure:

Figure 3: Architecture of classifier

# 6 Methodology

As mentioned in the introduction of this work, we will adopt two different approaches to see if they can contribute to improving the accuracy of a model. For each domain group, it will be labeled as:

$$d_{group}$$

Each dataset $X$ will be denoted as:

$$X_{text?}^{d_{group}}$$

For example, a dataset with cross-age training but no text features and a dataset without cross-age training but with text features will be denoted as:

$$X^{d_{adult}, d_{elderly}}, X_{text}^{d_{elderly}}$$

The details of each low-resource training technique will be explained in the following subsections.

## 6.1 Cross-Domain Learning

Cross-domain learning is a significant problem in speech emotion recognition because, to build multilingual speech recognition models, we must first solve the problem of domain mismatch between different languages. In this study, we have gathered a dataset that contains speakers of different sexes, and different ages. These speakers possess different acoustic features, such as the intensity and pitch of their speech[19].

The logic of cross-domain learning is that we believe a source dataset $X_s$ has some hidden relationship with a target dataset $X_t$, and this relationship can be learned through learning $X_s$[9]. That is, we expect that there is some intersection between the domain of $X_s$ and $X_t$, and by learning $X_s$, we can have a good enough approximation of the feature in $X_t$. For example, there is a significant relationship between Mandarin and Cantonese, so we train our transcribing model on Mandarin, hoping that it learns this relationship and performs well on Cantonese speech emotion recognition as well. A figure of the relation can be seen in the following figure:

In the context of our study, because we lack high-quality speech data from the elderly. To train a robust speech emotion recognition model that can accurately identify emotions from speech, we would need to adopt cross-domain learning. That is, we try to transfer the ability to recognize speech emotions in young adults to the elderly.

This will be done by only adding speech audio of young adults to the training data. Then, we will conduct the training as usual, that is, after extracting the representation vector, we train the classifier and try to predict the test set that consists of only elderly speakers. More details on the exact approach will be mentioned in the training approach section.
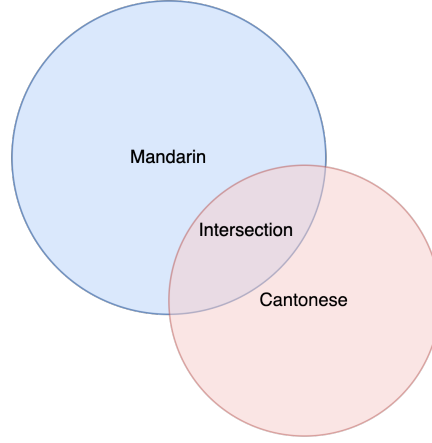
Figure 4: Relation between domains

## 6.2  Text-Based Features

Text-based features include many categories, text, word embeddings, and others. In the context of this study, text-based features are contextualized word embeddings that are learned by pre-trained models.

Word embeddings, also known as word vectors, are represented as a series of numbers, hence the name word vectors. These vectors allow models to understand the meaning of the words in a way that humans do not. Similar words should ideally be closer to each other. Meaning:

$$\|\overrightarrow{Food} - \overrightarrow{Ramen}\|_2 < \|\overrightarrow{Food} - \overrightarrow{Door}\|_2$$

should hold[2].

Word vectors have been used in recent studies in many fields: social science[2], sentiment analysis[11], and speech emotion recognition[16]. The number of studies embracing the use of word vectors can be an indicator of how powerful word vectors are in capturing low-level features of words.

In this study, word vectors will be extracted from words by a variation of fastText that is pre-trained on Cantonese data. To ensure the accuracy and the robustness of the vectors, we will be using a 300-dimension embedding. The word vectors are supplementary features to the acoustic features (speech representation) in hopes of boosting the accuracy of the model that is lacking training resources. More details about the training method will be mentioned in the training section below.

## 6.3  Training Approach

In order to extract low-level features for the classifier to learn, pre-trained models that can extract contextualized features will be utilized. For feature extrac-

tion, XLSR-Wav2Vec 2.0, Whisper, and fastText fine-tuned on Cantonese will be used.

Raw audio will be passed to the model, then be transformed into a representation of the audio with a dimension of 100.

In the case that both speech representations $s_i$ and text vectors $t_i$ will be used, Whisper will be used to transcribe the audio into text. Then, the text will be fed to fastText for feature extraction. Then the speech representation and the text embedding will be concatenated together as a new vector $r_i$ in the form:

$$r_i = [x_i^T, t_i^T]^T$$

the resulting vector from this operation will have 400 dimensions.

The classifier model will be trained by trying to solve the following minimization problem (Cross Entropy Loss):

$$\min \frac{\sum_{n=1}^{N} l_n}{N}$$

where

$$l_n = -\sum_{c=1}^{C} w_c \log \frac{\exp(x_{n,c})}{\sum_{i=1}^{C} \exp(x_{n,i})} y_{n,c}$$

and $x$ is the input vector, $y$ is the target, $w$ is the weight, $C$ is the number of classes, and $N$ is the number of minibatches of data. Where $x$ is a vector with 100 or 400 dimensions and $y$ is an integer.

The vectors pass through two linearly activated layers and a ReLU layer, then will finally be passed to a softmax layer to generate the probability density of the emotions.

The training hyperparameters are:

- Learning Rate: 5e-5

- Number of Epochs: 10

- Weight Decay: 0.01

The hyperparameters are carefully tested and selected to prevent overfitting of the classifier. It is also important to note that the feature extraction models will not be trained or modified in the course of the training.

# 7 Results and Discussions

## 7.1 Training Results

We have in total 6 variations of the dataset. The training results of each dataset can be seen in the following table:

| Dataset | Average Accuracy | Average F1-Score |
|---|---|---|
| $O^{d_{elderly}}$ (Own Data Only) | 32.38% | 0.2973 |
| $O^{d_{elderly}}_{text}$ | 35.20% | 0.3342 |
| $X^{d_{elderly}}$ | 37.53% | 0.3750 |
| $X^{d_{elderly}}_{text}$ | 34.30% | 0.3404 |
| $X^{d_{elderly},d_{adult}}$ | 49.39% | 0.4933 |
| $X^{d_{elderly},d_{adult}}_{text}$ | 48.16% | 0.4820 |

Table 1: Training Result of Classification Model

The statistics of each dataset is an average of a 3-fold cross-validation to ensure consistency and unbiasedness of the results.

## 7.2 Discussions

There are a few interesting things that were observed in the process of training the model:

- Relationship between accuracy and dataset size

- For 2 datasets, text-based features caused a drop in accuracy

- The Variance of accuracy between different folds of training increases when cross-domain learning is used

These are all interesting phenomena that occured in the training process of the classifier, and in the following sub-sections, we will try to explain why these phenomena occurred.

## 7.3 Relationship of Accuracy and Dataset Size

One interesting thing that happened during training is that we thought the accuracy of the model would grow somewhat linearly. After doing some exponential regression and plotting some graphs, we discovered the following:
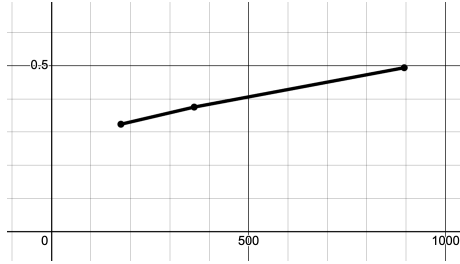
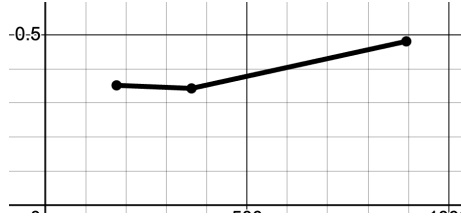Figure 5: Accuracy Graph for Acoustic Feature Only Datasets



Figure 6: Accuracy Graph for Text-Based Feature Databases

the black line in each graph is the lines connecting the three data points, that is the accuracy of $O^{d_{elderly}}, X^{d_{elderly}}, X^{d_{elderly},d_{adult}}$, and $O^{d_{elderly}}_{text}, X^{d_{elderly}}_{text}, X^{d_{elderly},d_{adult}}_{text}$.

It can be seen that there is some form of linearity in the growth of accuracy in the first graph, while in the second graph, the relationship looks more non-linear. While the sample size of this observation is very small, we still think that this is an interesting phenomenon because this may be an important step to improving prediction accuracy in low-resource domains. There can be many factors contributing to the non-linearity of the accuracy growth.

One thing that comes to mind is the power of word embeddings. Word embeddings are powerful in the sense that it capture a lot of low-layer or hidden features that humans can't comprehend. Even with some transcription errors (0.0972 CER), it can still harness a great portion of the sentimental information available in the text. When more similar word embeddings are fed to the classifier, it starts to generalize and apply the low-level features to future predictions, thus resulting in the non-linear growth. This might be a reminder to use related deep-learning-based features, as their correlation might be greater than we have imagined.

## 7.4 Accuracy Drop for Text-Based Features

In both $X^{d_{elderly}}$ and $X^{d_{elderly},d_{adult}}$, the accuracy of the model drops after using text-based features in the training process of the model. This can be caused by a few reasons: 300 dimensional word embeddings may be too hard to learn in

such a low-resource setting, and the linguistic content of the speech might not match the sentiment of the speech.

For example, speakers in the YueMotion dataset will say "你宜家幾多歲" in fear, which, in general, will not be the case. The model may be distracted by these misleading audio in the dataset and learn unrelated features. In our database, the linguistic content of the speech often highly correlates with the sentiment of the audio, and that is maybe why in our own dataset, the text-based feature performs better.

Also, word embedding with higher dimensions is harder to learn. In low-resource settings, high-dimension word embeddings can encounter the problem of data sparsity, meaning they are not correctly representing a word[15] due to lack of data to accurately describe in higher dimensions. Cantonese, being a low-resource language just as we have mentioned many times in this study, is prone to this problem of data sparsity. The data used to train the embedding model might not be dense enough to correctly transform words into high-dimensional vectors. Hence, the accuracy drops when text-based features are used.

## 7.5

## 8 Future Work

## 9 References

## References

[1] alvanlii. whisper-small-cantonese. https://huggingface.co/alvanlii/whisper-small-cantonese. Accessed: 2025-04-07.

[2] Alina Arseniev-Koehler. Theoretical foundations and limits of word embeddings: what types of meaning can they capture?, 2021.

[3] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020.

[4] Imen Baklouti, Olfa Ben Ahmed, and Christine Fernandez-Maloigne. Cross-lingual low-resources speech emotion recognition with domain adaptive transfer learning. In *Proceedings of the 13th International Conference on Data Science, Technology and Applications - Volume 1: DATA*, pages 118–128. INSTICC, SciTePress, 2024.

[5] Russell Beale and Christian Peter. *The Role of Affect and Emotion in HCI*, pages 1–11. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.

[6] Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N. Chang, Sungbok Lee, and Shrikanth S.

Narayanan. Iemocap: Interactive emotional dyadic motion capture database. *Language Resources and Evaluation*, 42(4):335–359, Nov 2008.

[7] Samuel Cahyawijaya, Holy Lovenia, Willy Chung, Rita Frieske, Zihan Liu, and Pascale Fung. Cross-lingual cross-age group adaptation for low-resource elderly speech emotion recognition, 2023.

[8] Census and Statistics Department of Hong Kong. Hong kong population projections 2020-2069. https://www.statistics.gov.hk/pub/B1120015082020XXXXB0100.pdf.

[9] Pin-Yu Chen. Model reprogramming: Resource-efficient cross-domain machine learning, 2023.

[10] ctl. wav2vec2-large-xlsr-cantonese. https://huggingface.co/ctl/wav2vec2-large-xlsr-cantonese. Accessed: 2025-04-07.

[11] Xian Fan, Xiaoge Li, Feihong Du, Xin Li, and Mian Wei. Apply word vectors for sentiment analysis of app reviews. In *2016 3rd International Conference on Systems and Informatics (ICSAI)*, pages 1062–1066, 2016.

[12] Fasih Haider, Senja Pollak, Pierre Albert, and Saturnino Luz. Emotion recognition in low-resource settings: An evaluation of automatic feature selection methods. *Computer Speech  Language*, 65:101119, 2021.

[13] Abid Haleem, Mohd Javaid, Ravi Pratap Singh, and Rajiv Suman. Telemedicine for healthcare: Capabilities, features, barriers, and applications. *Sensors International*, 2:100117, 2021.

[14] Kaixin Ma, Xinyu Wang, Xinru Yang, Mingtong Zhang, Jeffrey M Girard, and Louis-Philippe Morency. Elderreact: A multimodal dataset for recognizing emotional response in aging adults. In *2019 International Conference on Multimodal Interaction*, ICMI '19, page 349–357, New York, NY, USA, 2019. Association for Computing Machinery.

[15] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013.

[16] Leonardo Pepino, Pablo Riera, Luciana Ferrer, and Agustín Gravano. Fusion approaches for emotion recognition from speech using acoustic and text-based features. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6484–6488, 2020.

[17] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022.

[18] toastynews. hong-kong-fasttext. https://github.com/toastynews/hong-kong-fastText. Accessed: 2025-04-07.

[19] Peter Torre and Jessica A. Barlow. Age-related changes in acoustic characteristics of adult speech. *Journal of Communication Disorders*, 42(5):324–333, 2009.