# A Comparative Study of Speaker Role Identification in Air Traffic Communication Using Deep Learning Approaches

DONGYUE GUO, National Key Laboratory of Fundamental Science on Synthetic Vision, College of Computer Science, Sichuan University, China

JIANWEI ZHANG, BO YANG, and YI LIN*, College of Computer Science, Sichuan University, China

Automatic spoken instruction understanding (SIU) of the controller-pilot conversations in the air traffic control (ATC) requires not only recognizing the words and semantics of the speech but also determining the role of the speaker. However, few of the published works on the automatic understanding systems in air traffic communication focus on speaker role identification (SRI). In this paper, we formulate the SRI task of controller-pilot communication as a binary classification problem. Furthermore, the text-based, speech-based, and speech and text based multi-modal methods are proposed to achieve a comprehensive comparison of the SRI task. To ablate the impacts of the comparative approaches, various advanced neural network architectures are applied to optimize the implementation of text-based and speech-based methods. Most importantly, a multi-modal speaker role identification network (MMSRINet) is designed to achieve the SRI task by considering both the speech and textual modality features. To aggregate modality features, the modal fusion module is proposed to fuse and squeeze acoustic and textual representations by modal attention mechanism and self-attention pooling layer, respectively. Finally, the comparative approaches are validated on the ATCSpeech corpus collected from a real-world ATC environment. The experimental results demonstrate that all the comparative approaches are worked for the SRI task, and the proposed MMSRINet shows the competitive performance and robustness than the other methods on both seen and unseen data, achieving 98.56%, and 98.08% accuracy, respectively.

CCS Concepts: • **Computing methodologies** → *Information extraction*.

Additional Key Words and Phrases: speaker role identification, air traffic control, text classification, speech classification, spoken instruction understanding, multi-modal learning

## 1 INTRODUCTION

Speech communication between air traffic controllers (ATCOs) and pilots is one of the most important interaction ways in air traffic control (ATC) procedures. Recently, there is increasing interest in introducing automatic spoken instruction understanding (SIU) techniques to empower the ATC process [Lin 2021; Pardo et al. 2011]. In the past decades, it has been extensively investigated and widely applied on the safety detection system [Lin et al. 2020, 2019], the ATCOs training devices [Smídl et al. 2019], the ATC assistance systems [Lin et al. 2021b; Schulder et al. 2015].

In practice, ATC speech communication can be regarded as task-oriented conversations between ATCOs and pilots. Thus, the SIU system of ATC is usually a pipeline that consists of several subtasks, i.e., speech activity detection (SAD), automatic speech recognition (ASR), text instruction understanding (TIU), and speaker role identification (SRI). In this pipeline, firstly, the speech segment is captured from the real-time streaming by the SAD module, and then the ASR system translates it into human-readable texts. Subsequently, the TIU module converts the natural texts into predefined structured instructions that are further processed by the computer. Finally, the

---

*Corresponding author.

Authors' addresses: DONGYUE GUO, dongyueguo@stu.scu.edu.cn, National Key Laboratory of Fundamental Science on Synthetic Vision, College of Computer Science, Sichuan University, Chengdu, China, 610000; JIANWEI ZHANG, zhangjianwei@scu.edu.cn; BO YANG, boyang@scu.edu.cn; YI LIN, yilin@scu.edu.cn, College of Computer Science, Sichuan University, Chengdu, China, 610000.

computer-readable instructions and the speaker role (ATCO or pilot which was output by the SRI module) jointly provide a conversation context for other downstream applications.

As can be seen from mentioned illustrations, the SRI module is a critical component of the SIU system in the field of ATC. However, most of the existing research of the ATC-related SIU systems focuses on the ASR and TIU techniques [LIN 2021; Lin et al. 2021a; Oualil et al. 2017; Zuluaga-Gomez et al. 2020], and no detailed description of the SRI task was presented. A instruction understanding model and ATC communication rule-based methods for the SRI tasks were studied in [Lin et al. 2019], without providing SRI performance. To the best of our knowledge, none of the published works have reported complete approaches and results for the SRI tasks in the ATC domain.

Since the ATCO communicates with several pilots by radio in a single frequency, the role of the speaker cannot be distinguished from the communication data link. However, the speaker role is a kind of indispensable and important information in many ATC-related applications, such as safety detection systems, ATCO workload analysis systems. Therefore, the inability to identify the speaker role directly from communication brings a certain challenge to the ATC-related SIU tasks. Fortunately, there are two kinds of data that can be served as the potential entities for the SRI tasks.

- Text: On the one hand, according to the communication rules recommended by the international civil aviation organization (ICAO), the ATCOs should declare the call sign of the target aircraft before issuing the instructions, while the pilots read back the instructions firstly and then reporting their call sign. In general, most of the controller-pilot speech communication follow these rules, allowing the text classification to be a promising technology for SRI tasks.
- Speech: On the other hand, the speech can be considered as a representation of the speaker role from two aspects of signal and text. a) The controller-pilot speech communication presents distinctive features depending on the equipment and environment, such as a microphone, push-to-talk (PTT), background noise, radio. b) It implies the representation of its transcripts, which further provides more discriminative knowledge for the SRI task.

In this paper, we define the SRI task as a binary classification problem, i.e., all the instructions are classified into two classes: ATCO or pilot. Meanwhile, the SRI task is addressed by the data-driven approaches from three different inputs, i.e., text, speech, speech-text. To this end, the text classification, audio classification, and multi-modal classification approaches are proposed to achieve the SRI task. In this procedure, several popular network architectures are introduced to serve as backbone networks for each approach to eliminate the impact of the difference between network architectures. The BiLSTM [Zhou et al. 2016], TextCNN [Kim 2014], and Transformer [Vaswani et al. 2017] architecture are developed as the backbone network in the text-based methods, while x-vector [Snyder et al. 2018], SincNet [Ravanelli and Bengio 2018a], and CRNN [Choi et al. 2017] architecture are built for the speech-based methods. Most importantly, a multi-modal speaker role identification network (MMSRINet) is designed to learn the distinctive representations from both the speech and textual modalities for the speech-text based methods. Specifically, a modal attention mechanism is proposed to fuse the different representations to a joint feature vector. In addition, the self-attention pooling layer is applied to produce the joint vector by the weighted sum operations, which further be regarded as the multi-modal embedding. Finally, the multi-modal embedding is further fed into the classifier to generate the final probabilities of the speaker role.

All the proposed methods were validated on the ATCSpeech corpus [Yang et al. 2020] that was collected from a real-world ATC environment. In addition, in order to analyze and compare the performance and robustness of the model, we evaluate the trained model in two ways: 1) The model is validated on the test set of the ATCSpeech to evaluate the performance on the seen samples.

2) A supplement test set called test-s is used to verify the robustness of the model on the unseen samples collected in controller-pilot communication.

In summary, our contributions are listed as follows:

- A thorough comparison between the aforementioned deep learning based SRI techniques is investigated. To the best of our knowledge, this is the first work that investigates the SRI task in the ATC domain.
- The robustness and performance of the comparative models are comprehensively analyzed and discussed on the seen and unseen samples.
- A multi-modal SRI network, called MMSRINet, is proposed to achieve the ATC-related SRI task by considering both the speech and textual modal knowledge, which shows more competitive performance and robustness than other methods.

This paper is organized as follows. A brief review of related works is described in Section 2. Section 3 presents the detail of the proposed methods and architectures. The experimental configurations are provided in Section 4. The experimental results are reported and detailly discussed in Section 5. Finally, this paper is concluded in Section 6.

## 2 RELATED WORK

### 2.1 Text Classification

Text classification is a classical task in the field of natural language processing (NLP), which aims to classify a given text sequences into a certain class. In general, the approach can be grouped into two categories: rule-based methods and data-driven based methods. The rule-based approach usually requires a large number of predefined rules and is strongly dependent on domain knowledge, which can only be applied to limited scenarios due to poor flexibility.

Thanks to the development of deep learning techniques, the performance of data-driven methods has generally outperformed that of rule-based methods in recent years and has become the standard paradigm of text classification tasks [Minaee et al. 2021]. Zeng et al. [Kim 2014] utilized a convolutional neural network (CNN) [Lecun et al. 1998] to achieve the sentence classification tasks which makes representative progress in the NLP domain. To capture the long-term dependencies, the Att-BiLSTM model [Zhou et al. 2016] was built on a recurrent neural network (RNN) [Mikolov et al. 2010]. Currently, various improved methods based on the CNN or RNN block were proposed to achieve the text classification task, such as Character-level CNNs [Zhang et al. 2015], tree-based CNN [Mou et al. 2016], Tree-LSTM [Tai et al. 2015], Multi-Timescale LSTM [Liu et al. 2015]. With the successful applications of the Transformer architecture [Vaswani et al. 2017], many Transformer-based and pretrained language models were also proposed and achieved surprising performance [Devlin et al. 2019; Radford et al. 2018]. These methods achieved new state-of-the-art performance in text classification tasks by fine-tuning the pretrained models [Sun et al. 2019].

### 2.2 Audio Classification

Audio classification is widely applied in audio pattern recognition tasks, such as speaker identification [Ravanelli and Bengio 2018a; Snyder et al. 2018], acoustic event detection [Kumar and Raj 2016], accent classification [Hansen and Liu 2016; Lopez-Moreno et al. 2014], audio emotion recognition [Jermsittiparsert et al. 2020]. Recently, deep learning methods showed promising performance compared to traditional approaches for this task [Hershey et al. 2017]. Enormous works have been investigated to explore different model architectures and applications for audio classification. Shawn Hershey et al. [Hershey et al. 2017] demonstrated that the CNNs used in the image classification task, such as AlexNet [Krizhevsky et al. 2012], VGG [Simonyan and Zisserman 2014], and ResNet [He et al. 2016], achieved desired performance for the large-scale audio classification

task. Meanwhile, the convolutional recurrent neural network (CRNN) was proposed and to achieve music classification [Choi et al. 2017], audio event detection [Cakır et al. 2017], audio tagging [Xu et al. 2018], etc. In addition, in recent years, there are increasing interest in learning features from raw waveforms directly instead of handcraft features. Mirco Ravanelli et al. proposed the SincNet [Ravanelli and Bengio 2018a] to achieve speaker recognition which employs band-pass filters (based on the parametrized Sinc functions) in the first convolutional layer. Jee-weon Jung et al. proposed the RawNet [Jung et al. 2019] to improve the performance of the speaker verification from raw waveforms. In short, deep learning-based audio classification is still an interesting task in many applications.

## 2.3  Multi-modal Classification

With the explosive growth of multi-modal data in the digit world, multi-modal learning is attracting increasing research interest and shows powerful performance than that of unimodal modal methods [Ngiam et al. 2011]. Various modalities can be used to achieve classification tasks, including audio-video [Nagrani et al. 2018], audio-text [Mittal et al. 2020], image-text [Gallo et al. 2018; Kiela et al. 2018], etc. In general, the fusion strategy of the classification task can be implemented in the following two ways: early fusion and later fusion [Baltrušaitis et al. 2019]. Early fusion methods fuse the multi-modal feature vectors to a joint representation that is further fed into the classifier, while the later fusion makes a second decision on the output of two classifiers by an extra strategy. Since the advantages of early fusion in exploring the correlations and interactions between different modalities, in this paper, we introduce the early fusion methods to the SRI task. In an early work [Kiela and Bottou 2014], direct concatenation was employed to produce multi-modal joint representations. In order to identify the correlations of learned multi-modal features, a structural regularization was proposed in [Wu et al. 2014] to empower the deep neural network (DNN) based fusion layer, which also preserves the diversity of the different modality features.

In addition to the classification task, more powerful fusion methods were successfully integrated into the ASR and NLP architectures. Modality attention was proposed to fuse the audio-visual features for ASR tasks in [Zhou et al. 2019]. The works of [Fukui et al. 2016] and [Ovalle et al. 2017] used compact bilinear pooling and complex gating mechanisms to obtain multi-modal representations.

## 3  METHODOLOGY

### 3.1  The Unified Framework of Unimodal SRI

Since the SRI task is formulated as the binary classification problem, both the text-based and speech-based methods can be further refined as the classification task of sequential data with a variable length. In order to ensure the fairness of the comparison, a unified classification framework is designed for the text-based and speech-based methods. As illustrated in Fig. 1, the proposed framework includes an input module, feature encoder, and classifier. The detailed descriptions are shown below:

The input module is a conceptual component that can be compatible with both the text and speech input. In the text-based method, the input module is a word embedding layer to learn the representations of the transcripts. For the speech-based method, the input module is a speech signal preprocessor for raw waveforms, or a speech feature extractor for handcraft speech features.

The feature encoder consists of the backbone network and a pooling layer. The former extracts the high-level representations from the input sequence, while the latter squeezes the feature to a fixed dimensional vector. Specifically, let a sequence of features $x = \{x_1, x_2, ..., x_t\}$ that outputted by the input module, the backbone network is performed to generate the high-level feature representations
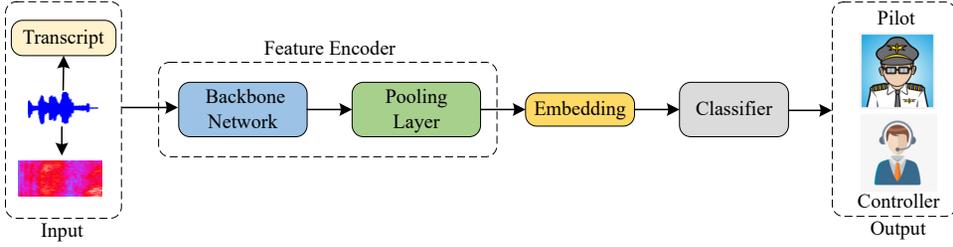
Fig. 1. A unified framework of text-based and speech-based methods towards SRI tasks.

$h = \{h_1, h_2, ..., h_t\}$ as Eq. (1), where $h_t \in \mathbb{R}^D$ is a D-dimensional feature vector. Then, as shown in Eq. (2), the pooling layer squeezes the $h$ to a fixed dimensional embedding vector $z \in \mathbb{R}^1$ on the temporal dimension by the pooling function $G$. Moreover, in this framework, the backbone network and the pooling layer are presented as the plugins, which are further adapted to inputs (text or speech) and facilitate the comparison of different algorithms.

$$h = Backbone(x) \tag{1}$$

$$z = G(h) \tag{2}$$

Three fully connected (FC) layers are designed in the classifier, i.e., two hidden layers, and an output layer, which further transforms the embedding into the probability of the speaker role. The ReLU is selected as the activation function for hidden layers, while the output layer adopts Softmax activation. In addition, batch normalization is designed before each activation.

## 3.2 Text-based method

*3.2.1 Transcription of ATC speech.* The text-based methods of the SRI task are often cascaded with the ASR system. In short, the transcripts of the ATC speech segments are recognized by the ASR system and serve as input for the SRI task. Therefore, in this case, some existing approaches can be introduced to solve this problem, such as text classification, grammar matching.

The core idea of the text-based SRI approach is based on the ATC rules issued by ICAO, both of the ATCOs and the pilots are expected to speak in a rigorous stylized sentence. Specifically, as group A in Table 1, the ATCOs must indicate the call sign of the target flight before speaking details of the instructions. On the contrary, the pilots report their call sign after reading back the instructions in the repetition procedure. However, in practice, some ATC instructions break the ATC rules, which causes extra burdens for text-based approaches. For example, as group B in Table 1, the pilot instruction starts with a call sign, which deviates from the ICAO recommendations.

Table 1. The example of the ATC speech transcripts.

| Group | Role | Transcript |
|---|---|---|
| A | Controller | *Air China four two three seven*, climb to eight thousand one meters |
| | Pilot | Climb to eight thousand one meters, *Air China four two three seven* |
| B | Pilot | *Cathay two five four*, request heading two seven zero |
| | Controller | *Cathay two five four*, heading two seven zero agreed |

In summary, the recommended ATC rules are the primary basis for the text-based SRI approach, which is able to achieve a preferred performance for most speech communications in the ATC environment. For the speech instructions that deviate from the ATC rules, the model is expected to learn discriminative features from a large amount of dataset.

*3.2.2  Network Design.* In the text-based methods, the word sequence of the ATC speech transcriptions is firstly fed into an embedding layer to generate the word representations. In succession, the backbone network is applied to learn high-level representations that are further converted into an embedding by the pooling layer. Finally, the classifier (as described in Section 3.1) is designed to estimate the probabilities of the speaker role.

In this work, three competitive neural architectures are selected as the backbone network to achieve text-based SRI tasks, i.e., bi-directional long short-term memory network (BiLSTM), CNN, and Transformer network, which are the well-known models for deep learning techniques. The detailed descriptions of the aforementioned backbone networks are summarized as follows:

- **BiLSTM**: In the RNN-based architecture, two BiLSTM Layers with 512 neurons are designed to serve as the backbone network. Compared to unidirectional LSTM networks, the BiLSTM networks is able to consider information from both the future and past dimensions to learn the temporal dependences of the input sequence.
- **CNN**: As referred to [Kim 2014], three CNN blocks (with variable kernel size) are applied to extract position-invariant representations from the embedding vector. Then, the feature map is produced by cascading the output of CNN blocks. The CNN block is constructed by concatenating the Conv2D layer, ReLU activation, and a Max Pooling layer. The size of CNN filters is set to (3, 4, 5), corresponding to different receptive fields.
- **Transformer**: The backbone network consists of 4 Transformer blocks, where the basic block consists of a masked multi-head attention module, a layer normalization, and a position-wise feed forward layer. In this work, we adopt 4 heads in the attention module, and the dimension of the feed forward layer is set to 512.

For the pooling layer, a self-attention based temporal pooling strategy is applied to the SRI model. Let the feature map $H = \{h_1, h_2, ..., h_t\}$ that generated by backbone network, where $t$ is the length of the input sequence. The embedding $e$ is produced by a weighted sum of the feature map $H$, in which the weight is calculated by the attention mechanism. The inference rule can be summarized as Eq. (3)-(5), where $H \in \mathbb{R}^{h^b \times t}$, $h^b$ is the dimension of the output vector, W is a trainable weight.

$$H^* = tanh(H) \tag{3}$$

$$\alpha = softmax(\mathrm{W}^T H^*) \tag{4}$$

$$e = tanh(H\alpha^T) \tag{5}$$

## 3.3  Speech-based method

*3.3.1  SRI-related features of the ATC Speech.* In the ATC procedure, the speech of controller-pilot communication is transmitted through very high frequency (VHF) radiotelephony, in which the ATCO's speech is land-to-air whereas the pilot's speech is air-to-land. Thus, the special representations in the radio, microphones, and background noise (control room and aircraft cockpit) will be presented in the speech signal. Fig. 2 shows the spectrogram of ATC speech, both the ATCO and the pilot speeches are collected from different speakers and control sectors.
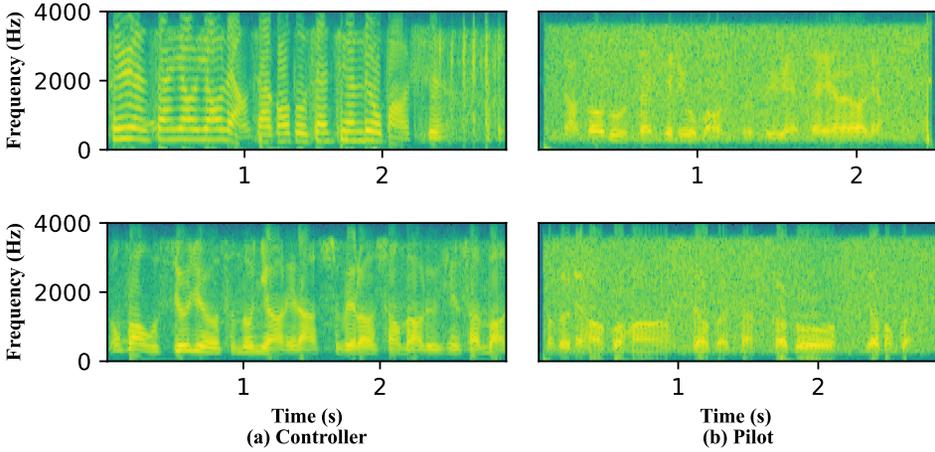
Fig. 2. Selective examples of ATC speech spectrogram. **(a)** The spectrogram of the controller speech and **(b)** The spectrogram of the pilot speech.

It is evident that the feature intensity distribute in different frequencies between ATCOs and pilots. For example, the frequency energy distribution of the ATCO speech is stronger than that of the pilot above 3000 Hz. In addition, different background noise models are presented for different speeches, specifically, the background noise distribution of the pilot speech is in a uniform manner whereas that of the ATCO is volatile. Moreover, the transcription implied in the speech is also the key information to support the SRI task. Therefore, the audio classification is expected to be a promising approach for the speech-based SRI task.

*3.3.2   Feature Encoder Design.* In the speech-based approaches, we also select three kinds of advanced neural networks of the speech-related tasks as the feature encoder, i.e., CRNN [Choi et al. 2017], X-vector [Snyder et al. 2018], SincNet [Ravanelli and Bengio 2018a]. Since the spectrogram is a kind of detailed and visualized handcraft features and has been widely used in audio classification tasks [Han et al. 2017; Zeng et al. 2019], the spectrogram of the raw speech serves as the input of CRNN, X-vector backbone networks. For the SincNet, due to its ability to process the speech signal, the raw waveform is directly fed into the model to achieve the SRI task. The detailed configurations of the feature encoder are described as follows:

- **CRNN**: CRNN is a popular architecture in the audio classification task. In this work, the CRNN-based backbone network is implemented by referring to [Choi et al. 2017]. The CRNN consists of 5 CNN blocks and 2-layer RNN with gated recurrent units (GRU). The CNN blocks are to recognize spatial-cross patterns by filters and downsample the features by the Max-pooling operation, whereas RNNs aim to summarize the learned patterns in the temporal dimension. Followed the backbone network, the self-attention pooling layer is applied to generate speech embeddings that are further fed into the classifier.
- **X-vector**: The X-vector system was proposed to extract DNN embeddings for speaker recognition [Snyder et al. 2018]. In this work, the front-end of the X-vector system is selected as the feature encoder, which includes 4 time-delay deep neural network (TDNN) layers and a statistics pooling layer. The frame-level representations are learned by TDNN layers, and the statistic pooling layer aggregates it to the sentence-level embedding.

- **SincNet**: The SincNet is a novel and effective CNN architecture for the speaker and speech recognition tasks with raw waveforms [Ravanelli and Bengio 2018a,b]. In the SincNet block, band-pass filters are applied to replace the standard CNN filters to convolves the waveform, which shows better model convergence and performance compared to the CNN block. In order to analyze the impact of the input features for the speech-based SRI tasks, instead of the handcraft features, a Sinc convolutional layer is designed before the TDNN module of X-vector models to conduct the SincNet-based backbone network.

### 3.4 MMSRINet

According to the aforementioned analysis, it is well received that the performance of text-based methods generally relies on ATC grammar, whereas the speech-based methods closely relate to the communication environment (equipment and background noise, etc.). The performance of text-based methods will be significantly reduced if the speech instructions deviate from the pre-defined ATC grammar. Similarly, the speech-based methods will suffer poor accuracy when it works on unseen data (i.e., communication environment not be covered by the training set). Following this idea, we believe that the discriminative knowledge learned by the text-based method and the speech-based method are complementary and can be further used to improve the accuracy of the SRI task. To this end, we present a multi-modal speaker role identification network, called MMSRINet, to consider both the textual modal knowledge and the signal modal characteristics of the ATC speech.

The architecture of the MMSRINet is illustrated in Fig. 3, consisting of the textual encoder, speech encoder, modal fusion module, and a classifier. Specifically, the high-dimensional textual and acoustic representations are learned by the textual encoder and speech encoder, respectively. Then, the representations from different modalities are aggregated and squeezed by the modal fusion module. Finally, the fused embedding is fed into the classifier to generate the final probability of the speaker role.

The textual encoder and speech encoder is implemented on the BiLSTM and CRNN, as described in Section 3.2.2 and Section 3.3.3. Furthermore, the composition of the CNN block also is shown in Fig. 3. The modal fusion module is constructed by the modal attention mechanism and the self-attention pooling layer, where the self-attention pooling operation is similar to that described in Section 3.2.2. The classifier is the same as that of in unimodal SRI framework described in Section 3.1. In this section, we mainly focus on fusing the representations between different modalities. Due to the heterogeneity of the speech and the text, aggregating these representations to a joint feature vector is still a unique computational challenge in the SRI task. In this work, a modal attention
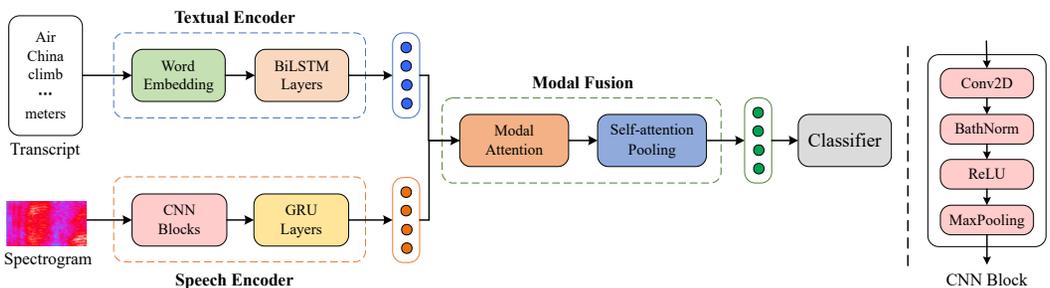


Fig. 3. The architecture of the proposed MMSRINet.

mechanism is designed to capture the correspondences between acoustic and textual modalities, as described below:

Let the high-dimensional representations output by the speech encoder $h^s = \{h_1^s, h_2^s, ..., h_n^s\}$, and the textual encoder $h^t = \{h_1^t, h_2^t, ..., h_m^t\}$, where the $m$, $n$ are the length of the textual and speech features, respectively. Eq. (6)-(9) show the fusion operation, firstly, the correlated score $e_{ij}$ between the speech feature of time step $i$ and textual feature of time step $j$ are calculated by the scoring function $Energy$ as Eq. (6), where the $W_a$ is trainable parameters. Secondly, the weights $\alpha_{ij}$ of the modality attention are generated by the $Softmax$ function (Eq. (7)). Thirdly, the correlated vector $c_i$ between the speech representation of time step $i$ and textual representation is computed by the weighted sum operation. Finally, the concatenation operator is applied to combine the correlated vector $c_i$ and speech feature $h_i^s$, and the final fused feature $f_i$ of time step $i$ can be produced as Eq. (9), where the $W_b$ is trainable parameters.

$$e_{ij} = Energy(h_i^s, h_j^t) = h_i^{sT} W_a h_j^t \tag{6}$$

$$\alpha_{ij} = Softmax\left(e_{ij}\right) = \frac{\exp\left(e_{ij}\right)}{\sum_{k=1}^{M} \exp\left(e_{ik}\right)} \tag{7}$$

$$c_i = \sum_{j=1}^{M} \alpha_{ij} h_j^t \tag{8}$$

$$f_i = tanh(W_b[c_i, h_i^s]) \tag{9}$$

The purpose of the modal attention mechanism is to calculate the correlations between textual representations and acoustic representations, so as to enhance the implicit textual information related to SRI tasks in acoustic representations, especially the position and semantic representations of the call signs. Further, the proposed modal attention mechanism also relieves the impacts of the communication environment, which further improves the robustness of the model.

## 4 EXPERIMENTAL SETUP

### 4.1 The ATC Corpus

In this work, the experiment is conducted on the ATCSpeech corpus [Yang et al. 2020] that was collected from real-world ATC communications. The transcripts and speaker roles were manually labeled to support the approach validation. There are about 26.52 hours (25,765 utterances) of ATCO speech and 31.29 hours (35,895 utterances) of pilot speech. Some minor utterances (about 0.013 hours) without the speaker role information during the annotation were labeled *unknown* and be excluded in our experiment. The sample rate of all samples in the ATCSpeech corpus are 8,000 Hz. In the training phase, all the proposed models are developed on the train set and tuned parameters on the dev set.

To further evaluate the performance and robustness of the different methods, two test datasets are used in the evaluation phase, i.e., **test** set in ATCSpeech and a supplementary test set called **test-s**. The test-s dataset was collected from the Chengdu ATC center, which not be covered by the training dataset. The total duration of the test-s dataset is about 2 hours (1,930 utterances), and the labels of the speech are also annotated manually. The main purpose of introducing the test-s set is to evaluate the robustness of the model for the unseen ATC environment.

## 4.2 Comparison models

Since the pretraining models have been shown competitive performance in language modeling and speech processing, the pretraining methods have become a new paradigm of deep learning. In order to present a comprehensive comparison of the deep learning based SRI models and validate the performance of the proposed unimodal SRI framework and the MMSRINet, the pretraining models of the speech and the text are cascaded with the classifier to develop the SRI task in the experiment. The detailed descriptions of the pretraining based SRI models are listed as follows.

- **Pretrained-T**: The pretrained-T SRI model is a text-based approach that is finetuned with a pretrained BERT model [Devlin et al. 2019] on the ATCSpeech corpus. In this work, the BERT model is trained on the ATCSpeech-Large corpus which contains 120,000+ transcriptions of the ATC communication utterances.
- **Pretrained-S**: The pretrained-S is a speech-based SRI model which is finetuned with a pretrained Wav2vec 2.0 model [Baevski et al. 2020] on the ATCSpeech corpus. The Wav2vec 2.0 model is trained on an unlabeled speech corpus called RUD which includes about 2600+ hours of speech, 2.87 million utterances.

The details of the RUD and the ATCSpeech-large corpus can be found in our previous work [LIN 2021].

## 4.3 Experimental Configurations

To ensure the fairness of comparison between different models and methods, we used the same configuration in common blocks. Specifically, the dimension of the word embedding is set to 512, and the neurons of two FC layers in the classifier are set to 256 and 2, respectively. Both the dimension of the embedding outputted by the feature encoder (SRI framework) and modal fusion module (MMSRINet) are set to 512.

In the experiment of speech-based methods, the spectrogram is computed by 80 linearly spaced log-filterbanks with 25ms windows and 15ms overlaps. In addition, for the SincNet, the raw speech waveforms are performed normalization and fed into the model directly. Due to the multilingual nature of the ATCspeech corpus, the Chinese character and the English word are served as the basic vocabulary tokens in the text-based methods. There are 1284 tokens in our vocabulary, includes 698 Chinese characters, 584 English words, and two special tokens, i.e., <PAD>, <UNK>.

In this work, we construct and train all the models with the open-source deep learning framework PyTorch 1.4.0. The training server was configured as follows: Ubuntu 16.04 operating system with 2*NVIDIA GeForce RTX 2080Ti GPU, Intel Xeon E5-2630 CPU, and 128 GB memory. The Adam optimizer with $10^{-4}$ initial learning rate is applied to all training tasks and the cross-entropy is selected as the loss function. An early stopping strategy is performed to terminate the training procedure by observing the SRI performance on the dev dataset.

Since the samples of ATCO and pilot are imbalanced in the ATCSpeech corpus, the performance of all models is measured by the accuracy (ACC %), precision, recall, F1-score (F1), and area under the receiver operating characteristic curve (AUC) on both the test set and test-s set. Specifically, ACC considers the ratio of samples that are correctly predicted, while AUC investigates the ability of the classifier to identify the ATCO utterance and pilot utterance as different thresholds are selected. Precision, recall, F1 are commonly used metrics for binary classifiers, the pilot utterance serves as the positive class when these metrics are calculated in the experiment.

## 5 RESULTS AND DISCUSSIONS

### 5.1 Results of text-based Methods

The results of the text-based methods are reported in Table 2. For the proposed unimodal SRI framework, compared to the TextCNN and Transformer approach, BiLSTM based backbone achieves better performance on both the test and test-s sets, reaches 97.05% and 94.97% accuracy, respectively. The results benefit from the powerful temporal modeling of RNN architectures in time series data. The Transformer model suffers from poor accuracy and the lowest AUC in the test-s set while it obtains the higher AUC in the test set. It can be attributed to the Transformer model does not perform its full capacity on small-scale data sets, and optimizing it with enormous samples may be a promising way to achieve desired performance improvement. Overall, the three text-based models achieve comparable performance, and the accuracy is 96%-97% in the test set. However, the accuracy, precision, F1-score are significantly reduced on the test-s set, i.e., about 2%. In general, the design of the backbone network is not the key factor that affects the model performance in the text-based SRI approaches.

For the Pretrained-T model, it obtained better performance than the proposed unimodal models in most evaluation metrics, especially in AUC. It can be attributed that the pretrained model learned more discriminative features about the SRI task and achieve a stronger classification ability. We also recognize that the Pretrained-T model obtained comparable performance on test and test-s set and achieved higher robustness than the proposed SRI framework. However, these models are usually pretrained on a large amount of data and are not suitable for low-resource conditions.

Table 2. The results of the text-based methods. The Proc. and Reca. represent precision and recall, respectively.

| Methods | Test | | | | | Test-s | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ACC % | AUC % | Prec. % | Reca. % | F1 % | ACC % | AUC % | Prec. % | Reca. % | F1 % |
| **BiLSTM** | 97.05 | 77.86 | **97.39** | 97.59 | 97.49 | 94.97 | 75.61 | **95.39** | 96.58 | 95.98 |
| **TextCNN** | 96.94 | 78.70 | 96.99 | 97.90 | 97.44 | 94.61 | 79.98 | 93.84 | 97.75 | 95.75 |
| **Transformer** | 96.13 | 81.55 | 96.27 | 97.38 | 96.77 | 94.87 | 67.71 | 94.35 | 97.58 | 95.94 |
| **Pretrained-T** | **97.13** | **97.26** | 95.05 | **97.99** | **97.56** | **97.46** | **97.74** | 94.62 | **98.90** | **97.92** |

Table 3. Examples of different types of prediction error samples. In this table, *akube*, *scooter* are OOV tokens, and the probabilities were output by the BiLSTM model.

| Type | No. | Transcript | Role | Probability |
|---|---|---|---|---|
| OOV | 1 | Hainan seven four five two, direct to *akube* | Controller | 0.27 |
| | 2 | *Scooter* one two five, climb maintain seven thousand two hundred meters | Controller | 0.17 |
| DFG | 3 | Descend maintain eight thousand one hundred meters | Pilot | 0.48 |
| | 4 | Bohai seven four nine seven, confirm bemta one xray | Pilot | 0.39 |
| | 5 | Confirm, bohai seven four nine seven | Controller | 0.34 |

By analyzing the experimental results, two kinds of samples contribute to the performance degradation, i.e., out of vocabulary (OOV) and deviate from grammar (DFG). Among the misclassification samples of the unimodal SRI framework in the test-s set, the OOV samples account for about 30%, while DFG sentences account for 70%. The performance improvements of the Pretrained-T model

Table 4. The results of the speech-based methods and MMSRINet.

| Methods | Test | | | | | Test-s | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | ACC % | AUC % | Prec. % | Reca. % | F1 % | ACC % | AUC % | Prec. % | Reca. % | F1 % |
| CRNN | 97.62 | 68.44 | 98.21 | 97.80 | 98.0 | 94.30 | 69.87 | 92.91 | 98.33 | 95.54 |
| X-vector | 97.44 | 68.89 | 96.54 | **99.26** | 97.88 | 93.78 | 72.31 | 91.15 | **99.58** | 95.18 |
| SincNet | 97.06 | 70.35 | 96.23 | 98.95 | 97.57 | 95.54 | 71.20 | **98.94** | 93.83 | 96.32 |
| Pretrained-S | 98.13 | **97.82** | 98.46 | 96.57 | 98.50 | 97.56 | **97.23** | 98.41 | 95.51 | 97.98 |
| MMSRINet | **98.56** | 70.49 | **98.83** | 98.91 | **98.87** | **98.08** | 82.03 | 97.63 | 99.16 | **98.39** |

on the test-s set are mainly benefited from the robustness of the pretrained model to the OVV tokens. Some selective samples with different prediction errors are listed in Table 3. The OOV tokens are mainly derived from the route waypoints (Table 3 No.1) and airline call signs (Table 3 No.2) that are unseen in the ATCSpeech corpus. Especially for airline call signs, the model is easy to be confused by the location of call signs and reports the error results in the prediction. The samples of DFG sentences can be divide into two types: the speech without call signs and spoken call signs break the ICAO recommendation. The former occurs mainly in the instructions repetition (as shown in Table 3 No.3), while the latter tends to occur in the consultative conversation between the ATCOs and pilots (Table 3 No.4 and No.5).

In conclusion, processing the speech that contains OOV tokens and DFG sentence is still a huge challenge for the text-based SRI task. Indeed, identifying the speaker role of the DFG sentence in the ATC environment based on transcript alone and without conversation context is also a difficult task for human understanding, which makes a huge limitation of the text-based methods in the SRI task. In addition, the text-based methods are often cascaded with an ASR system in real-time SIU applications, so the final accuracy is also impacted by the ASR performance.

## 5.2 Results of Speech-based Methods and MMSRINet

In this work, the core idea of the proposed modal attention mechanism is to utilize transcriptions to empower the implicit textual representations of the speech. Therefore, the MMSRINet can be considered as the variant of the speech-based methods. In this section, we report and discuss the results of speech-based methods and MMSRINet together.

The results are reported in Table 4. For the unimodal SRI framework, the performance of speech-based methods is generally better than that of text-based methods, i.e., all of them achieve over 97% ACC. However, the AUC of the speech-based methods is lower than that of text-based methods. As can be seen from the results, the CRNN have obtained the higher ACC in the test set, while the SincNet presented a more competitive performance in the test-s set, reaching 97.62% and 95.54%, respectively. For the Pretrained-S model, as discussed in Section 5.1, it also obtained better performance than that without the pretraining process and harvest the highest AUC in speech-based SRI models.

We further explore the incorrect samples and two key factors are regarded to affect the performance of the speech-based models. Firstly, compared to text-based methods, the percentage of DFG speech is significantly reduced in incorrect samples, about only 10%. The noise of the speech is the primary reason to limit the model performance. Secondly, since the sample of the ATCOs and pilots are imbalanced in the training set, the speech-based model of the unimodal SRI framework is prone to be overfitted. The model prefers to predict the sample as the class most appear in the training set (pilot), which is also the reason for the lower AUC. The AUC is greatly improved in the

Pretrained-S model which might be learned a more robust speaker role related representation in the pretrained process.

The results of the MMSRINet demonstrated that utilizing both the speech and textual modality features is a feasible approach for the SRI tasks. The proposed MMSRINet reports the best accuracy and F1-score on both the test set and test-s set, achieving 98.56%, 98.08% ACC, and 98.87%, 98.39% F1-score, respectively. Actually, the acoustic features and textual representations of the ATC speech are complementary knowledge for the SRI tasks. In addition, by the proposed modal attention mechanism, the acoustic characteristics and ATC grammars of the speech can be properly considered in the proposed MMSRINet.

In summary, the following conclusions can be obtained from the experimental results:

(1) The text-based methods usually rely on the ATC grammars and are suitable for the general conversation contexts. The OOV issue is a challenge when the model is migrated to an unseen ATC environment.

(2) Speech-based methods obtain better ACC in the test set than text-based methods. The representation difference in the acoustic features between ATCOs and pilots is regarded as a piece of key knowledge for the SRI tasks. Since acoustic features are prone to be changed and distorted in unseen or noisy channels, the performance of the speech-based methods will be limited by the real-time communication environment.

(3) Pretraining based models can learn more useful SRI-related features on pretraining process and improve the robustness and classification ability of the SRI models on unseen data. But it requires large amounts of data to fit the pretrained model in pretraining process, not suitable for the low-resource conditions.

(4) The experimental result demonstrates that the multi-modal approach is a competitive solution for the SRI task. It learns the complementarity knowledge from both the acoustic features and textual grammar of the ATC speech, providing more discriminative information for the classifier. Thanks to the multi-modal inputs, the model presented better performance and robustness which is beneficial for unseen datasets.

## 5.3 Ablation Study

In this section, we design kinds of ablation experiments to verify the effectiveness of the proposed self-attention pooling strategy and modal attention mechanism in the MMSRINet. All experiments are conducted in the ATCSpeech corpus and evaluated on the test set.

*5.3.1 Self-attention pooling vs. other pooling strategy.* To validate the effectiveness of the proposed self-attention pooling strategy, the sum and average operation on the temporal dimension are applied to replace the self-attention pooling in the MMSRINet. The results are presented in Table 5, it can be seen from the experimental results that the sum operation is slightly inferior to the self-attention pooling strategy and obtained the 98.37 % accuracy and 98.64 % F1-score. The performance of the average pooling operation only achieve 97.62 % accuracy in the test set. This is because the average pooling strategy assumes that the features of each time step make the same contributions of the SRI task. In contrast, the proposed self-attention pooling strategy generates a weight coefficient at each time step, which is conducive to highlighting SRI-related features and improving performance of the SRI task.

*5.3.2 Modal attention vs. concatenation.* As described in Section 3.4, the motivation of the proposed modal attention mechanism is to utilize transcriptions to empower the implicit textual representations of the speech. To validate the effectiveness of the proposed modal attention mechanism,

Table 5. The results of the different pooling strategy in MMSRINet.

| Pooling strategy | ACC% | AUC% | Precision % | Recall % | F1 % |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Average | 97.62 | 64.59 | 96.64 | 99.47 | 98.04 |
| Sum | 98.37 | 69.21 | 98.43 | 98.84 | 98.64 |
| Self-attention (Ours) | 98.56 | 70.49 | 98.83 | 98.91 | 98.87 |

the concatenate operation is applied to alternate the modal attention and generated the joint representation vector.

The results are shown in Table 6, the performance of proposed modal attention is superior to that of the concatenate operation in most evaluation metrics. By analyzing the misclassification samples, we found that the performance will be reduced using concatenate operation when the transcription can not provide discriminative features for the SRI task (Such as the DFG sentence No.3 and No.5 in Table 3.). Note that it also harvests better performance than unimodal methods when concatenating the speech and textual representation directly. Thus, it is an effective technique to achieve the high accuracy SRI tasks by considering multi-modal features that further support our motivation.

Table 6. The results of concatenation operation vs. the proposed modal attention module.

| Fusion methods | ACC% | AUC% | Precision % | Recall % | F1 % |
|:---:|:---:|:---:|:---:|:---:|:---:|
| Concatenation | 98.25 | 65.50 | 97.83 | 99.26 | 98.54 |
| Modal attention (Ours) | 98.56 | 70.49 | 98.83 | 98.91 | 98.87 |

In short, the above ablation experiments demonstrate that the proposed self-attention pooling strategy and modal attention module are helpful to improve the performance of the SRI models.

## 6 CONCLUSIONS

In this paper, we presented a comprehensively comparative study for the SRI tasks using deep learning approaches in the ATC domain. Three kinds of methods with different inputs were investigated to solve the problems of the SRI tasks, i.e., text-based methods, speech-based methods, and multi-modal methods. Firstly, we formulated the SRI task as the binary classification problem, and further refine the above methods as text classification, speech classification, and multi-modal classification task. Secondly, the efficacy of the above methods is confirmed by theoretical and experimental demonstrations. Finally, the experiments demonstrated that the proposed MMSRINet is a competitive approach that achieves the best performance and robustness in the seen and unseen ATC environments.

In the future, we plan to explore more efficient and effective approaches for the SRI tasks. In addition, the fusion and application of the multi-modal data in the ATC environment would be also an interesting research topic.

## REFERENCES

2021. ATCSpeechNet: A multilingual end-to-end speech recognition framework for air traffic control systems. *Applied Soft Computing* 112 (2021), 107847. https://doi.org/10.1016/j.asoc.2021.107847

Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.*

Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. 2019. Multimodal Machine Learning: A Survey and Taxonomy. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 2 (Feb. 2019), 423–443. https://doi.org/10.1109/TPAMI.2018.2798607

Emre Cakır, Giambattista Parascandolo, Toni Heittola, Heikki Huttunen, and Tuomas Virtanen. 2017. Convolutional Recurrent Neural Networks for Polyphonic Sound Event Detection. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 25, 6 (2017), 1291–1303.

Keunwoo Choi, György Fazekas, Mark Sandler, and Kyunghyun Cho. 2017. Convolutional Recurrent Neural Networks for Music Classification. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2392–2396.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. https://doi.org/10.18653/v1/n19-1423

Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. 2016. Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016*. The Association for Computational Linguistics, 457–468. https://doi.org/10.18653/v1/d16-1044

I. Gallo, A. Calefati, S. Nawaz, and M. K. Janjua. 2018. Image and Encoded Text Fusion for Multi-Modal Classification. In *2018 Digital Image Computing: Techniques and Applications (DICTA)*. 1–7. https://doi.org/10.1109/DICTA.2018.8615789

Yoonchang Han, Jae-Hun Kim, and Kyogu Lee. 2017. Deep Convolutional Neural Networks for Predominant Instrument Recognition in Polyphonic Music. *IEEE ACM Trans. Audio Speech Lang. Process.* 25, 1 (2017), 208–221. https://doi.org/10.1109/TASLP.2016.2632307

John HL Hansen and Gang Liu. 2016. Unsupervised Accent Classification for Deep Data Fusion of Accent and Language Information. *Speech Communication* 78 (2016), 19–33.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Las Vegas, NV, USA, 770–778. https://doi.org/10.1109/CVPR.2016.90

Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss, and Kevin Wilson. 2017. CNN Architectures for Large-Scale Audio Classification. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 131–135. https://doi.org/10.1109/ICASSP.2017.7952132

Kittisak Jermsittiparsert, Abdurrahman Abdurrahman, Parinya Siriattakul, Ludmila A Sundeeva, Wahidah Hashim, Robbi Rahim, and Andino Maseleno. 2020. Pattern Recognition and Features Selection for Speech Emotion Recognition Model Using Deep Learning. *International Journal of Speech Technology* 23, 4 (2020), 799–806.

Jee-weon Jung, Hee-Soo Heo, Ju-ho Kim, Hye-jin Shim, and Ha-Jin Yu. 2019. RawNet: Advanced End-to-End Deep Neural Network Using Raw Waveforms for Text-Independent Speaker Verification. *arXiv:1904.08104 [cs, eess]* (July 2019). arXiv:1904.08104 [cs, eess]

Douwe Kiela and Léon Bottou. 2014. Learning Image Embeddings Using Convolutional Neural Networks for Improved Multi-Modal Semantics. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 36–45.

Douwe Kiela, Edouard Grave, Armand Joulin, and Tomás Mikolov. 2018. Efficient Large-Scale Multi-Modal Classification. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18)*. AAAI Press, 5198–5204.

Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A Meeting of SIGDAT, a Special Interest Group of the ACL*. ACL, 1746–1751. https://doi.org/10.3115/v1/d14-1181

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet Classification with Deep Convolutional Neural Networks. *Advances in neural information processing systems* 25 (2012), 1097–1105.

Anurag Kumar and Bhiksha Raj. 2016. Audio Event Detection Using Weakly Labeled Data. In *Proceedings of the 24th ACM International Conference on Multimedia*. 1038–1047.

Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998. Gradient-Based Learning Applied to Document Recognition. *Proc. IEEE* 86, 11 (1998), 2278–2324. https://doi.org/10.1109/5.726791

Yi Lin. 2021. Spoken Instruction Understanding in Air Traffic Control: Challenge, Technique, and Application. *Aerospace* 8, 3 (2021), 65.

Yi Lin, Linjie Deng, Zhengmao Chen, Xiping Wu, Jianwei Zhang, and Bo Yang. 2020. A Real-Time ATC Safety Monitoring Framework Using a Deep Learning Approach. *IEEE Trans. Intell. Transp. Syst.* 21, 11 (2020), 4572–4581. https://doi.org/10.1109/TITS.2019.2940992

Yi Lin, Dongyue Guo, Jianwei Zhang, Zhengmao Chen, and Bo Yang. 2021a. A Unified Framework for Multilingual Speech Recognition in Air Traffic Control Systems. *IEEE Transactions on Neural Networks and Learning Systems* 32, 8 (2021), 3608–3620. https://doi.org/10.1109/TNNLS.2020.3015830

Yi Lin, Xianlong Tan, Bo Yang, Kai Yang, Jianwei Zhang, and Jing Yu. 2019. Real-time Controlling Dynamics Sensing in Air Traffic System. *Sensors* 19, 3 (2019), 679. https://doi.org/10.3390/s19030679

Yi Lin, YuanKai Wu, Dongyue Guo, Pan Zhang, Changyu Yin, Bo Yang, and Jianwei Zhang. 2021b. A Deep Learning Framework of Autonomous Pilot Agent for Air Traffic Controller Training. *IEEE Transactions on Human-Machine Systems* (2021), 1–9. https://doi.org/10.1109/THMS.2021.3102827

Pengfei Liu, Xipeng Qiu, Xinchi Chen, Shiyu Wu, and Xuanjing Huang. 2015. Multi-Timescale Long Short-Term Memory Neural Network for Modelling Sentences and Documents. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015*. 2326–2335. https://doi.org/10.18653/v1/d15-1280

Ignacio Lopez-Moreno, Javier Gonzalez-Dominguez, Oldrich Plchot, David Martinez, Joaquin Gonzalez-Rodriguez, and Pedro Moreno. 2014. Automatic Language Identification Using Deep Neural Networks. In *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 5337–5341. https://doi.org/10.1109/ICASSP.2014.6854622

Tomás Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association*. ISCA, Makuhari, Chiba, Japan, 1045–1048.

Shervin Minaee, Nal Kalchbrenner, Erik Cambria, Narjes Nikzad, Meysam Chenaghlu, and Jianfeng Gao. 2021. Deep Learning Based Text Classification: A Comprehensive Review. *arXiv:2004.03705 [cs, stat]* (Jan. 2021). arXiv:2004.03705 [cs, stat]

Trisha Mittal, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. 2020. M3er: Multiplicative Multimodal Emotion Recognition Using Facial, Textual, and Speech Cues. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 1359–1367.

Lili Mou, Rui Men, Ge Li, Yan Xu, Lu Zhang, Rui Yan, and Zhi Jin. 2016. Natural Language Inference by Tree-Based Convolution and Heuristic Matching. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016*, Vol. 2. The Association for Computer Linguistics. https://doi.org/10.18653/v1/p16-2022

Arsha Nagrani, Samuel Albanie, and Andrew Zisserman. 2018. Seeing Voices and Hearing Faces: Cross-Modal Biometric Matching. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 8427–8436.

Jiquan Ngiam, Aditya Khosla, Mingyu Kim, Juhan Nam, Honglak Lee, and Andrew Y. Ng. 2011. Multimodal Deep Learning. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011*, Lise Getoor and Tobias Scheffer (Eds.). Omnipress, 689–696.

Youssef Oualil, Dietrich Klakow, György Szaszák, Ajay Srinivasamurthy, Hartmut Helmke, and Petr Motlícek. 2017. A context-aware speech recognition and understanding system for air traffic control domain. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2017*. IEEE, 404–408. https://doi.org/10.1109/ASRU.2017.8268964

John Edison Arevalo Ovalle, Thamar Solorio, Manuel Montes-y-Gómez, and Fabio A. González. 2017. Gated Multimodal Units for Information Fusion. In *5th International Conference on Learning Representations, ICLR 2017*. https://openreview.net/forum?id=S12_nquOe

José Manuel Pardo, Javier Ferreiros, Fernando Fernández Martínez, Valentín Sama Rojo, Ricardo de Córdoba, Javier Macías Guarasa, Juan Manuel Montero, Rubén San-Segundo-Hernández, Luis Fernando D'Haro, and Germán González. 2011. Automatic Understanding of ATC Speech: Study of Prospectives and Field Experiments for Several Controller Positions. *IEEE Trans. Aerosp. Electron. Syst.* 47, 4 (2011), 2709–2730. https://doi.org/10.1109/TAES.2011.6034660

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving Language Understanding by Generative Pre-Training.

Mirco Ravanelli and Yoshua Bengio. 2018a. Speaker Recognition from Raw Waveform with SincNet. In *2018 IEEE Spoken Language Technology Workshop (SLT)*. IEEE, Athens, Greece, 1021–1028. https://doi.org/10.1109/SLT.2018.8639585

Mirco Ravanelli and Yoshua Bengio. 2018b. Speech and Speaker Recognition from Raw Waveform with SincNet. *CoRR* abs/1812.05920 (2018). arXiv:1812.05920 http://arxiv.org/abs/1812.05920

M. Schulder, H. Helmke, Y. Oualil, J. Rataj, and D. Klakow. 2015. Assistant-Based Speech Recognition for ATM Applications. In *USA/Europe Air Traffic Management Research and Development Seminar*.

Karen Simonyan and Andrew Zisserman. 2014. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv preprint arXiv:1409.1556* (2014). arXiv:1409.1556

Lubos Smídl, Jan Svec, Daniel Tihelka, Jindrich Matousek, Jan Romportl, and Pavel Ircing. 2019. Air traffic control communication (ATCC) speech corpora and their use for ASR and TTS development. *Lang. Resour. Evaluation* 53, 3 (2019), 449–464. https://doi.org/10.1007/s10579-019-09449-5

David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur. 2018. X-Vectors: Robust DNN Embeddings for Speaker Recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, Calgary, AB, 5329–5333. https://doi.org/10.1109/ICASSP.2018.8461375

Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. How to Fine-Tune BERT for Text Classification?. In *Chinese Computational Linguistics - 18th China National Conference, CCL 2019 (Lecture Notes in Computer Science, Vol. 11856)*. Springer, Kunming, China, 194–206. https://doi.org/10.1007/978-3-030-32381-3_16

Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 2015. Improved Semantic Representations from Tree-Structured Long Short-Term Memory Networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015*, Vol. 1. 1556–1566. https://doi.org/10.3115/v1/p15-1150

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention Is All You Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*. 5998–6008.

Zuxuan Wu, Yu-Gang Jiang, Jun Wang, Jian Pu, and Xiangyang Xue. 2014. Exploring Inter-Feature and Inter-Class Relationships with Deep Neural Networks for Video Classification. In *Proceedings of the 22nd ACM International Conference on Multimedia*. 167–176.

Yong Xu, Qiuqiang Kong, Wenwu Wang, and Mark D. Plumbley. 2018. Large-Scale Weakly Supervised Audio Classification Using Gated Convolutional Neural Network. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 121–125. https://doi.org/10.1109/ICASSP.2018.8461975

Bo Yang, Xianlong Tan, Zhengmao Chen, Bing Wang, Min Ruan, Dan Li, Zhongping Yang, Xiping Wu, and Yi Lin. 2020. ATCSpeech: A Multilingual Pilot-Controller Speech Corpus from Real Air Traffic Control Environment. In *Interspeech 2020*. ISCA, Shanghai, China, 399–403. https://doi.org/10.21437/Interspeech.2020-1020

Yuni Zeng, Hua Mao, Dezhong Peng, and Zhang Yi. 2019. Spectrogram based multi-task audio classification. *Multim. Tools Appl.* 78, 3 (2019), 3705–3722. https://doi.org/10.1007/s11042-017-5539-3

Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. Character-level Convolutional Networks for Text Classification. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015*. 649–657.

Peng Zhou, Wei Shi, Jun Tian, Zhenyu Qi, Bingchen Li, Hongwei Hao, and Bo Xu. 2016. Attention-Based Bidirectional Long Short-Term Memory Networks for Relation Classification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016*, Vol. 2. The Association for Computer Linguistics, Berlin, Germany. https://doi.org/10.18653/v1/p16-2034

Pan Zhou, Wenwen Yang, Wei Chen, Yanfeng Wang, and Jia Jia. 2019. Modality Attention for End-to-End Audio-Visual Speech Recognition. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 6565–6569. https://doi.org/10.1109/ICASSP.2019.8683733

Juan Zuluaga-Gomez, Petr Motlícek, Qingran Zhan, Karel Veselý, and Rudolf A. Braun. 2020. Automatic Speech Recognition Benchmark for Air-Traffic Communications. In *Interspeech 2020, 21st Annual Conference of the International Speech Communication Association*. ISCA, 2297–2301. https://doi.org/10.21437/Interspeech.2020-2173