# Bridging the Gap between Audio and Text using Parallel-attention for User-defined Keyword Spotting

Youkyum Kim, Jaemin Jung, Jihwan Park, Byeong-Yeol Kim, and Joon Son Chung

*Abstract*—This paper proposes a novel user-defined keyword spotting framework that accurately detects audio keywords based on text enrollment. Since audio data possesses additional acoustic information compared to text, there are discrepancies between these two modalities. To address this challenge, we present ParallelKWS, which utilises self- and cross-attention in a parallel architecture to effectively capture information both within and across the two modalities. We further propose a phoneme duration-based alignment loss that enforces the sequential correspondence between audio and text features. Extensive experimental results demonstrate that our proposed method achieves state-of-the-art performance on several benchmark datasets in both seen and unseen domains, without incorporating extra data beyond the dataset used in previous studies.

*Index Terms*—attention mechanism, multi-modal fusion, user-defined keyword spotting

## I. INTRODUCTION

**K**EYWORD spotting (KWS) plays a crucial role as an entry point for initiating voice-activated services on smart devices, which have recently been in growing demand. Earlier KWS systems [1], [2], [3], [4], [5] based on deep learning primarily focused on detecting only pre-defined keywords. With the rapid advancement of artificial intelligence services and the need for enhanced user experience, there has been a shift towards user-defined keyword spotting (UDKWS) systems. These systems allow users to set their own keywords, broadening the scope and applicability of KWS.

Previous works [6], [7], [8], [9], [10] have predominantly concentrated on UDKWS systems where an audio sample is used for pre-enrolling the keyword, known as query-by-example (QbyE) methods. The performance of QbyE methods is highly variable, mainly due to discrepancies between the pre-enrolled audio and the input spoken utterance. In response to the disparities in audio samples and to enhance user convenience, UDKWS systems have incorporated a method for text-based keyword enrollment. However, text lacks acoustic information compared to audio, making it challenging to reduce the distinctions between these two modalities [11].

To address this issue, current research in UDKWS with text-based enrollment predominantly focuses on reducing the discrepancy between audio and text modalities. Establishing a phoneme-to-vector database by converting phonemes into

The first two authors contributed equally to this work.

Youkyum Kim, Jaemin Jung, and Joon Son Chung are with School of Electric Engineering, Korea Advanced Institute of Science and Technology, Daejeon 34141, Republic of Korea (e-mail:dbrua1998@kaist.ac.kr; jjm5811@kaist.ac.kr; joonson@kaist.ac.kr)

Jihwan Park and Byeong-Yeol Kim are with 42dot Inc., Seoul 06620, Republic of Korea (e-mail: jihwan.park@42dot.ai; byeongyeol.kim@42dot.ai)

averages of frame-level audio embeddings effectively reduces the mismatch between audio and text embedding spaces on in-domain datasets [12]. However, its performance in unseen domain datasets remains suboptimal. Other methods assess the similarity between audio and text embeddings using attention-based modules. Shin et al. [13] leverage audio embeddings as the key and value, and text embeddings as the query to the cross-attention module to evaluate the similarity between two modalities at the utterance level. Lee et al. [14] suggest a self-attention-based framework that merges audio and text embeddings into a singular representation.

To strengthen the correlation between two different modalities (audio and text), we present ParallelKWS, a UDKWS framework that adopts both self- and cross-attention mechanisms [15]. The effectiveness of modality fusion using both self- and cross-attention has been reported in various deep learning fields, including speech emotion recognition [16], [17], [18] and feature matching [19], [20], [21]. However, it has not yet been explored in the context of keyword spotting. The self-attention module captures both inter- and intra-modal information by processing concatenated audio and text embeddings as input [14]. To enrich the inter-modal information influenced by each respective modality, ParallelKWS also incorporates two cross-attention modules using audio and text embeddings as their queries, respectively.

Furthermore, we propose a phoneme duration-based alignment loss as an auxiliary training objective to obtain a fine-grained alignment between the embeddings from audio and text modalities. We employ a pre-trained speech embedder as a component of the audio encoder. As this embedder is trained with phoneme-level connectionist temporal classification (CTC) loss [22], the model inherently provides phoneme duration information of the audio samples at frame-level. A target matrix, generated from this duration information, is utilised to enforce sequential correspondence between audio and text. By aligning the phonetic timing of the spoken words with the corresponding textual representation, this approach improves how the model associates varied speech patterns with their textual counterparts. Experimental results show that our approach outperforms comparable previous works on most benchmark datasets, and demonstrate the effectiveness of our proposed framework.

## II. PROPOSED METHOD

In this section, we describe our proposed framework including model architecture and training objective. The overall framework is illustrated in Fig. 1. It comprises two distinct encoders that capture features from the audio and text
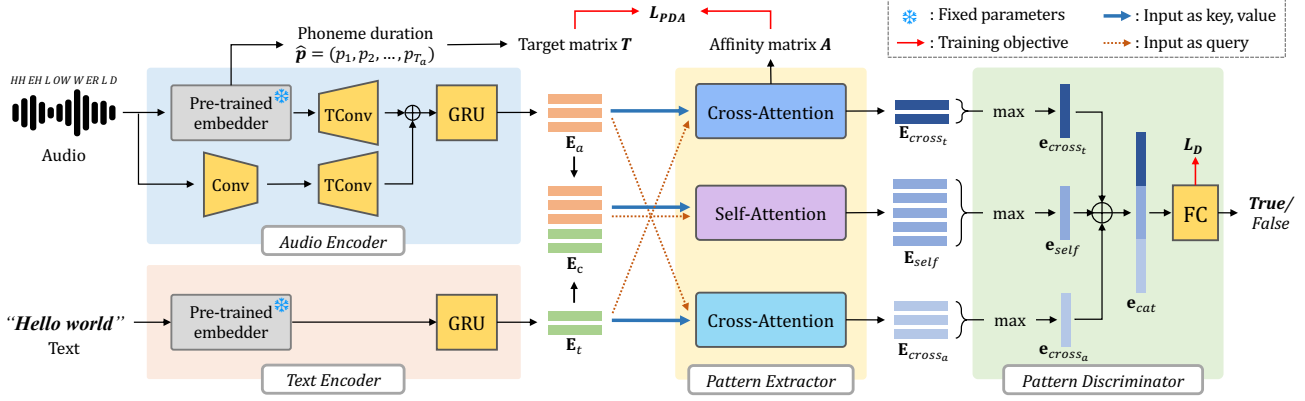
Fig. 1. Overall framework of ParallelKWS. "TConv" denotes "Transposed convolution". "max" and "FC" stand for "Max-pooling" and "Fully-connected layer," respectively. Audio embedding $\mathbf{E}_a$ is used as query in one cross-attention module, with text embedding $\mathbf{E}_t$ serving as key and value. In the other cross-attention module, this setup is reversed. Concatenated embedding $\mathbf{E}_c$ is input to self-attention module. Outputs from each attention module are max-pooled and then concatenated. Finally, the concatenated output is passed through FC to produce the final logit.

modalities. The framework also includes a pattern extractor, which combines these audio and text features, and a pattern discriminator responsible for determining the presence of the keyword.

### A. Model Architecture

**Audio encoder.** The audio encoder is composed of a dual-path feature extractor [14] followed by a single GRU [23] layer. One path of the feature extractor includes a pre-trained speech embedder, a small conformer [24] optimised using phoneme-level CTC loss, and a 1-D transposed convolution with a kernel size of 5 and a stride of 4. Following [12], the conformer is structured with 6 encoder layers, an encoder dimension of 144, a convolution kernel size of 3, and 4 attention heads. The other path consists of a 1-D convolution with a kernel size of 3 and stride of 2, followed by a 1-D transposed convolution with a kernel size of 3 and stride of 2. The outputs from both paths are concatenated along the feature dimension and then fed into a GRU layer, yielding the final audio embeddings. Audio embeddings are denoted as $\mathbf{E}_a \in \mathbb{R}^{T_a \times d}$, where $T_a$ and $d$ represent the lengths of the audio features and the dimension of the embeddings, respectively. $d$ is set to 128 in this study.

**Text encoder.** To reduce the mismatch with the output of the phoneme-based audio encoder, we use a pre-trained grapheme-to-phoneme (G2P) [25] model as a text encoder, followed by a single GRU layer. The G2P embeddings are derived from the last hidden states of the encoder [14]. The text embeddings are denoted as $\mathbf{E}_t \in \mathbb{R}^{T_t \times d}$, where $T_t$ refers to the lengths of the text features.

**Pattern extractor.** To effectively fuse audio and text information, we construct a pattern extractor using the Parallel-attention, which combines cross-attention and self-attention modules in parallel. The attention mechanism calculates the weighted sum of the values ($V$), based on the similarity scores between the queries ($Q$) and keys ($K$):

$$\text{Attn}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) V. \qquad (1)$$

$\mathbf{E}_{cross_t}$ is the output embedding of the cross-attention module, with the text embedding $\mathbf{E}_t$ as the query and the audio embedding $\mathbf{E}_a$ as the key and value. Conversely, $\mathbf{E}_{cross_a}$ is obtained from the other cross-attention module where the audio embedding is utilised as the query:

$$\mathbf{E}_{cross_t} = \text{Attn}(\mathbf{E}_t \mathbf{W}_t^Q, \mathbf{E}_a \mathbf{W}_a^K, \mathbf{E}_a \mathbf{W}_a^V), \qquad (2)$$

$$\mathbf{E}_{cross_a} = \text{Attn}(\mathbf{E}_a \mathbf{W}_a^Q, \mathbf{E}_t \mathbf{W}_t^K, \mathbf{E}_t \mathbf{W}_t^V). \qquad (3)$$

In the self-attention mechanism, the unimodal embeddings $\mathbf{E}_a$ and $\mathbf{E}_t$ are concatenated across the time dimension to form the concatenated embedding $\mathbf{E}_c$, which is utilized as the query, key, and value to obtain the self-attention output $\mathbf{E}_{self}$:

$$\mathbf{E}_{self} = \text{Attn}(\mathbf{E}_c \mathbf{W}_c^Q, \mathbf{E}_c \mathbf{W}_c^K, \mathbf{E}_c \mathbf{W}_c^V), \qquad (4)$$

where $\mathbf{W}^Q$, $\mathbf{W}^K$, and $\mathbf{W}^V$ are projection matrices of the query, key, and value, respectively.

**Pattern discriminator.** The pattern discriminator determines whether the keyword is detected. We first apply a max-pooling layer along the time axis to condense the outputs from both the cross- and self-attention modules. These condensed outputs are then concatenated along the feature axis to create the integrated features. Finally, we employ a fully-connected layer with a sigmoid activation. The process is summarised as follows:

$$\mathbf{e}_{cat} = \text{Concat}(\mathbf{e}_{cross_t}, \mathbf{e}_{cross_a}, \mathbf{e}_{self}) \qquad (5)$$

$$\hat{y} = \sigma(\mathbf{W} \cdot \mathbf{e}_{cat} + \mathbf{b}) \qquad (6)$$

where $\mathbf{e}_{cross_t}$, $\mathbf{e}_{cross_a}$, and $\mathbf{e}_{self}$ in $\mathbb{R}^d$ represent the condensed features of $\mathbf{E}_{cross_t}$, $\mathbf{E}_{cross_a}$, and $\mathbf{E}_{self}$, and $\mathbf{W}$, $\mathbf{b}$, $\sigma$ are the weights, biases, and the sigmoid function of the fully-connected layer, respectively.

### B. Training Objective

**Phoneme duration-based alignment loss.** We propose a novel training objective that enforces the model to learn the sequential correspondence between audio and text modalities based on phoneme duration information extracted from a

TABLE I
COMPARISON OF MODEL PERFORMANCES AND ABLATION STUDY ON THE TRAINING OBJECTIVE.
THE RESULTS FOR † ARE AS REPORTED IN PRIOR WORKS. THE BEST RESULTS ARE IN BOLD.

| Method | EER (%) ↓ | | | | AUC (%) ↑ | | | |
|---|---|---|---|---|---|---|---|---|
| | G | Q | LP$_E$ | LP$_H$ | G | Q | LP$_E$ | LP$_H$ |
| *Baselines* | | | | | | | | |
| CMCD [13]† | 27.25 | 12.15 | 8.42 | 32.90 | 81.06 | 94.51 | 96.70 | 73.58 |
| FlexiKWS [12]† | 14.05 | - | 0.8 | 18.4 | 93.16 | - | 99.94 | 89.2 |
| PhonMatchNet (re-impl.) [14] | 14.04 ±1.08 | 11.72 ±1.28 | 0.48 ±0.03 | 18.77 ±0.22 | 93.81 ±0.83 | 95.50 ±0.84 | 99.80 ±0.01 | 88.01 ±0.16 |
| *Baselines (with additional data)* | | | | | | | | |
| FlexiKWS w/ neg. [12]† | 13.45 | - | 1.7 | 14.4 | 93.94 | - | 99.84 | **92.7** |
| PhonMatchNet [14]† | **6.77** | 4.75 | 2.80 | 18.82 | **98.11** | 98.90 | 99.29 | 88.52 |
| *ParallelKWS (ours)* | | | | | | | | |
| Only detection loss | 8.78 ±0.27 | 2.90 ±0.11 | 0.11 ±0.01 | 14.80 ±0.10 | 97.31 ±0.18 | 99.66 ±0.03 | 99.96 ±0.00 | 91.29 ±0.10 |
| Detection loss + MM loss | 8.02 ±0.25 | 3.38 ±0.83 | 0.13 ±0.01 | 15.46 ±0.13 | 97.72 ±0.16 | 99.50 ±0.24 | 99.95 ±0.01 | 90.72 ±0.15 |
| Detection loss + PDA loss | 7.78 ±0.40 | **2.61** ±0.18 | **0.09** ±0.01 | **14.36** ±0.07 | 97.75 ±0.13 | **99.67** ±0.05 | **99.97** ±0.01 | 91.68 ±0.10 |

pre-trained speech embedder. Inspired by [13], we align the sequential information from audio and text embeddings by matching the affinity matrix with the phoneme duration-based target matrix. Here, the attention map from the cross-attention module with the query of text embeddings is used as the affinity matrix. For negative pairs, we utilise a target matrix derived from normally distributed random noise, following [13].

For a given positive audio-text pair, $\hat{\boldsymbol{p}} = (p_1, p_2, ..., p_{T_a})$ represents a vector of phoneme predictions from the pre-trained speech embedder. As these predictions are at the frame level, the number of consecutive identical phoneme predictions likely contains information about the phoneme duration of the audio sample. We assign a group index to each $p_i$ in $\hat{\boldsymbol{p}}$, incrementing the index whenever a new phoneme prediction appears, thus grouping consecutive identical phoneme predictions. The resulting consecutive index vector can be denoted as $\hat{\boldsymbol{c}} = (c_1, c_2, ..., c_{T_a})$, where $c_i$ is defined as follows:

$$c_i = \begin{cases} 1, & \text{if } i = 1 \\ c_{i-1}, & \text{if } p_i = p_{i-1} \\ c_{i-1} + 1, & \text{if } p_i \neq p_{i-1}. \end{cases} \quad (7)$$

Using the above index vector, we define matrix $\boldsymbol{D} = [d_{ij}] \in \mathbb{R}^{T_a \times T_t}$, where $d_{ij} = j - c_i$ for all $i, j \in \mathbb{N}$, $1 \leq i \leq T_a$, and $1 \leq j \leq T_t$. The phoneme duration-based target matrix $\boldsymbol{T} = [t_{ij}] \in \mathbb{R}^{T_a \times T_t}$ is obtained through the following equation.

$$x_{ij} = -\frac{(d_{ij}/T_t)^2}{2g^2}. \quad (8)$$

$$t_{ij} = \frac{\exp(x_{ij})}{\sum_i \exp(x_{ij})}. \quad (9)$$

Here, $g$ is a hyperparameter that determines the gradient of the exponential function and is set to $0.1$ in this work. The phoneme duration-based alignment loss is defined as the mean square error between the affinity matrix $\boldsymbol{A}$ and the phoneme duration-based target matrix $\boldsymbol{T}$:

$$L_{PDA} = ||\boldsymbol{A} - \boldsymbol{T}||^2. \quad (10)$$

**Detection loss.** To assess if the input audio sample and the input text correspond to the same keyword, we use binary cross-entropy loss on the logits from the pattern discriminator.

Since this detection loss addresses both the entire audio sample and the phonemes within it, the network is trained to recognise similarities between audio and text at the level of entire utterances.

$$L_D = -(y \cdot \log \hat{y} + (1 - y) \cdot \log(1 - \hat{y})), \quad (11)$$

where $\hat{y}$ and $y$ denote the predicted probability and the ground truth label, respectively.

Finally, we formulate the overall loss ($L_{total}$) as follows:

$$L_{total} = \lambda \cdot L_{PDA} + L_D, \quad (12)$$

where $\lambda$ is a weight factor, and is set to 0.3.

## III. EXPERIMENTS

### A. Datasets and Evaluation Methods

We employ the LibriPhrase [13] dataset, which comprises phrases ranging from 1 to 4 words, and is divided into a training set and a test set, derived from distinct splits of the LibriSpeech [26] dataset: *train-clean* and *train-other*. We use 800k phrases for training, evenly distributed with 200k phrases for each word length, in line with [12], [13], [27]. Additionally, we use LibriSpeech *train-clean* dataset along with LibriPhrase training set to train the conformer with phoneme-level CTC loss. This training involves an initial phase on the LibriSpeech *train-clean* dataset, followed by fine-tuning using shorter audio segments from the LibriPhrase training set. Input audio data augmentation is performed using various noises from the MUSAN [28] dataset and room impulse response filters. The entire network is then trained on the LibriPhrase training set, without updating the parameters in the pre-trained conformer and G2P model. Audio features are extracted using 80-channel filterbanks with a 25ms window and a 10ms frame shift.

The LibriPhrase test set is categorised based on the Levenshtein distance [29] between negative pairs, where a lower distance indicates higher phonetic similarity and greater difficulty in discrimination. The test set with hard negative pairs and easy negative pairs are labeled as LibriPhrase-hard (**LP$_H$**) and LibriPhrase-easy (**LP$_E$**), respectively. For evaluation, four distinct KWS benchmark datasets are used: LibriPhrase-easy, LibriPhrase-hard, Google Speech Commands V1 (**G**) [30],

TABLE II
EFFECTIVENESS OF THE ATTENTION MODULES IN PATTERN EXTRACTOR.
THE NUMBER OF PARAMETERS (# PARAMS.) ONLY REFLECTS THE
TRAINABLE PARAMETERS. PRE-TRAINED EMBEDDERS CONTAIN 2.33M
PARAMETERS. THE BEST RESULTS ARE IN BOLD.

| Model | | | EER (%) ↓ | | | |
|---|---|---|---|---|---|---|
| Self | Cross | # params. | G | Q | $LP_E$ | $LP_H$ |
| ✓ | | $0.55M$ | $9.38 \pm 0.32$ | $4.81 \pm 1.57$ | $0.23 \pm 0.03$ | $18.56 \pm 0.30$ |
| | ✓ | $0.61M$ | $9.32 \pm 0.21$ | $3.24 \pm 0.63$ | $0.16 \pm 0.02$ | $15.80 \pm 0.16$ |
| ✓ | ✓ | $0.68M$ | $\mathbf{8.78} \pm 0.27$ | $\mathbf{2.90} \pm 0.11$ | $\mathbf{0.11} \pm 0.01$ | $\mathbf{14.80} \pm 0.10$ |

and Qualcomm Keyword Speech (**Q**) [31], with the official split for $LP_E$ and $LP_H$ as provided in [13]. For **G** and **Q**, we adhere to the testing protocol in [14] to maintain fairness in comparison, considering all keywords except the anchor keyword as negatives. We report the Equal Error Rate (EER) and Area Under the ROC Curve (AUC) scores for each benchmark dataset. We present the average performance and standard deviation across three experiments, each conducted with a distinct random seed for reliability.

### B. Implementation Details

The network is optimised for 100 epochs using the Adam optimizer [32], set to a fixed learning rate of 1e-3. For evaluation, we select the model with the lowest EER on the test sets. We establish the batch size at 2048, and the training process takes approximately one day on a single A5000 GPU which has a memory size of 24GB. The framework for our model is implemented using the PyTorch library.

## IV. RESULTS

### A. Comparison with Baselines

In Table I, we report the performance of our proposed framework, ParallelKWS, alongside that of existing baselines. CMCD [13] and FlexiKWS [12] utilise the same training dataset as our study. However, PhonMatchNet [14] employs a KWS model pre-trained on various external domain data ($200M$ audio clips) collected from YouTube [33] as a speech embedder. To ensure a fair comparison, we also present the performance of PhonMatchNet re-implemented using the same speech embedder as ours. When trained on the same dataset, ParallelKWS outperforms all existing baselines in terms of both EER and AUC scores. Notably, our method significantly improves performance on the $LP_E$ dataset by 99.0%, 89.2%, and 82.0% over CMCD, FlexiKWS, and the re-implemented PhonMatchNet, respectively.

We also compare our model with those trained using additional datasets. FlexiKWS uses phonetically confusable keywords as additional training data, labeled as *FlexiKWS w/ neg.* in Table I. Nevertheless, ParallelKWS demonstrates improved performance on all test sets, except for a slight (1.1%) decrease in AUC score on the $LP_H$ dataset. Compared to PhonMatchNet, which includes a speech embedder pre-trained on large-scale external data and employs phoneme-level detection loss to enhance the capability of distinguishing similar pronunciations, ParallelKWS shows improved performance on the datasets except for **G**. Especially on the $LP_H$ dataset, our

framework demonstrates a significant improvement of 23.7%. Through the comparison of performance with the baselines, we highlight the effectiveness of ParallelKWS in capturing data dependencies from unseen domains (**G** and **Q**) without additional training processes, while maintaining its capability with in-domain data ($LP_E$ and $LP_H$).

### B. Ablation Study

**Effectiveness of parallel-attention architecture.** We demonstrate the impact of parallel-attention architecture through ablation studies on the attention modules within the pattern extractor. As shown in Table II, using either self-attention or cross-attention results in performance degradations across all test sets compared to their parallel connection. Notably, parallel-attention significantly improves performance in the $LP_H$ dataset, which contains marginally distinguishable pronunciations, as well as in the test sets from unseen domains, **G** and **Q**. These results indicate that the parallel-attention architecture precisely captures both inter- and intra-modal information from audio and text, effectively merging the two modalities.

**Effectiveness of phoneme duration-based alignment loss.** We conduct an ablation study to assess the functionality of the proposed phoneme duration-based alignment loss (PDA loss). This study aims to confirm the effectiveness of incorporating phoneme duration information. We compare our approach with the monotonic matching loss (MM loss) proposed in [13], applying it to our architecture. The monotonic matching approach aligns the affinity matrix with a target matrix that is organised in a monotonic order across the audio and text sequences. However, this target matrix lacks intrinsic duration-related information.

As indicated in the seventh row of Table I, using MM loss leads to decreased performance in the **Q**, $LP_E$, and $LP_H$ datasets, while it shows improvement in the **G** dataset. In contrast, using PDA loss results in an EER reduction of 11.4%, 10.0%, 21.2%, and 3.0% in the **G**, **Q**, $LP_E$, and $LP_H$ datasets, respectively. This result emphasises that aligning audio and text sequences in the affinity matrix based on phoneme duration encourages the preceding modules to be trained with an enhanced capability to capture duration-related sequential information.

## V. CONCLUSION

In this paper, we introduce a framework for user-defined keyword spotting with text-based enrollment, effectively integrating audio and text information. Our framework employs cross- and self-attention modules in a parallel architecture to capture both inter- and intra-modal information, thus improving the capability of the model to fuse audio and text modalities. We also implement an alignment loss that utilises phoneme duration information, derived from the pre-trained speech embedder, to enhance the alignment of sequential information between audio and text embeddings. Experimental results show that our framework outperforms previous models on most test sets, achieving this without any training on data from external domains.

REFERENCES

[1] G. Chen, C. Parada, and G. Heigold, "Small-footprint keyword spotting using deep neural networks," in *Proc. ICASSP*, 2014.

[2] J. Hou, Y. Shi, M. Ostendorf, M.-Y. Hwang, and L. Xie, "Region proposal network based small-footprint keyword spotting," *IEEE Signal Processing Letters*, 2019.

[3] R. Tang and J. Lin, "Deep residual learning for small-footprint keyword spotting," in *Proc. ICASSP*, 2018.

[4] A. Berg, M. O'Connor, and M. T. Cruz, "Keyword transformer: A self-attention model for keyword spotting," in *Proc. Interspeech*, 2021.

[5] J. Huh, M. Lee, H. Heo, S. Mun, and J. S. Chung, "Metric learning for keyword spotting," in *IEEE Spoken Language Technology workshop*, 2021.

[6] G. Chen, C. Parada, and T. N. Sainath, "Query-by-example keyword spotting using long short-term memory networks," in *Proc. ICASSP*, 2015.

[7] J. Huang, W. Gharbieh, H. S. Shim, and E. Kim, "Query-by-example keyword spotting system using multi-head attention and soft-triple loss," in *Proc. ICASSP*, 2021.

[8] M. Mazumder, C. Banbury, J. Meyer, P. Warden, and V. J. Reddi, "Few-Shot Keyword Spotting in Any Language," in *Proc. Interspeech*, 2021.

[9] A. Parnami and M. Lee, "Few-shot keyword spotting with prototypical networks," in *International Conference on Machine Learning Technologies*, 2022.

[10] J. Jung, Y. Kim, J. Park, Y. Lim, B.-Y. Kim, Y. Jang, and J. S. Chung, "Metric learning for user-defined keyword spotting," in *Proc. ICASSP*, 2023.

[11] X. Tan, T. Qin, F. Soong, and T.-Y. Liu, "A survey on neural speech synthesis," *arXiv preprint arXiv:2106.15561*, 2021.

[12] K. Nishu, M. Cho, P. Dixon, and D. Naik, "Flexible keyword spotting based on homogeneous audio-text embedding," in *Proc. ICASSP*, 2023.

[13] H.-K. Shin, H. Han, D. Kim, S.-W. Chung, and H.-G. Kang, "Learning audio-text agreement for open-vocabulary keyword spotting," in *Proc. Interspeech*, 2022.

[14] Y.-H. Lee and N. Cho, "Phonmatchnet: phoneme-guided zero-shot keyword spotting for user-defined keywords," in *Proc. Interspeech*, 2023.

[15] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Proc. NeurIPS*, 2017.

[16] L. Sun, B. Liu, J. Tao, and Z. Lian, "Multimodal cross-and self-attention network for speech emotion recognition," in *Proc. ICASSP*, 2021.

[17] D. Yang, S. Huang, Y. Liu, and L. Zhang, "Contextual and cross-modal interaction for multi-modal speech emotion recognition," *IEEE Signal Processing Letters*, vol. 29, pp. 2093–2097, 2022.

[18] M. Luo, H. Phan, and J. Reiss, "cross-modal fusion techniques for utterance-level emotion recognition from text and speech," in *Proc. ICASSP*, 2023.

[19] P.-E. Sarlin, D. DeTone, T. Malisiewicz, and A. Rabinovich, "Superglue: Learning feature matching with graph neural networks," in *Proc. CVPR*, 2020.

[20] Q. Wang, J. Zhang, K. Yang, K. Peng, and R. Stiefelhagen, "Match-former: Interleaving attention in transformers for feature matching," in *Proc. ACCV*, 2022.

[21] X. Lu, Y. Yan, B. Kang, and S. Du, "Paraformer: Parallel attention transformer for efficient feature matching," in *Proc. AAAI*, 2023.

[22] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber, "Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks," in *Proc. ICML*, 2006.

[23] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," in *Proc. EMNLP*, 2014.

[24] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu *et al.*, "Conformer: Convolution-augmented transformer for speech recognition," in *Proc. Interspeech*, 2020.

[25] K. Park and J. Kim, "g2pe," https://github.com/Kyubyong/g2p, 2019.

[26] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *Proc. ICASSP*, 2015.

[27] K. Nishu, M. Cho, and D. Naik, "Matching Latent Encoding for Audio-Text based Keyword Spotting," in *Proc. Interspeech*, 2023.

[28] D. Snyder, G. Chen, and D. Povey, "Musan: A music, speech, and noise corpus," *arXiv preprint arXiv:1510.08484*, 2015.

[29] V. I. Levenshtein *et al.*, "Binary codes capable of correcting deletions, insertions, and reversals," in *Soviet physics doklady*, 1966.

[30] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *arXiv preprint arXiv:1804.03209*, 2018.

[31] B. Kim, M. Lee, J. Lee, Y. Kim, and K. Hwang, "Query-by-example on-device keyword spotting," in *IEEE Automatic Speech Recognition and Understanding Workshop*, 2019.

[32] D. P. Kingma and J. L. Ba, "Adam: A method for stochastic optimization," in *Proc. ICLR*, 2014.

[33] J. Lin, K. Kilgour, D. Roblek, and M. Sharifi, "Training keyword spotters with limited and synthesized speech data," in *Proc. ICASSP*, 2020.