

U2-KWS: UNIFIED TWO-PASS OPEN-VOCABULARY KEYWORD SPOTTING WITH KEYWORD BIAS

Ao Zhang¹, Pan Zhou², Kaixun Huang¹, Yong Zou², Ming Liu², Lei Xie^{1*}

¹Audio, Speech and Language Processing Group (ASLP@NPU), School of Computer Science, Northwestern Polytechnical University, Xian, China

²Space AI, Li Auto

ABSTRACT

Open-vocabulary keyword spotting (KWS), which allows users to customize keywords, has attracted increasingly more interest. However, existing methods based on acoustic models and post-processing train the acoustic model with ASR training criteria to model all phonemes, making the acoustic model under-optimized for the KWS task. To solve this problem, we propose a novel unified two-pass open-vocabulary KWS (U2-KWS) framework inspired by the two-pass ASR model U2. Specifically, we employ the CTC branch as the first stage model to detect potential keyword candidates and the decoder branch as the second stage model to validate candidates. In order to enhance any customized keywords, we redesign the U2 training procedure for U2-KWS and add keyword information by audio and text cross-attention into both branches. We perform experiments on our internal dataset and Aishell-1. The results show that U2-KWS can achieve a significant relative wake-up rate improvement of 41% compared to the traditional customized KWS systems when the false alarm rate is fixed to 0.5 times per hour.

Index Terms— Open-vocabulary keyword spotting, U2-KWS, customized keyword bias, multi-task learning

1. INTRODUCTION

Keyword spotting (KWS) is the task of detecting predefined keywords from a consecutive audio stream, which is widely applied in edge devices with a speech interface. Traditional KWS approaches are based on keyword/filler models which involve modeling keyword and non-keyword segments with hidden Markov model (HMM) [1, 2]. With the advance of deep learning, many works are proposed to build systems based on a single neural network without HMM and directly predict the keyword or sub-word tokens of the keyword [3, 4, 5, 6, 7]. Nowadays, the majority of keyword spotting systems primarily rely on such deep learning approaches. While these systems can achieve high precision on the pre-defined keywords, they require a large amount of specific training data

that contains those pre-defined keywords. Additionally, keywords cannot be changed after training, which requires the user to remember specific keywords. To deal with such limitations, *open-vocabulary* keyword spotting, which allows users to customize keywords, has gained popularity in recent years.

Many previous works for open-vocabulary KWS employ query-by-example (QbyE) methods, which use only audio signals as input [8, 9]. QbyE methods enroll reference keyword speech and compare it with new input speech queries. These methods require users to record the speech of keywords during customization and their performance is influenced apparently by the reference keyword speech. In contrast, customizing keywords based on text makes it easier to set keywords and is stable for different speakers. A common text-based open-vocabulary KWS method usually adopts an acoustic model from ASR to transform speech into phonetic posteriorgrams and then leverage a post-processing technique such as HMM to predict keyword existence [10, 11, 12, 13]. In practice, the system with a single streaming acoustic model produces a large number of false alarms [14, 15]. Therefore, the multi-stage strategy has been previously adopted to reduce false alarms [14, 16, 17, 18]. In general, the first stage is a lightweight always-on keyword detector. Once a keyword candidate is detected, the corresponding audio segment is sent to the following stage(s) for further verification. Multi-stage systems have several modules and complex structures, making them hard to build. Recently, there has been increasing interest in unifying multi-stage modules into one single model. In this direction, Cascaded Transducer-Transformer (CATT-KWS) uses two-pass models, which unify streaming and non-streaming ASR approaches [19, 20], to unify multi-stage KWS into one model [21]. Specifically, it uses the streaming part, which is originally used to generate streaming hypotheses, as the first-stage model to detect possible keywords, and then uses the non-streaming parts, which are originally used to re-score streaming hypotheses, as the validation stages for further verification of keywords detected in the first stage. Although this approach changes the inference of two-pass models from the original recognition & rescoring to the multi-stage KWS, the two-pass model is still trained

*: Corresponding author.

according to the criterion of ASR, and the model does not explicitly utilize the information of the keywords. The model without knowledge of keywords is optimized for the accuracy of all words, thus resulting in the under-optimization problem for the KWS task.

Therefore, it is important to particularly integrate the keyword information into the acoustic model to make it more sensitive to the user-defined keywords. Recent works have employed cross-attention techniques [22] to incorporate keyword representation in the context of acoustic representation. These approaches can be categorized into two categories based on the type of query employed for cross-attention: acoustic-query attention and keyword-query attention. The acoustic-query attention employs the acoustic representation as the query, utilizing the keyword information to bias the acoustic model towards the user-defined keywords [23, 14, 15, 24]. This attention method can process audio frame by frame, and the output of the attention mechanism is used to predict phonetic posteriorgrams. Because the information in one frame is relatively limited, such attention cannot adequately model the relationship between the entire audio and the keyword. In contrast, the keyword-query attention employs the keyword representation as the query and performs cross-attention with the entire audio. This method effectively models the dependencies between keywords and audio and the attention output represents the temporal correlation patterns between acoustic representations and keywords, which can be used for direct keyword prediction [25, 26, 27, 28]. This approach queries entire audio and thus better models global relationships. However, due to its global attention mechanism, it is not suitable for streaming inference.

In this paper, we propose the unified two-pass open-vocabulary keyword spotting framework (U2-KWS). Unlike the previously mentioned multi-stage KWS model, our approach focuses on improving the two-pass model specifically for the KWS task rather than relying on an acoustic model trained with the criterion of ASR. First, we employ attention-based keyword bias methods to bias the model towards user-defined keywords. Unlike previous works that use single-type attention, we apply the acoustic-query attention in the streaming branch and the keyword-query attention in the decoder branch within our framework. This allows us to leverage the strengths of each method, considering both streaming capability and accuracy. Besides, we redesign the training and inference of the two-pass model. We only feed the attention decoder with the keywords and corresponding acoustic representations, thus training it as a specific keyword rescoring module focusing on the task of keyword verification. We simulate the process of customizing keywords by sampling keywords from transcripts to generate positive and negative samples during training. We evaluate the proposed U2-KWS framework on the AISHELL-1 dataset [29] and an internal in-car speech communication dataset. Results

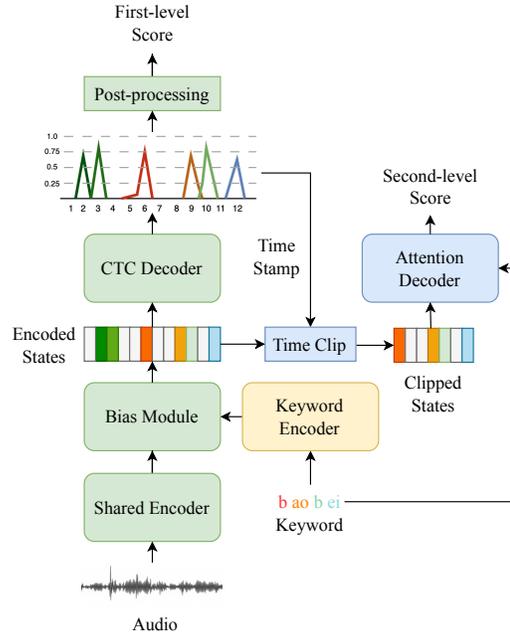


Fig. 1. Block diagram of the proposed U2-KWS.

on the proposed approach demonstrate an impressively high wake-up rate and low false alarm rate.

2. U2-KWS

In this section, we provide a detailed introduction to our proposed U2-KWS framework. First, we provide an overview of the entire architecture of the two-pass keyword spotting framework. Next, we discuss the use of keyword-bias methods in both streaming and non-streaming branches. Finally, we present the newly designed training and inference procedures for the two-pass keyword spotting.

2.1. Model Architecture

The proposed model architecture is shown in Fig. 1. It contains four parts: a shared encoder, a bias module, a CTC decoder, and an attention decoder. The green part of the figure is the streaming CTC branch in the two-pass model, which is employed as the first stage model for detecting potential keyword candidates in the audio, while the blue part is the non-streaming branch in the two-pass model, which works as the second-stage model by verifying the detected candidates with an attention decoder to reduce false alarms. The yellow part represents the keyword encoder which models relationships within keywords and encodes keywords into high-level representations.

The shared encoder consists of multiple Conformer [30] layers and is trained with the dynamic chunk strategy [31] which enables latency control and can forward in different chunk sizes. The keyword encoder is a single LSTM. The bias module consists of a multi-head attention for model bias and

a linear layer for dimensional transformation. The CTC decoder consists of a linear layer and a log-softmax layer, which outputs the phonetic posteriorgrams. The attention decoder is a Transformer decoder [22] but gets keywords instead of the hypotheses as input and does cross-attention with acoustic representations corresponding to the keyword phrase.

2.2. Customized Keyword Bias

In this section, we will provide a detailed explanation of the keyword bias methods based on acoustic-query attention and keyword-query attention. Afterward, we will discuss how we unify these methods into our two-pass model.

2.2.1. Attention-based Keyword Bias

The cross-attention mechanism is a widely used and efficient approach for modeling inter-modal relationships. In this paper, we formalize the cross-attention mechanism as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (1)$$

Here, Q represents the query vector, K represents the key vector, and V represents the value vector. We will now discuss the attention-based keyword bias method.

The acoustic-query attention method employs the acoustic representation h_a as the query and the keyword representation h_k as the key and value, which can be formulated as:

$$h_a \xrightarrow{\text{linear}} Q_a; h_k \xrightarrow{\text{linear}} K_k, V_k, \quad (2)$$

$$\tilde{h}_a = \text{Attention}(Q_a, K_k, V_k). \quad (3)$$

The resulting attention output \tilde{h}_a retains the same dimension as the acoustic representation and can be combined with h_a for the prediction of phonetic posteriorgrams.

The keyword-query attention method utilizes the keyword representation h_k as the query and the acoustic representation h_a as the key and value. The resulting attention output \tilde{h}_k represents the agreement between the audio and text [27]. This process can be described as follows:

$$h_k \xrightarrow{\text{linear}} Q_k; h_a \xrightarrow{\text{linear}} K_a, V_a, \quad (4)$$

$$\tilde{h}_k = \text{Attention}(Q_k, K_a, V_a). \quad (5)$$

2.2.2. Two-pass Keyword Bias

The applications of acoustic-query attention and keyword-query attention in our framework are shown in Fig. 2. To incorporate the acoustic-query attention, we integrate it into the streaming branch of the two-pass model. Firstly, we encode the customized keyword into an embedding using the keyword encoder. Next, we utilize the acoustic-query attention in the bias module to integrate the keyword information into the acoustic representation encoded by the shared encoder. The output of the attention module is concatenated with the

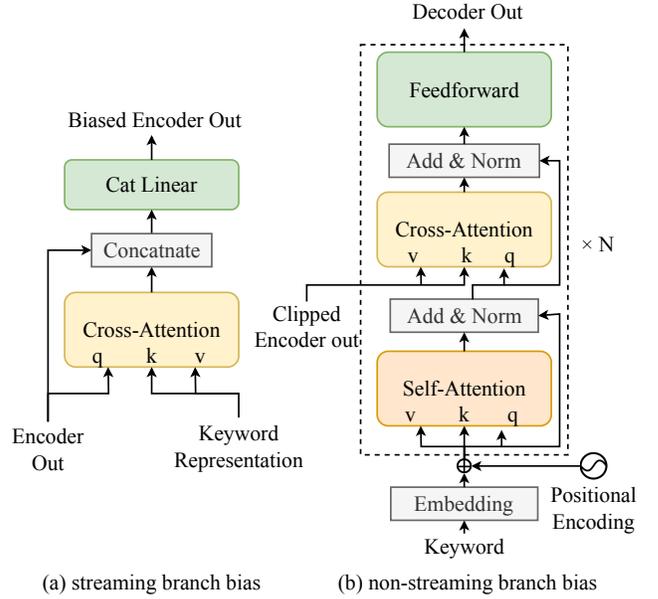


Fig. 2. Block diagram of keyword bias in U2-KWS.

original acoustic representation and then projected back to the original dimension. The resulting biased encoder output includes both the acoustic and keyword information and is used for subsequent inference.

For the keyword-query attention, we use it in the modified transformer decoder to bias the non-streaming branch. We keep the original transformer decoder structure and fix its input as keywords instead of the decoding results of the streaming branch, thus making the decoder a specialized keyword re-scoring module. In this case, the self-attention module encodes keywords into high-level representations and the cross-attention is the keyword-query attention. The output of the decoder is the token-level posterior probability, and we calculate the possibility of the keyword path as the score of the second stage model. This approach more adequately models the correlation between keywords and speech but is also more resource-intensive so we only used this structure for the non-streaming branch.

2.3. Training and Inference

2.3.1. Training of First Stage

We train the streaming branch first. The streaming branch is optimized with CTC loss, and since we introduced the keyword-bias module to generate keyword samples during training. To generate positive and negative samples, we employ a sampling technique on the transcripts to simulate user-defined keywords. Specifically, for a given utterance, we randomly select a consecutive word sub-sequence from its transcript as the keyword to create a positive sample. For generating a negative sample, we randomly combine words from the lexicon that are not present in the transcribed text of

the utterance. To introduce a training error for keyword spotting and encourage the model to pay attention to the keyword information, we add a $\langle\text{eok}\rangle$ token after the sampled keyword, which is similar to the label-augment method proposed for context bias in CLAS [32]. To provide a clearer understanding of this process, we take the sample ‘‘Call you Jarvis’’ as an example. The positive and negative samples produced for this sample are illustrated in Table 1. To be concise, we employ words as an instance for illustrating purposes, as we actually convert words into phonemes during training and inference.

Table 1. Example of keyword sampling

Class	Positive	Negative
Text	Call you Jarvis	
Keyword	Jarivs	Alex
Keyword	Jarvis $\langle\text{eok}\rangle$	Alex $\langle\text{eok}\rangle$
Encoder Input	Jarvis $\langle\text{eok}\rangle$	Alex $\langle\text{eok}\rangle$
CTC Target	Call you Jarvis $\langle\text{eok}\rangle$	Call you Jarvis
Decoder Input	$\langle\text{sos}\rangle$ Jarvis	$\langle\text{sos}\rangle$ Alex
Decoder Target	Jarvis $\langle\text{eok}\rangle$	$\langle\text{eos}\rangle\langle\text{eos}\rangle$

2.3.2. Training of Second Stage

After the training of the streaming branch, we jointly optimize the whole two-pass model with attention loss and CTC loss. In the training phase, we follow the same method to generate positive and negative samples as before. The overall forward and loss calculations of decoder training remain the same with the attention decoder in the ASR task, but we make modifications to the input token and target. Specifically, we configure the input of the decoder to be a fixed sequence consisting of ‘‘ $\langle\text{sos}\rangle$ + keyword’’. For positive samples, the target is set as ‘‘keyword + $\langle\text{eok}\rangle$ ’’, while for negative samples, the model directly predicts $\langle\text{eos}\rangle$, as shown in Table 1. In this way, the decoder focuses on the prediction of keyword sequence during training and becomes a token-level classifier with keyword discrimination capability.

Inspired by the end-to-end segmentation method used in the two-pass ASR model [33], we use timestamp information from the streaming branch to clip the encoder output for the decoder branch. This allows the cross attention of the decoder to only focus on the acoustic representation that contains the keyword. By doing so, we can reduce the complexity of the task for the decoder and improve keyword spotting performance.

To achieve the clip of the encoder output, we require the start and end times of the keyword. With the $\langle\text{eok}\rangle$ token in the first-level model, we can identify the frame with the highest posterior probability of $\langle\text{eok}\rangle$ as the end frame. To esti-

mate the start frame, we use a method similar to the spike trigger CTC [34], where we consider a spike as any non-blank token with a posterior probability exceeding a pre-set threshold. We count the number of spikes from the end frame back to the start of the speech, and when the number of spikes exceeds the length of the keyword sequence, we consider the last counted spike as the starting frame. This method doesn’t need extra alignment information and can output stable lengths compared with the alignment method based on CTC posterior probability.

The joint loss function for the entire two-pass model can be expressed as:

$$L = \lambda L_{ctc} + (1 - \lambda)L_{att}, \quad (6)$$

where L_{ctc} and L_{att} represent the CTC loss and the attention loss, respectively. The hyper-parameter λ controls the contribution of each loss to the overall training objective.

2.3.3. Inference

In the initialization phase, the keyword encoder generates a high-level representation of the customized keyword entered by the user, which is marked in yellow in Figure 1. The first stage model, highlighted in green in the same figure, processes the input audio stream chunk by chunk to get the posterior probability of CTC. We search the keyword path on the post probability to get the path with the highest probability and use the probability as the score of the first stage. If the score surpasses the threshold set for the first stage, we clip the acoustic representation based on the $\langle\text{eok}\rangle$ token and non-blank token spikes of CTC. Then, we employ the decoder to conduct keyword rescoring by sequentially feeding the keyword sequence token into the decoder and computing the probability of the next token of the keyword. This score is compared to the threshold for the final prediction. Since our encoder is trained with dynamic chunks, the encoder can infer in full-chunk mode to fully utilize the context. We can re-feed audio segments containing keyword candidates detected in the first stage to the encoder again in full-chunk mode to get more informative acoustic representations for the rescoring of the decoder in the second stage. Thus we further improve performance by only introducing minimal latency. We call the direct use of streaming encoder output the causal decoder, and the re-feeding method the full decoder.

3. EXPERIMENTS

In this section, we introduce the corpus and describe the experimental setup including our model configuration. Experimental results and analysis are also presented at last.

3.1. Corpus

Internal Corpus: Models are trained on a Mandarin ASR corpus comprising 1,000 hours of speech collected in hybrid

electric vehicles. We randomly shuffled the development set from the training set. To evaluate the accuracy of our models, we record the positive test set in the same environment as the training set. This test set contains 30 different keywords, each including 200 samples, for a total of 3,000 positive samples. These keywords are composed of two to four Chinese characters. To evaluate false alarms, we used a separate audio set of 60 hours, which mainly contained chat, command, and radio broadcasts, as the negative test set. All the data mentioned above, including the training set, development set, positive test set, and negative test set, are anonymous and hand-transcribed.

AISHELL-1: We also conduct experiments on the public Chinese Mandarin speech corpus AISHELL-1 [29]. For training and validation, we use the train set and dev set of AISHELL-1. For evaluation, we construct a test set consisting of 7,176 pairs of positive and negative samples from the AISHELL-1 test set using the keyword sampling method adopted in training.

3.2. Configuration

Our model used 80-dimensional log Mel-filter banks with a 25ms window and a 10ms shift. SpecAugment [35] is applied 2 frequency masks with maximum frequency mask ($F = 10$), and 2 time masks with maximum time mask ($T = 50$) to alleviate over-fitting. Two convolution sub-sampling layers with kernel size 3×3 and stride 2 are used in the front of the encoder.

The baseline CTC model we used consists of a 12-layer conformer with 128-dimensional input, 4 self-attention heads, and 256 linear units. The larger baseline expands the input dim to 256 and the number of linear units to 2048 which makes it nine times larger than the baseline model. For the keyword bias, the keyword encoder comprises an embedding layer and one layer LSTM with a dimensionality of 128. The bias module is a multi-head attention layer of 128 hidden states and 4 attention heads. For the decoder branch, we employ 2 transformer layers with 512 linear units as the decoder. All models are trained with dynamic chunks and conducted inference in chunks of 8.

The output units include 210 context-independent (CI) phones, a $\langle \text{sos/eos} \rangle$ symbol, and a $\langle \text{eok} \rangle$ symbol. The model consists of multiple components with their respective parameter sizes. The encoder has a size of 3.75M, the CTC decoder has a size of 0.2M, the attention decoder has a size of 0.6M, the keyword encoder has a size of 0.16M, and the bias module has a size of 0.04 M. In total, the model has 4.75M parameters. It is important to note that the keyword encoder only works for the system initialization, and the decoder is activated only when the streaming branch detects keyword candidates. This allows the two-pass model to have only 0.04M extra always-on parameters compared to the baseline model, making it suitable for deployment on edge devices.

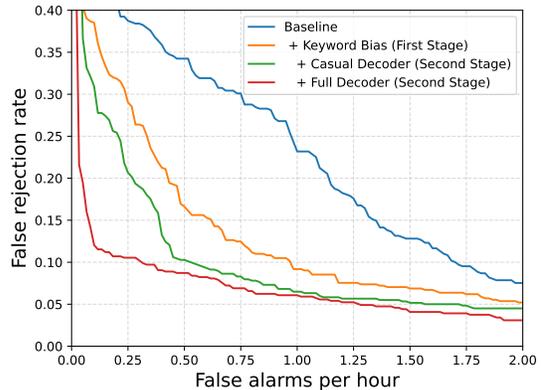


Fig. 3. ROC curves for different systems.

3.3. Evaluation Metrics

We evaluate the performance in terms of the receiver operating characteristic (ROC) curve and F1-score. To generate the ROC curves for each keyword, we utilize corresponding positive samples and the 60-hour negative test set. To evaluate the overall performance of different keywords, we average the false rejection rates of all keywords at the same false alarm rate. This allows us to plot overall ROC curves, providing a comprehensive evaluation of models. Similar to the process of ROC curves, we compute the F1-score for each keyword and average them to evaluate the overall performance.

3.4. Performance of Two-pass KWS Framework

In Fig. 3 and Table 2, we show the results of our proposed framework and baseline model on the internal test sets. The baseline system consists of an acoustic encoder with the CTC decoder. Integrating keyword information with keyword bias into the baseline model we get the first stage model. Compared with the baseline model, the first stage model increases the wake-up rate by 20% under the condition of 0.5 false alarms per hour. This shows that integrating keyword information through attention is a very effective method for open-vocabulary keyword spotting. Both the causal decoder and the full decoder introduced in Sec 2.3.3 can further improve performance on the first stage model. The causal decoder utilizes the output of the streaming encoder for rescoring, making better use of contextual information and achieving a 6% improvement in wake-up rate under a 0.5 false wake-up condition. On the other hand, the full decoder uses the full-chunk encoder to obtain more accurate acoustic representations, resulting in a 10% improvement under a 0.25 false wake-up condition.

Table 2. F1-scores of different systems on varying lengths

Model	F1-score in different word length			
	Overall	2	3	4
Baseline	0.790	0.711	0.813	0.847
+ Encoder Integration	0.872	0.791	0.893	0.932
+ Casual Decoder	0.894	0.793	0.926	0.963
+ Full Decoder	0.910	0.807	0.946	0.977

3.5. Impact of Varying Keyword Length

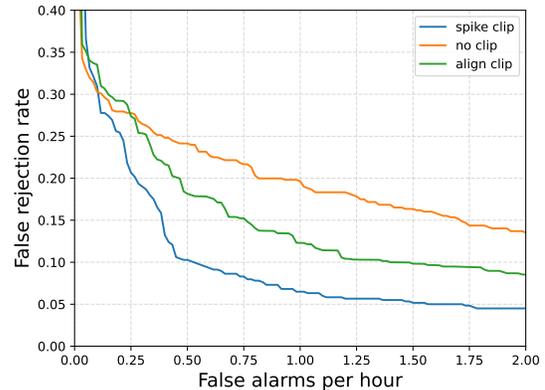
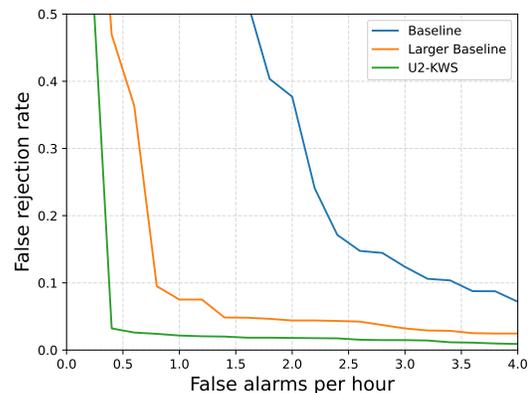
To investigate the impact of keyword lengths on system performance, we conducted experiments by grouping the test sets based on the character number of keywords and evaluated them separately. In Table 2, we present the average F1-scores for each keyword length. Our findings indicate that longer keywords tend to result in better overall performance. This can be attributed to the fact that shorter keywords are more prone to generating false alarms. We can observe that the longer the keywords are, the greater the improvement the decoder branch brings. There are two reasons for this. First, for short sequences, the structural benefit of the decoder which makes full use of context becomes smaller, and the second is that long sequences need to calculate more steps in the decoder, which can give full play to the distinguishing ability of the decoder.

3.6. Impact of Decoder Strategy

We conduct comparative experiments with the causal decoder to explore various methods of clipping the encoder output during decoder training and inference. The results of these experiments are depicted in Fig 4. The method without clipping assigns the task of searching keyword-related acoustic representations to the cross-attention of the decoder, which increases the complexity of decoder tasks. Using timestamp information in the first stage to clip the output of the encoder can make the decoder focus on the keyword validation task, so the clip method is better than the non-clip method. Although the method based on CTC alignment can get accurate timestamps for clean positive samples, it will output too long or too short clipping ranges for negative samples and difficult positive samples, resulting in decoder instability. In contrast, the spike-based methods provide a more stable clipping range by deriving the timestamp based on the number of peaks rather than specific content. Although this method may clip redundant parts, the cross-attention of the decoder is resilient to a small amount of redundant information. Thus, the spike-based method yields the best results.

3.7. Performance on AISHELL-1

To ensure our findings are reproducible, we carry out additional experiments using the publicly available Mandarin

**Fig. 4.** ROC curves for systems with different clip strategies.**Fig. 5.** ROC curves for different systems on AISHELL-1.

speech corpus AISHELL-1. The ROC curve presented in Fig. 5 clearly demonstrates that our proposed framework U2-KWS significantly outperforms the baseline. These results are consistent with the conclusions drawn from our analysis of the internal dataset.

4. CONCLUSIONS

In this paper, we propose a novel two-pass open-vocabulary keyword spotting framework U2-KWS. The framework uses the streaming branch as the first stage model to detect keyword candidates, while the non-streaming branch is activated to make further validation only when the first stage model has detected keyword candidates. To make the model sensitive to the keywords, we introduce keyword bias into both branches and design a two-pass training process for KWS with keyword sampling and encoder clipping. Our method outperforms the baseline by a significant margin on both internal and open datasets. We also set up additional experiments to explore the effectiveness of the decoder strategy and the impact of different keyword lengths.

5. REFERENCES

- [1] Fengpei Ge and Yonghong Yan, “Deep neural network based wake-up-word speech recognition with two-stage detection,” in *Proc. ICASSP*. 2017, pp. 2761–2765, IEEE.
- [2] Minhua Wu, Sankaran Panchapagesan, Ming Sun, Jiacheng Gu, Ryan Thomas, Shiv Naga Prasad Vitaladevuni, Björn Hoffmeister, and Arindam Mandal, “Monophone-based background modeling for two-stage on-device wake word detection,” in *Proc. ICASSP*. 2018, IEEE.
- [3] Tara N. Sainath and Carolina Parada, “Convolutional neural networks for small-footprint keyword spotting,” in *Proc. Interspeech*. 2015, pp. 1478–1482, ISCA.
- [4] Guoguo Chen, Carolina Parada, and Georg Heigold, “Small-footprint keyword spotting using deep neural networks,” in *Proc. ICASSP*. IEEE, 2014, pp. 4087–4091.
- [5] Bo Wei, Meirong Yang, Tao Zhang, Xiao Tang, Xing Huang, Kyuhong Kim, Jaeyun Lee, Kiho Cho, and Sung-Un Park, “End-to-end transformer-based open-vocabulary keyword spotting with location-guided local attention,” in *Proc. Interspeech*. 2021, pp. 361–365, ISCA.
- [6] Christin Jose, Yuriy Mishchenko, Thibaud Sénéchal, Anish Shah, Alex Escott, and Shiv Naga Prasad Vitaladevuni, “Accurate detection of wake word start and end using a CNN,” in *Proc. Interspeech*. 2020, pp. 3346–3350, ISCA.
- [7] Axel Berg, Mark O’Connor, and Miguel Tairum Cruz, “Keyword transformer: A self-attention model for keyword spotting,” in *Proc. Interspeech*. 2021, pp. 4249–4253, ISCA.
- [8] Guoguo Chen, Carolina Parada, and Tara N. Sainath, “Query-by-example keyword spotting using long short-term memory networks,” in *Proc. ICASSP*. 2015, pp. 5236–5240, IEEE.
- [9] Jinmiao Huang, Waseem Gharbieh, Han Suk Shim, and Eugene Kim, “Query-by-example keyword spotting system using multi-head attention and soft-triple loss,” in *Proc. ICASSP*. 2021, pp. 6858–6862, IEEE.
- [10] Théodore Bluche, Maël Primet, and Thibault Gisselbrecht, “Small-footprint open-vocabulary keyword spotting with quantized LSTM networks,” *arXiv preprint arXiv:2002.10851*, 2020.
- [11] Kyuyeon Hwang, Minjae Lee, and Wonyong Sung, “Online keyword spotting with a character-level recurrent neural network,” *arXiv preprint arXiv:1512.08903*, 2015.
- [12] Christopher T. Lengerich and Awni Y. Hannun, “An end-to-end architecture for keyword spotting and voice activity detection,” in *Proc. NIPS*, 2016.
- [13] Yimeng Zhuang, Xuankai Chang, Yanmin Qian, and Kai Yu, “Unrestricted vocabulary keyword spotting using LSTM-CTC,” in *Proc. Interspeech*. 2016, pp. 938–942, ISCA.
- [14] Zuozhen Liu, Ta Li, and Pengyuan Zhang, “RNN-T based open-vocabulary keyword spotting in mandarin with multi-level detection,” in *Proc. ICASSP*. IEEE, 2021, pp. 5649–5653.
- [15] Yao Tian, Haitao Yao, Meng Cai, Yaming Liu, and Zejun Ma, “Improving rnn transducer modeling for small-footprint keyword spotting,” in *Proc. ICASSP*. IEEE, 2021, pp. 5624–5628.
- [16] Siddharth Sigtia, Pascal Clark, Rob Haynes, Hywel Richards, and John Bridle, “Multi-task learning for voice trigger detection,” in *Proc. ICASSP*. 2020, pp. 7449–7453, IEEE.
- [17] Runyan Yang, Gaofeng Cheng, Haoran Miao, Ta Li, Pengyuan Zhang, and Yonghong Yan, “Keyword search using attention-based end-to-end asr and frame-synchronous phoneme alignments,” *IEEE ACM Trans. Audio Speech Lang. Process.*, vol. 29, pp. 3202–3215, 2021.
- [18] Siddharth Sigtia, John Bridle, Hywel Richards, Pascal Clark, Erik Marchi, and Vineet Garg, “Progressive voice trigger detection: Accuracy vs latency,” in *Proc. ICASSP*. IEEE, 2021, pp. 6843–6847.
- [19] Binbin Zhang, Di Wu, Zhuoyuan Yao, Xiong Wang, Fan Yu, Chao Yang, Liyong Guo, Yaguang Hu, Lei Xie, and Xin Lei, “Unified streaming and non-streaming two-pass end-to-end model for speech recognition,” *arXiv preprint arXiv:2012.05481*, 2020.
- [20] Arun Narayanan, Tara N. Sainath, Ruoming Pang, Jiahui Yu, Chung-Cheng Chiu, Rohit Prabhavalkar, Ehsan Variani, and Trevor Strohman, “Cascaded encoders for unifying streaming and non-streaming ASR,” in *Proc. ICASSP*. 2021, pp. 5629–5633, IEEE.
- [21] Zhanheng Yang, Sining Sun, Jin Li, Xiaoming Zhang, Xiong Wang, Long Ma, and Lei Xie, “CaTT-KWS: A multi-stage customized keyword spotting framework based on cascaded transducer-transformer,” in *Proc. Interspeech*. 2022, pp. 1681–1685, ISCA.

- [22] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *Proc. NIPS*, 2017, pp. 5998–6008.
- [23] Yanzhang He, Rohit Prabhavalkar, Kanishka Rao, Wei Li, Anton Bakhtin, and Ian McGraw, “Streaming small-footprint keyword spotting using sequence-to-sequence models,” in *Proc. ASRU*. IEEE, 2017, pp. 474–481.
- [24] Yu Xi, Tian Tan, Wangyou Zhang, Baochen Yang, and Kai Yu, “Text adaptive detection for customizable keyword spotting,” in *Proc. ICASSP*. 2022, pp. 6652–6656, IEEE.
- [25] Théodore Bluche and Thibault Gisselbrecht, “Predicting detection filters for small footprint open-vocabulary keyword spotting,” in *Proc. Interspeech*. 2020, pp. 2552–2556, ISCA.
- [26] Bo Wei, Meirong Yang, Tao Zhang, Xiao Tang, Xing Huang, Kyuhong Kim, Jaeyun Lee, Kiho Cho, and Sung-Un Park, “End-to-End transformer-based open-vocabulary keyword spotting with location-guided local attention,” in *Proc. Interspeech*. 2021, pp. 361–365, ISCA.
- [27] Hyeon-Kyeong Shin, Hyewon Han, Doyeon Kim, Soo-Whan Chung, and Hong-Goo Kang, “Learning audio-text agreement for open-vocabulary keyword spotting,” in *Proc. Interspeech*. 2022, pp. 1871–1875, ISCA.
- [28] Yujie Yang, Kun Zhang, Zhiyong Wu, and Helen Meng, “Keyword-specific acoustic model pruning for open-vocabulary keyword spotting,” in *Proc. ICASSP 2023*. IEEE, 2023, pp. 1–5.
- [29] Hui Bu, Jiayu Du, Xingyu Na, Bengu Wu, and Hao Zheng, “AISHELL-1: an open-source mandarin speech corpus and a speech recognition baseline,” in *Proc. O-COCOSDA*. 2017, pp. 1–5, IEEE.
- [30] Anmol Gulati, James Qin, Chung-Cheng Chiu, Niki Parmar, Yu Zhang, Jiahui Yu, Wei Han, Shibo Wang, Zhengdong Zhang, Yonghui Wu, and Ruoming Pang, “Conformer: Convolution-augmented transformer for speech recognition,” in *Proc. Interspeech*. 2020, pp. 5036–5040, ISCA.
- [31] Zhuoyuan Yao, Di Wu, Xiong Wang, Binbin Zhang, Fan Yu, Chao Yang, Zhendong Peng, Xiaoyu Chen, Lei Xie, and Xin Lei, “WeNet: production oriented streaming and non-streaming end-to-end speech recognition toolkit,” in *Proc. Interspeech*. 2021, pp. 4054–4058, ISCA.
- [32] Golan Pundak, Tara N. Sainath, Rohit Prabhavalkar, Anjali Kannan, and Ding Zhao, “Deep Context: end-to-end contextual speech recognition,” in *Proc. SLT*. 2018, pp. 418–425, IEEE.
- [33] W Ronny Huang, Shuo-Yiin Chang, Tara N Sainath, Yanzhang He, David Rybach, Robert David, Rohit Prabhavalkar, Cyril Allauzen, Cal Peyser, and Trevor D Strohman, “E2E segmentation in a two-pass cascaded encoder asr model,” in *Proc. ICASSP*. 2023, pp. 1–5, IEEE.
- [34] Zhengkun Tian, Jiangyan Yi, Jianhua Tao, Ye Bai, Shuai Zhang, and Zhengqi Wen, “Spike-triggered non-autoregressive transformer for end-to-end speech recognition,” in *Proc. Interspeech*, 2020, pp. 5026–5030.
- [35] Daniel S. Park, William Chan, Yu Zhang, Chung-Cheng Chiu, Barret Zoph, Ekin D. Cubuk, and Quoc V. Le, “SpecAugment: A simple data augmentation method for automatic speech recognition,” in *Proc. Interspeech*. 2019, pp. 2613–2617, ISCA.