

QUERY-BY-EXAMPLE SPOKEN TERM DETECTION USING ATTENTION-BASED MULTI-HOP NETWORKS

Chia-Wei Ao and Hung-yi Lee

Graduate Institute of Communication Engineering, National Taiwan University
r04942094@ntu.edu.tw, hungyilee@ntu.edu.tw

ABSTRACT

Retrieving spoken content with spoken queries, or query-by-example spoken term detection (STD), is attractive because it makes possible the matching of signals directly on the acoustic level without transcribing them into text. Here, we propose an end-to-end query-by-example STD model based on an attention-based multi-hop network, whose input is a spoken query and an audio segment containing several utterances; the output states whether the audio segment includes the query. The model can be trained in either a supervised scenario using labeled data, or in an unsupervised fashion. In the supervised scenario, we find that the attention mechanism and multiple hops improve performance, and that the attention weights indicate the time span of the detected terms. In the unsupervised setting, the model mimics the behavior of DTW, and it performs as well as DTW but with a lower run-time complexity.

Index Terms— Attention-based Multi-hop Network

1. INTRODUCTION

Retrieving spoken content with spoken queries, also known as query-by-example spoken term detection (STD) [1–6], is attractive because hand-held or wearable devices make spoken queries a natural choice. The most intuitive way to search over spoken content for a spoken query is to directly match the audio signals to find those audio snippets that sound like the spoken query, and dynamic time warping (DTW) [7] is widely used. Despite DTW’s wide use, it has several drawbacks. As typical DTW does not have trainable parameters, even in an online system that collects the training data from user feedback, the data cannot be directly used to improve the algorithm. In addition, the time complexity of DTW is usually proportional to the product of the lengths of the spoken queries and audio segments, which for real applications is usually excessive.

Query-by-example STD by representing each word segment as a vector [8–12] is much more efficient than the conventional Dynamic Time Warping (DTW) based approaches, because only the similarities between two single vectors are needed, in addition to the significantly better retrieval performance obtained [11]. Audio segment representation is still an open problem. Several approaches have been successfully used in STD [9, 13–15], but these approaches were developed primarily in more heuristic ways, rather than deep learning. By learning RNN with an audio segment as the input and the corresponding word as the target, the outputs of the hidden layer at the last few time steps can be taken as the representation of the input segment [16]. Audio segment embedding can also be jointly learned with their corresponding character sequences

by multi-view approach [17]. Sequence-to-sequence Autoencoder is used to represent variable-length audio segments by vectors with fixed dimensionality, which is referred to as Audio Word2Vec [11]. This previous approach assumes that speech segments to be retrieved have been pre-segmented at word boundaries, which is not realistic. However, it was shown that neural embeddings learned from pre-segmented audio can be applied for embedding arbitrary segments [18]. In this paper, we use attention mechanism to locate the time span of the input query in the utterances to be retrieved, so word boundary segmentation is not needed at both the training and testing stages, and attention is shown to improve the query-by-example performance.

The target of this paper is to develop an end-to-end deep learning model for query-by-example STD. On STD with text query, end-to-end approaches have been explored, in which a function which can map the acoustic features of an utterance and a text query to a confidence score is developed. Along this direction, encouraging results have been obtained based on structured support vector machine (SVM) [19–21]. However, learning structured SVM is computationally intensive, so this approach is hard to scale. An end-to-end deep learning based system for text query STD has been proposed [22]. Attention mechanism and multiple hops are not used in the model, which is different from the proposed approach in this paper.

In this paper, we propose an end-to-end query-by-example STD model. The model is an attention-based multi-hop network, the input of which is a query and an audio segment (containing several utterances), and the output a confidence score representing whether the audio segment includes the query term. In the network, the input spoken query is represented as a vector by an LSTM encoder. We use the attention mechanism to locate the time span of the query term in the audio segment. Similar to query expansion, multiple hops are used to update the spoken queries via information extracted from the audio segment. Then a key term detector determines whether the query term exists in the input audio segment. These network components are all learned end-to-end, and the model can be learned in a supervised or unsupervised setting. In a supervised setting, the model is learned from a set of labeled data, which can be collected by user feedback in real applications. In an unsupervised setting, the neural network mimics the behavior of DTW, and it performs as well as DTW but with a lower run-time complexity.

2. FRAMEWORK AND TRAINING SCENARIO

Shown in the upper half of Fig 1 is the framework of query-by-example STD using an end-to-end network. The network input is the spoken query and an audio segment in the database to be retrieved. Both the spoken query and the audio segment are represented by acoustic features such as MFCCs. In this paper, the audio

This work was sponsored by Ministry of Science and Technology, R.O.C.

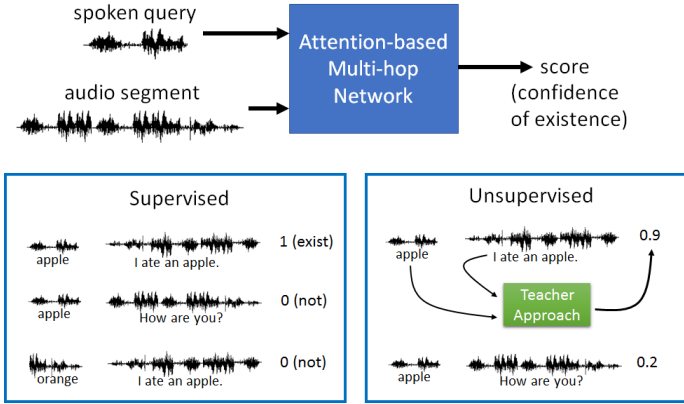


Fig. 1. Framework and training scenarios

segments cover several utterances, and as such are much longer than the spoken queries. The output of the network is a scalar. The scalar represents the confidence that the term in the spoken query exists in the audio segment. Given a spoken query, the system ranks the audio segments in the database according to the confidence scores, and yields the search results. It may seem that the proposed approach only extracts the target audio segments as opposed to locating the time spans of the query terms in the segments; in truth, the time spans are determined based on the attention mechanism in the network. This is shown in the experimental results.

As shown in the lower half of Fig. 1, the network is trained in two different scenarios. In the first scenario (left lower corner), labeled examples are collected, each pair of which is composed of a spoken query and an audio segment, including a label indicating whether the segment contains the term in the spoken query. In this case, network training is cast as a binary classification problem. Since labeled data is used in the first scenario, the network achieves better performance than query-by-example approaches such as DTW that use no labeled data. In the second scenario (right lower corner), given a set of query-segment pairs, an existing query-by-example STD approach referred to here as the teacher approach (any method with good performance could be used here) assigns a score to each example pair. This score represents the confidence that the segment includes the query. The network thus learns to predict the confidence scores of the teacher approach given the same example pairs; this is hence a regression task. In the second scenario, as no extra labeled data is needed, it is unsupervised. Since the network is learned from an existing approach, it cannot outperform its teacher. However, if the network performance is equivalent to that of the teacher approach, and if the network is faster than the teacher, it may be reasonable to use the network at testing time instead of the teacher approach.

3. ATTENTION-BASED MULTI-HOP NETWORK

In this section we describe the model architecture of the attention-based multi-hop network.

3.1. Query Representation

Fig. 2 (A) illustrates the encoding of the input spoken query into a vector representation V_Q . The input query is a sequence of T vectors, x_1, x_2, \dots, x_T , each vector x_i of which is an acoustic fea-

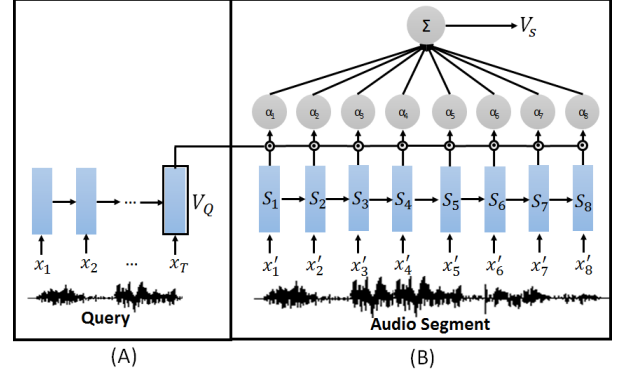


Fig. 2. Attention mechanism

ture vector such as MFCC. In Fig. 2 (A), a long short-term memory (LSTM) network [23] takes one frame from the input spoken query sequentially at a time. After going through all the frames in the query, the query vector representation V_Q is the hidden layer output of the LSTM network at the last time index.

3.2. Audio Segment Representation with Attention

Fig. 2 (B) shows an audio segment (containing several utterances) in the database to be retrieved; although this is a lengthy acoustic feature sequence, we show only eight features for simplicity. The LSTM in Fig 2 (B) goes through the whole document and encodes each frame¹. The vector representation of the t -th frame S_t is the hidden layer outputs of the LSTM network. This process can be completed off-line, before the spoken query is submitted.

Then the attention value α_t for each time index t is the cosine similarity between the query vector V_Q (obtained in Fig. 2 (A)) and the vector representation S_t of each frame, $\alpha_t = S_t \odot V_Q$, where symbol \odot denotes cosine similarity between two vectors. We normalize attention values α_t as α'_t . The score list is normalized using the softmax activation function:

$$\alpha'_i = \frac{\exp(\alpha_i)}{\sum_{i=1}^T \exp(\alpha_i)} \quad (1)$$

This has been widely used in many existing neural attention frameworks [24–27]. Then vectors S_t from the LSTM network for every frame in the audio segment are weighted with this normalized attention value α'_t and summed to yield the segment representation vector $V_S = \sum_t \alpha'_t S_t$, which is used to determine the confidence score for spoken query V_Q . To ensure a time complexity linear to the length of the input audio segment, we do not use more sophisticated attention models [28]; thus the approach is faster than DTW.

3.3. Hopping

Fig. 3 illustrates hopping: the input spoken query is first converted into a vector V_{Q_1} by the module in Fig. 2 (A), after which the module in (B) uses this V_{Q_1} to compute the attention values α_t to obtain the story vector V_{S_1} . Then V_{Q_1} and V_{S_1} are summed to form the new question vector V_{Q_2} . This process is the first hop (hop 1). The output of the first hop V_{Q_2} can be used to compute the new attention values to obtain a new story vector V_{S_2} . This can be seen as

¹The LSTMs used in Figs 2 (A) and (B) are the same.

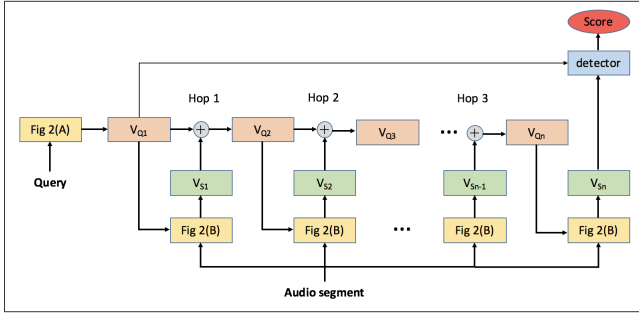


Fig. 3. The hopping operation

relevance feedback [13, 29], where the machine goes over the audio segment again, extracting information to expand the query to form a new query vector. Again, V_{Q_1} and V_{S_1} are summed to form V_{Q_2} (hop 2). After n hops (n is predefined), the output of the last hop V_{S_n} is used to calculate the confidence score.

3.4. Keyword Detection

Finally, as shown in the upper half of Fig 3, a detector determines the confidence score based on query vector representation V_Q and utterance vector representation V_{S_n} . Here we use three ways to calculate this score: (1) Use the cosine similarity between V_Q and V_{S_n} as the score; (2) Use the detector – a connected feedforward neuron network – taking V_Q and V_{S_n} as input, and output a scalar as the confidence score; (3) Combine approaches (1) and (2): a neural network takes as input the query vector V_Q , utterance vector V_{S_n} , and cosine similarity, and outputs a score.

4. EXPERIMENTAL SETUP

We used the LibriSpeech corpus [30] as the data for the experiments. To train the attention-based multi-hop network, some query-segment pairs were needed as training examples. 70,000 training examples were used in the experiments, including 500 different spoken queries; all audio segments were from the LibriSpeech *train-clean-360* set. In the supervised scenario, the label for each example (a query-segment pair) specified whether the audio segment contains the spoken query². In the unsupervised scenario, each example is labeled using the score from the DTW algorithm [31]. There are three testing sets. In all testing sets, the audio segments were from the *train-other-500* set in LibriSpeech; thus the audio segments in the training and testing sets did not overlap. As described below, the spoken queries were different.

- Testing Set 1: There were 1,500 query-segment pairs, including 30 different spoken queries (each query has 50 examples in average). The spoken queries were all from the training set.
- Testing Set 2: As with testing set 1, this set also had 1,500 query-segment pairs with 30 different spoken queries (each query has 50 examples in average). The spoken queries in this set had the same text form as testing set 1, but did not come from the training set.

²This is easily determined using the manual transcriptions of the audio segments available from the LibriSpeech corpus.

- Testing Set 3: This set had 10,000 query-segment pairs with 100 different spoken queries (each query has 100 examples in average). In this set, the spoken queries were not from the training set, and the text form of the spoken queries never appeared in the training queries.

All the spoken queries corresponds to single words, but the proposed approach can also be applied on phrases. 39-dimension MFCCs were used as the acoustic features. Both attention-based multi-hop network and DTW baseline used the same set of features, so they can be fairly compared. Mean average precision (MAP) was used as the evaluation measure.

The network structure and hyper-parameters were set as below without further tuning if not specified. The LSTM encoder consisted of two hidden layers with 128 LSTM units. The networks were trained for 100 epochs using ADAM [32]. The keyword detector was a network with four hidden layers with 128, 64, 32 and 2 neurons respectively. In the supervised scenario, the attention-based multi-hop network was a binary classifier trained using cross-entropy loss; in the unsupervised scenario, mean square error was the loss function.

5. EXPERIMENTAL RESULTS

In Sections 5.1, 5.2 and 5.3, we consider the supervised scenario. The results of the unsupervised scenario are presented in Section 5.4.

5.1. Attention-based Model

Table 1. Results of attention-based network with a single hop. Rows (A) and (B) are baselines. Part (C) is attention-based networks. NN, Cos, and NN+Cos are respectively the three detectors from Section 3.4.

Approach		Test 1	Test 2	Test 3
(A): DTW		0.6173	0.5778	0.5678
(B): Network without Attention		0.5935	0.5563	0.5468
(C): Attention-based Network	(1) NN	0.6523	0.6246	0.5754
	(2) Cosine	0.6331	0.6043	0.5746
	(3) NN+Cos	0.6268	0.6370	0.5759
(D): (A)+(C)	(1) NN	0.6720	0.6340	0.5868
	(2) Cosine	0.6433	0.6002	0.5843
	(3) NN+Cos	0.6451	0.6309	0.5808

Table 1 shows the results of the attention-based model with a single hop. Rows (A) and (B) are two baselines: Row (A) is the MAP of the search results ranked according to DTW similarities on the three testing sets, and row (B) is the model without attention mechanism. We used an LSTM to encode both the spoken query and audio segment as a vector representation by taking the hidden layer output of the LSTM network at the last time index. Next, we input to the neural network key term detector the query vector and audio segment representations. The detector then outputs a score representing the confidence that the query appears in the audio segment. We find that without the attention mechanism, even though the networks are learned from labeled training data, they are outperformed by DTW, which needs no training data (rows (B) v.s. (A)).

We observe that the results of attention-based networks outperform those without attention and DTW (part (C) v.s. rows (A), (B)). From Table 1, we note that compared to DTW, the attention model yields larger improvements on test sets 1 and 2. This shows that even though the training and testing queries are from different speakers, the attention-based model still learns the keyword acoustic patterns,

which are speaker-independent. However, the attention-based network yields little improvement on test set 3; it is more difficult for the network to transfer what it has learned to keywords it has never seen before. We evaluated the three detectors mentioned in Section 3.4, denoted in Table 1 as NN, Cos, and NN+Cos. Regardless of the model, using the network for keyterm detection always works better than simply computing cosine similarity (NN v.s Cos). Taking cosine similarity as another input to the keyterm detection network (NN+Cos) does not improve performance on test set 1, but does improve the performance in some cases on test sets 2 and 3.

Part (D) shows the integration of the DTW output in row (A) and the attention-based model in part (C), for which the integration weight is 0.4 for DTW and 0.6 for the attention-based model. We observe improvements for test sets 1 and 3. This shows that DTW and the attention-based model are complementary.

5.2. Attention analysis

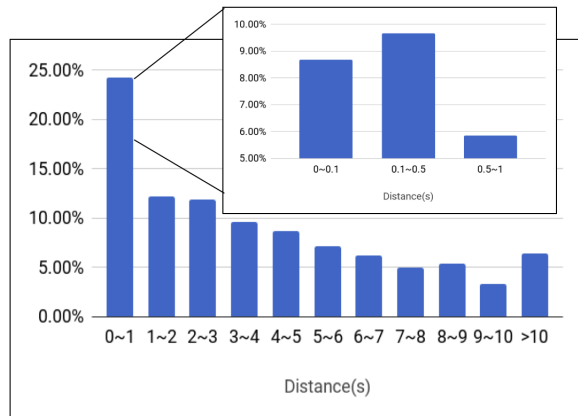


Fig. 4. Time differences between maximum attention weight and end of query.

In spoken term detection, we seek to determine not only whether the query term exists in the audio segments, but sometimes also – if they exist – the time spans of these query terms within the segments. We find that the attention weights reveal the time spans of the query terms. In Fig. 4 we present the analysis of test set 1. For all query terms and for the audio segments containing the query terms, we compute the time difference between the position with the highest attention weight and the end of the query. The horizontal axis in the figure shows the time difference, while the vertical axis is the percentage among all the audio segments considered. From the figure, we find that distances under one second accounted for 25% of the cases; that is, in these cases the attention mechanism located the query term with less than a one-second error. Further analysis in the enclosed subfigure shows the time duration from 0 to 1 seconds. We find that most of the time differences fall between 0.1 and 0.5; this shows that the highest peak of the attention weights are quite close to the query word. This suggests that attention yields a precise focus on the end location of the query.

5.3. Multiple Hopping

Table 2 shows the results when using multiple hops to generate audio segment representations. The results in row (B) are the results

Table 2. Multiple-hop results

	Test set 1	Test set 2	Test set 3
(A): DTW	0.6173	0.5778	0.5678
(B): 1-hop	0.6523	0.6246	0.5754
(C): 2-hop	0.6472	0.6430	0.5842
(D): 3-hop	0.6676	0.6404	0.5837
(E): 4-hop	0.6417	0.6476	0.5792
(F): (A)+(D)	0.6789	0.6430	0.5830

without multiple hops, also shown in row (C-1) of Table 1; the results with 2 to 4 hops are those in rows (C) to (E). Multiple hops outperform single hops (rows (C) to (E) v.s. (B)), except for 1 and 3 hops on test set 1. This shows that hopping improves model generality because in test sets 2 and 3 the training and testing data are mismatched. In row (F), we also integrated the DTW and 3-hop results, yielding further improvements to the MAP score on test set 1.

5.4. Unsupervised Scenario

Table 3. Results of attention-based multi-hop network learning from a teacher approach (DTW).

	Test set 1	Test set 2	Test set 3
(A) DTW	0.6173	0.5778	0.5678
(B) Attention + 1-hop	0.6128	0.5893	0.5548
(C) Attention + 3-hop	0.6141	0.5964	0.5702

Here the network is learned from a teacher approach, so no extra label data is needed. We use DTW as the teacher approach, and normalize the DTW similarity scores between 0 and 1 as the target of regression. The results are shown in Table 3. From the table, we find that the performance of the attention-based network without multiple hops is comparable to DTW (rows (B) v.s. (A)), and that the 3-hop network outperforms DTW on test sets 2 and 3 (rows (C) v.s. (A)). Here we emphasize that the time complexity of the network during testing is far less than DTW: given a document length of M and a query length of N , the time complexity of DTW is $O(M \times N)$, while the time complexity of the network is $O(M \times n)$, where n is the number of hops³. Therefore, it is reasonable to replace DTW with a network learned from it.

6. CONCLUSION

In this paper, we propose an end-to-end query-by-example STD model based on an attention-based multi-hop network. The model can be trained in either an supervised or unsupervised fashion. In the supervised scenario, we show that attention and multiple hops are both very helpful, and that the attention weights of the proposed model reveal the time span of the input keyterm. In the unsupervised setting, the neural network mimics DTW behavior, and achieves performance comparable to DTW with shorter runtimes. In the future, we will explore more new attention-based models, and investigate new models which directly output time spans instead of a confidence score. We will also compare the performance of the proposed approach and DTW on posterior features and cross-lingual bottleneck features.

³We implemented DTW in C++ and the network using Tensorflow. On average, without using a GPU, the proposed approach was 7 times faster than DTW.

7. REFERENCES

- [1] Florian Metze, Xavier Anguera, Etienne Barnard, Marelle Davel, and Guillaume Gravier, “Language independent search in MediaEval’s spoken web search task,” *Computer Speech & Language*, vol. 28, no. 5, pp. 1066 – 1082, 2014.
- [2] Xavier Anguera, Florian Metze, Andi Buzo, Igor Szoke, and Luis Javier Rodriguez-Fuentes, “The spoken web search task,” in *MediaEval 2013 Workshop*, 2013.
- [3] Xavier Anguera, Luis-Javier Rodriguez Fuentes, Igor Szke, Andi Buzo, and Florian Metze, “Query by example search on speech at mediaeval 2014,” in *MediaEval 2014 Workshop*, 2014.
- [4] Andi Buzo, Xavier Anguera, Florian Metze, Jorge Proenca, Martin Lojka, and Xiao Xiong, “Query by example search on speech at mediaeval 2015,” .
- [5] Igor Szke, Miroslav Skcel, Luk Burget, and Jan ernock, “Coping with channel mismatch in query-by-example - but QUESST 2014,” in *ICASSP*, 2015.
- [6] Cheung-Chi Leung, Lei Wang, Haihua Xu, Jingyong Hou, Van Tung Pham, Hang Lv, Lei Xie, Xiong Xiao, Chongjia Ni, Bin Ma, Eng Siong Chng, and Haizhou Li, “Toward high-performance language-independent query-by-example spoken term detection for MediaEval 2015: Post-evaluation analysis,” in *INTERSPEECH*, 2015.
- [7] Yaodong Zhang and James R Glass, “Unsupervised spoken keyword spotting via segmental dtw on gaussian posteriorgrams,” in *Automatic Speech Recognition & Understanding, 2009. ASRU 2009. IEEE Workshop on*. IEEE, 2009, pp. 398–403.
- [8] Keith Levin, Katharine Henry, Aren Jansen, and Karen Livescu, “Fixed-dimensional acoustic embeddings of variable-length segments in low-resource settings,” in *Automatic Speech Recognition and Understanding (ASRU), 2013 IEEE Workshop on*. IEEE, 2013, pp. 410–415.
- [9] Keith Levin, Aren Jansen, and Benjamin Van Durme, “Segmental acoustic indexing for zero resource keyword search,” in *ICASSP*, 2015.
- [10] Herman Kamper, Weiran Wang, and Karen Livescu, “Deep convolutional acoustic word embeddings using word-pair side information,” in *ICASSP*, 2016.
- [11] Yu-An Chung, Chao-Chung Wu, Chia-Hao Shen, Hung-Yi Lee, and Lin-Shan Lee, “Audio word2vec: Unsupervised learning of audio segment representations using sequence-to-sequence autoencoder,” *arXiv preprint arXiv:1603.00982*, 2016.
- [12] Yougen Yuan, Cheung-Chi Leung, Lei Xie, Bin Ma, and Haizhou Li, “Learning neural network representations using cross-lingual bottleneck features with word-pair information,” in *INTERSPEECH*, 2016.
- [13] Hung-Yi Lee and Lin-Shan Lee, “Enhanced spoken term detection using support vector machines and weighted pseudo examples,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 21, no. 6, pp. 1272–1284, 2013.
- [14] I.-F. Chen and C.-H. Lee, “A hybrid HMM/DNN approach to keyword spotting of short words,” in *INTERSPEECH*, 2013.
- [15] A. Norouzi, A. Jansen, R. Rose, and S. Thomas, “Exploiting discriminative point process models for spoken term detection,” in *INTERSPEECH*, 2012.
- [16] Guoguo Chen, Carolina Parada, and Tara N. Sainath, “Query-by-example keyword spotting using long short-term memory networks,” in *ICASSP*, 2015.
- [17] Wanjia He, Weiran Wang, and Karen Livescu, “Multi-view recurrent neural acoustic word embeddings,” *ICLR*, 2017.
- [18] Shane Settle, Keith Levin, Herman Kamper, and Karen Livescu, “Query-by-example search with discriminative neural acoustic word embeddings,” *arXiv preprint arXiv:1706.03818*, 2017.
- [19] Rohit Prabhavalkar, Karen Livescu, Eric Fosler-Lussier, and Joseph Keshet, “Discriminative articulatory models for spoken term detection in low-resource conversational settings,” in *ICASSP*, 2013.
- [20] M. Wollmer, F. Eyben, J. Keshet, A. Graves, B. Schuller, and G. Rigoll, “Robust discriminative keyword spotting for emotionally colored spontaneous speech using bidirectional LSTM networks,” in *ICASSP*, 2009.
- [21] Joseph Keshet, David Grangier, and Samy Bengio, “Discriminative keyword spotting,” *Speech Communication*, vol. 51, pp. 317 – 329, 2009.
- [22] Kartik Audhkhasi, Andrew Rosenberg, Abhinav Sethy, Bhuvana Ramabhadran, and Brian Kingsbury, “End-to-end ASR-free keyword search from speech,” *arXiv preprint arXiv:1701.04313*, 2017.
- [23] Sepp Hochreiter and Jürgen Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [24] Huijuan Xu and Kate Saenko, “Ask, attend and answer: Exploring question-guided spatial attention for visual question answering,” *arXiv preprint arXiv:1511.05234*, 2015.
- [25] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” *arXiv preprint arXiv:1502.03044*, 2015.
- [26] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [27] Jason Weston, Sumit Chopra, and Antoine Bordes, “Memory networks,” *arXiv preprint arXiv:1410.3916*, 2014.
- [28] Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi, “Bidirectional attention flow for machine comprehension,” *arXiv preprint arXiv:1611.01603*, 2016.
- [29] Haipeng Wang, Tan Lee, Cheung-Chi Leung, Bin Ma, and Haizhou Li, “Using parallel tokenizers with DTW matrix combination for low-resource spoken term detection,” in *ICASSP*, 2013.
- [30] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur, “Librispeech: an ASR corpus based on public domain audio books,” in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*. IEEE, 2015, pp. 5206–5210.
- [31] “dynamic time warping approach,” github, <https://github.com/chunan/libdtw>.

- [32] Diederik Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.