
A state of the art end-to-end speech recognition algorithm with a homogenous structure.

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 TODO(chanw.com) Revise the abstract. This paper proposes a new end-to-end
2 speech recognition system referred to as Recurrent neural network Sequence Classi-
3 fication (RSC). Recently, end-to-end speech recognition systems have been gaining
4 more attention from researchers and several different structures have been proposed.
5 Typical examples include Connectionist Temporal Classification (CTC), Recurrent
6 Neural Network Transducer (RNN-T), and Sequence-to-Sequence Modeling using
7 the attention mechanism. The CTC-based approaches are using the CTC loss
8 assuming the conditional independence property. Although the structure of CTC
9 is simpler than other end-to-end systems, the performance is worse due to this
10 conditional independence assumption. RNN-T and attention-based approaches
11 have distinct components such as the prediction network and encoder in RNN-T,
12 and the encoder, the decoder, and the attention layer in the attention-based model.
13 Compared to these models, FSC has a homogeneous structure from the bottom
14 to the top layer. LSTM and max-pooling layer are repeated followed by the top
15 softmax layer. The embedded softmax output is fed back as input of each RNN
16 layer. The training starts with flat initialization and alignment is made after every
17 epoch. The loss is the Cross Entropy (CE) loss using this alignment result. In
18 our experimental results, FSC algorithm has shown better performance than more
19 complicated attention-based end-to-end speech recognition system. Another major
20 advantage is training is significantly faster.

1 Introduction

22 Recently, there has been tremendous improvements in speech recognition systems fueled by advances
23 in deep neural networks [7, 12, 14, 16, 17, 18].

24 TODO(chanw.com) Talk about the AI speakers.

25 TODO(chanw.com) Talk about the advancements of the end-to-end systems.

26 Recently, it has been frequently observed that if sequence-to-sequence end-to-end ASR systems are
27 trained on sufficiently large amounts of acoustic training data, they can outperform conventional
28 HMM-DNN/RNN hybrid systems [5, 1].

29 Recently, we observed that training with large-scale noisy data generated by a *Room Simulator* [4]
30 improves speech recognition accuracy dramatically. This system has been successfully employed for
31 training acoustic models for Google Home or Google voice search [4].

2 Review on sequence-to-sequence speech recognition algorithms

In this section, we review well-known *sequence-to-sequence* algorithms used in speech recognition. A sequence-to-sequence model maps a sequence of input acoustic features into a sequence of graphemes or words [15]. We denote a sequence of input acoustic feature vectors by \mathcal{X}_0^{M-1} and a sequence of target labels by \mathcal{Y}_0^{L-1} :

$$\mathcal{X}_0 = \{\vec{x}[0], \vec{x}[1], \dots, \vec{x}[M-1]\}, \quad (1a)$$

$$\mathcal{Y}_0 = \{y_0, y_1, \dots, y_{L-1}\}, \quad (1b)$$

where M is the number of frames in the input feature sequence and L is the number of labels in the output target sequence. In Fig. 1a ~ 1c, we show block diagrams of CTC [9], RNN-T [8, 10], and the attention-based model [6, 3] respectively. CTC in Fig. 1a and RNN-T in Fig. 1b are *frame-synchronous*, which means that the ASR system generates the output target for each input frame.

2.1 Connectionist Temporal Classification

Fig. 1a shows the entire structure of the Connectionist Temporal Classification (CTC) [9]. In CTC, we use the following CTC loss [9, 11]:

$$P(\mathcal{Y}_0^{L-1} | \mathcal{X}_0^{M-1}) = \sum_{\mathcal{Z}_0^{M-1} \in \mathcal{A}_{CTC}(\mathcal{X}_0^{M-1}, \mathcal{Y}_0^{L-1})} \prod_{m=0}^{M-1} P(\vec{z}[m] | \mathcal{X}_0^{M-1}), \quad (2)$$

where $\mathcal{A}_{CTC}(\mathcal{X}_0^{M-1}, \mathcal{Y}_0^{L-1})$ correspond to frame-level alignments of length M such that removing blanks and repeated symbols from \mathcal{Z}_0^{M-1} yields \mathcal{Y}_0^{L-1} .

3 Feedback Sequence Classifier

3.1 Neural Network Structure

Fig. 1d shows the entire structure of the FSC. As shown in this figure, LSTM layers and max-pool layers are interleaved from the bottom layer up to the third max-pool layer.

In Fig. 1d, speech features $\vec{x}[m]$ are sampled at every 10 ms, which is the usual speech feature frame rate in speech recognition [13, 2]. Using three two-to-one max-pool layers, $\vec{z}[p]$ has a frame period of 80 ms. The relationship between the temporal index p in $\vec{z}[p]$ and m in $\vec{x}[m]$ is given by:

$$p = \lceil m/8 \rceil \quad (3)$$

The reason for using the one dimensional max-pool layer is to make the neural network output have a similar rate to the that of each *phone* in utterances. Note that a *phone* is any distinct speech sound serving as a phonetic unit.

It has been observed that the average phone duration is between 50 ms ~ 100 ms [19, 20].

3.2 Training of the Feedback Sequence Classifier

In this section, we describe how to train FSC. The first step is making an alignment using the N-best alignment. The second step is updating the neural network parameters using the Cross Entropy (CE) criterion.

3.3 Training procedure

The first step is finding the optimal frame-synchronous sequence \mathcal{Z}

$$\mathcal{Z} = A(\mathcal{X} | \Theta)$$

Updates the model.

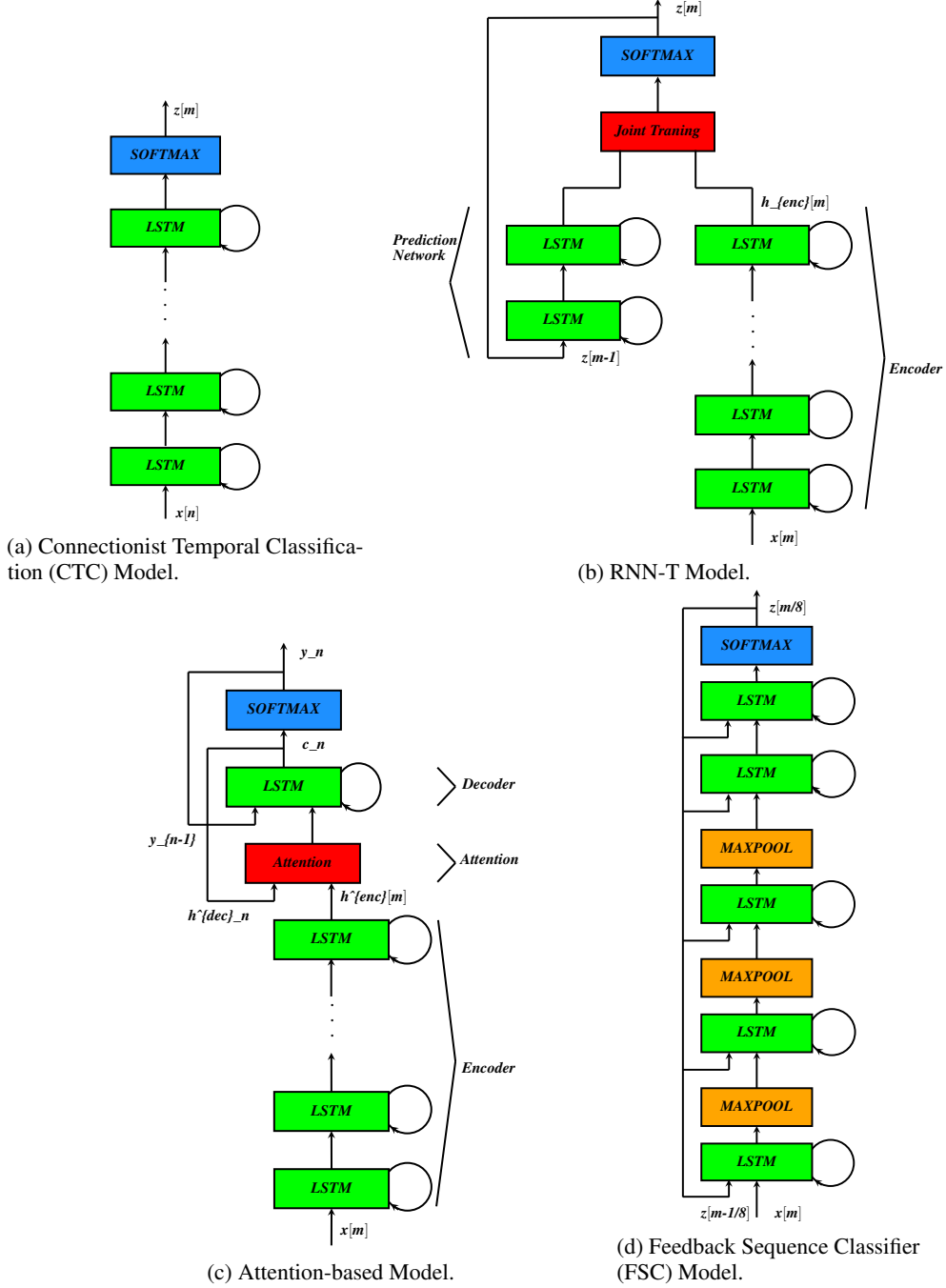


Figure 1: Comparison of block diagrams of different sequence-to-sequence speech recognition approaches. The proposed Feedback Sequence Classifier (FSC) is shown in Fig. 1d.

3.4 Viterbi alignment and N-best alignment

In conventional frame-wise CE-training, Viterbi alignment (a.k.a forced alignment) has been performed to obtain the frame-level acoustic unit boundaries TODO(chanwcom). Even though the Viterbi algorithm has advantages in simplicity and efficiency, it is based on the conditional independence assumption.

In TFSC, we propose the following N-best alignment approach rather than the Viterbi alignment algorithm to find the alignment information.

```

73   for  $m = 0, \dots, M - 1$  do
74       for  $l = 0, \dots, N_b - 1$  do
75            $\pi^{(l)}[m]$ 
76       end for
77   end for
78   if  $i \geq maxval$  then
79        $i \leftarrow 0$ 
80   else
81       if  $i + k \leq maxval$  then
82            $i \leftarrow i + k$ 
83       end if
84   end if

```

85 4 Experimental Results

86 Acknowledgments

87 References

88 References follow the acknowledgments. Use unnumbered first-level heading for the references. Any
89 choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font
90 size to small (9 point) when listing the references. **Remember that you can go over 8 pages as**
91 **long as the subsequent ones contain *only* cited references.**

References

- [1] Chanwoo Kim, Minkyu Shin, Abhinav Garg, and Dhananjaya N. Gowda. Improved vocal tract length perturbation for a state-of-the-art end-to-end speech recognition system. Sept. 2019 (submitted).
- [2] Chanwoo Kim and Richard. M. Stern. Power-Normalized Cepstral Coefficients (PNCC) for Robust Speech Recognition. *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, pages 1315–1329, July 2016.
- [3] William Chan, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4960–4964, March 2016.
- [4] Chanwoo Kim, A. Misra, K.K. Chin, T. Hughes, A. Narayanan, T. Sainath, and M. Bacchiani. Generation of simulated utterances in virtual rooms to train deep-neural networks for far-field speech recognition in Google Home. In *INTERSPEECH-2017*, pages 379–383, Aug. 2017.
- [5] Chung-Cheng Chiu, Tara N. Sainath, Yonghui Wu, Rohit Prabhavalkar, Patrick Nguyen, Zhifeng Chen, Anjuli Kannan, Ron J. Weiss, Kanishka Rao, Ekaterina Gonina, Navdeep Jaitly, Bo Li, Jan Chorowski, and Michiel Bacchiani. State-of-the-art speech recognition with sequence-to-sequence models. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4774–4778, April 2018.
- [6] Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. Attention-based models for speech recognition. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems 28*, pages 577–585. Curran Associates, Inc., 2015.
- [7] D. Yu, M. L. Seltzer, J. Li, J.-T. Huang, and F. Seide. Feature learning in deep neural networks - studies on speech recognition tasks. In *Proceedings of the International Conference on Learning Representations*, 2013.
- [8] Alex Graves. Sequence transduction with recurrent neural networks. *CoRR*, abs/1211.3711, 2012.
- [9] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine Learning, ICML '06*, pages 369–376, New York, NY, USA, 2006. ACM.
- [10] Alex Graves, Abdel rahman Mohamed, and Geoffrey Hinton. Speech recognition with deep recurrent neural networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 6645–6649, May 2013.
- [11] Yanzhang He, Tara N. Sainath, Rohit Prabhavalkar, Ian McGraw, Raziell Alvarez, Ding Zhao, David Rybach, Anjuli Kannan, Yonghui Wu, Ruoming Pang, Qiao Liang, Deepti Bhatia, Yuan Shangguan, Bo Li, Golan Pundak, Khe Chai Sim, Tom Bagby, Shuo yin Chang, Kanishka Rao, and Alexander Gruenstein. Streaming end-to-end speech recognition for mobile devices. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6381–6385, May 2019.
- [12] Geoffrey Hinton, Li Deng, Dong Yu, George E. Dahl, Abdel rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N. Sainath, and Brian Kingsbury. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, Nov 2012.
- [13] Xuedong Huang, Alex Acero, and Hsiao-Wuen Hon. *Spoken Language Processing: A Guide to Theory, Algorithm, and System Development*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition, 2001.
- [14] M. Seltzer, D. Yu, and Y.-Q. Wang. An investigation of deep neural networks for noise robust speech recognition. In *Int. Conf. Acoust. Speech, and Signal Processing*, pages 7398–7402, 2013.
- [15] Rohit Prabhavalkar, Kanishka Rao, Tara N. Sainath, Bo Li, Leif Johnson, and Navdeep Jaitly. A comparison of sequence-to-sequence models for speech recognition. In *Proc. Interspeech 2017*, pages 939–943, 2017.
- [16] Tara N. Sainath, Ron J. Weiss, Kevin W. Wilson, Bo Li, Arun Narayanan, Ehsan Variani, Michiel Bacchiani, Izhak Shafran, Andrew Senior, Kean Chin, Ananya Misra, and Chanwoo Kim. Multichannel signal processing with deep neural networks for automatic speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(5):965–979, May 2017.
- [17] Tara N. Sainath, Ron J. Weiss, Kevin W. Wilson, Arun Narayanan, Michiel Bacchiani, Bo Li, Ehsan Variani, Izhak Shafran, Andrew Senior, Kean Chin, Ananya Misra, and Chanwoo Kim.

- 144 [18] Vincent Vanhoucke, Andrew Senior, and Mark Z. Mao. improving the speed of neural networks on cpus.
145 In *Deep Learning and Unsupervised Feature Learning Workshop, NIPS 2011*, 2011.
- 146 [19] X. Wang, L. C. W. Pols, and L. F. M Bosh. Integration of context-dependent durational knowledge into
147 hmm-based speech recognition. In *ICSLP-1996*, pages 1073–1076, Oct. 1996.
- 148 [20] B. Ziółko and M. Ziółko. Time durations of phonemes in polish language for speech and speaker
149 recognition. In Zygmunt Vetulani, editor, *Human Language Technology. Challenges for Computer Science*
150 *and Linguistics*, pages 105–114, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.