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

21 1 Introduction

22 Recently, deep neural network has significantly improved the performance of speech recognition
23 systems [].

24 2 Review of end-to-end speech recognition systems

25 Recently, there have been growing interests in end-to-end speech recognition systems. [CITE] Fig.
26 1 shows the entire structure of the Connectionist Temporal Classification (CTC) and the attention-
27 based end-to-end speech recognition system.

28 3 Sequence Homogeneously-Structured Classifier

29 Figure ?? shows the entire structure of the SHC. Speech-recognition task can be also considered as
30 a

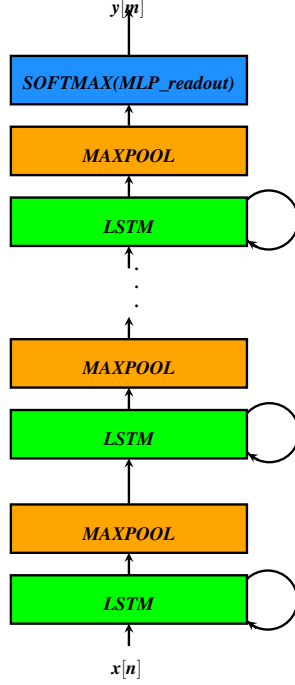


Figure 1: The entire structure of the sequence homogeneously-structured classifier (SHC).

3.1 The structure of the Sequence Homogeneously-Structured Classifier

Fig. shows the block diagram for enhancing the noisy feature.

Fig. 2 shows the structure of the SHC algorithm. The structure is very simple and somewhat similar to that of CTC shown in TODO

3.2 Training of the Sequence Homogeneously-Structured Classifier

In this section, we describe how to train SHC. The procedure is summarized in the following table.

3.3 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.

```

44  for  $m = 0, \dots, M - 1$  do
45    for  $l = 0, \dots, N_b - 1$  do
46       $\pi^{(l)}[m]$ 
47    end for
48  end for
49  if  $i \geq maxval$  then
50     $i \leftarrow 0$ 
51  else
52    if  $i + k \leq maxval$  then
53       $i \leftarrow i + k$ 
54    end if
55  end if

```

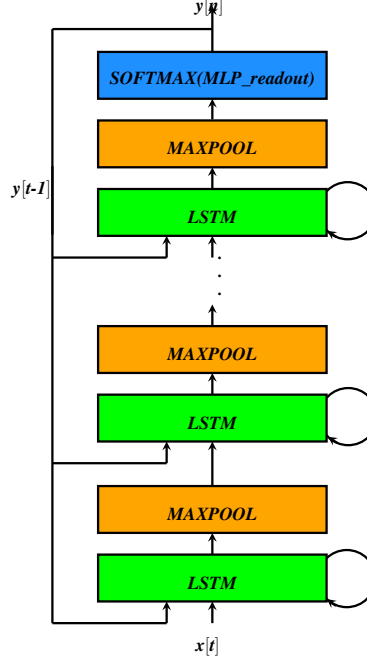


Figure 2: The entire structure of the sequence homogeneously-structured classifier (SHC).

3.4 Feature extraction

Even though the log-mel coefficients have been widely used as features for speech recognition [5], the log non-linearity has a disadvantage in that the value diverges to negative infinity as the mel filterbank coefficient approaches zero [1]. Thus, we use the power-law nonlinearity rather than the log-nonlinearity:

$$q[m, l] = (p[m, l])^{\frac{1}{15}} \quad (1)$$

We use the power coefficient of $\frac{1}{15}$ as in [1, 3].

In this section, we discuss how to predict the clean feature and the error ratio from the corrupt feature. Let us denote the clean target feature, the corrupt feature by \mathbf{x}, \mathbf{y} respectively.

Using the room simulation system described in [2], we create a pair of clean and corrupt utterances. Log-mel [XX] features are calculated from these original clean and corrupt utterances. Let us denote the clean feature by \mathbf{t}_i and the simulated corrupt feature by \mathbf{x}_i , respectively where i is the utterance index. The training set is represented by the following set:

$$\mathcal{T} = \{ \langle \mathbf{x}_i, \mathbf{t}_i \rangle \mid 0 \leq i \leq N - 1 \} \quad (2)$$

where N is the number of training examples.

In the ECE algorithm, we first estimate the true target \mathbf{t}_i . The estimation error is defined by the following equation:

$$\mathbf{e}_i = \mathbf{t}_i - \mathbf{x}_i, \quad 0 \leq i \leq N - 1 \quad (3)$$

The norm of the error vector \mathbf{e}_i would be generally large when the norm of the estimated clean feature is large.

Thus, instead of trying to estimate the error itself, the neural network tries to estimate the error ratio given as follows:

$$r_i = \quad (4)$$

75 4 Generation of simulated clean and noisy data sets

76 To train neural networks described in section ??, we need a training set consisting of pairs of clean
 77 speech and noisy speech. Unfortunately, it is not easy to have such a training set from real speech
 78 utterances. Thus, we use the simulation system described in [2] to synthetically generate noisy
 79 speech utterances from the clean speech utterances.

80 In our application, for clean training set, we used Wall Street Journal (WSJ) si-284. For noise set, we
 81 used the Resource Management 1 (RM1) [4] 50 % of time and DEMAND noise [TODO(chanwcom)
 82 Adds reference] for the remaining 50 % of time.

83 TODO(chanwcom) Cite DEMAND noise set.

84 5 Estimation of the optimal weight in the EFW algorithm

85 In this section, we describe the procedure of feature weighting approach. Suppose that the feature
 86 before enhancement is represented by \vec{x} . Let us represent the inference neural network to estimate
 87 the target by \mathcal{T} .

$$\mathcal{X} = \{\}$$
 (5)

88 Instead of directly using the enhanced feature $\hat{\vec{y}}$, we consider the following interpolated feature \vec{z}
 89 from the original corrupt input \vec{x} and .

$$\vec{z} = \vec{w} \odot \hat{\vec{y}} + (1 - \vec{w}) \odot \vec{x}$$
 (6)

90 where \odot denotes the Hadamard product (entry-wise product).

91 Let us denote the estimated variance vector by $\hat{\vec{v}}$:

$$\begin{aligned} \hat{\vec{v}} &= \text{Var}[\hat{\vec{y}}] \\ &= E[(\hat{y} - y)^2 | \vec{x}[0], \vec{x}[1], \dots, \vec{x}[T]] \end{aligned}$$
 (7)

92 In our discussion, let us assume that the expectation of $\hat{\vec{y}}$ is the same as the true target \vec{y} .

$$E[\hat{\vec{y}}] = \vec{y}$$
 (8)

93 Now, let us obtain the expected value of \vec{z} from (6) and (8)

$$E[\vec{z}] = (1 - \vec{w}) \odot (\vec{x} - \vec{y})$$
 (9)

94 Thus, the bias of \vec{z} is given by:

$$\begin{aligned} \text{Bias}_{\vec{z}} &= E[\vec{z}] - \vec{t} \\ &= \vec{w} \odot (\hat{\vec{y}} - \vec{y}) \end{aligned}$$
 (10)

95 From (6) and (7), the variance of \vec{z} is given by:

$$\text{Var}_{\vec{z}} = \vec{v} \odot \vec{w} \odot \vec{w}$$
 (11)

96 The mean squared error of the i -th element of the \vec{z} is given by:

$$\begin{aligned} \text{MSE}_{\vec{z}_i} &= \text{Bias}_{\vec{z}_i}^2 + \text{Var}_{\vec{z}_i} \\ &= [(\vec{x}_i - \vec{y}_i)^2 + \vec{v}_i^2] \vec{w}_i^2 - 2(\vec{x}_i - \vec{y}_i)^2 \vec{w}_i + (\vec{x}_i - \vec{y}_i)^2 \end{aligned}$$
 (12)

97 Since the above equation (12) is a quadratic equation with respect to \vec{w}_i , we may find that the
 98 minimum value of $\text{MSE}_{\vec{z}_i}$ is obtained when \vec{w}_i has the following value:

$$\vec{w}_i = \frac{(\vec{x}_i - \vec{y}_i)^2}{(\vec{x}_i - \vec{y}_i)^2 + \vec{v}_i}$$
 (13)

99 Thus, the final form of the estimated vector $\hat{\vec{z}}$ becomes:

$$\hat{\vec{z}} = \hat{\vec{w}} \odot \hat{\vec{y}} + (1 - \hat{\vec{w}}) \odot \vec{x} \quad (14)$$

100 where

$$\hat{\vec{w}} = (\vec{x} - \vec{y})^{\odot 2} [(\vec{x} - \vec{y})^{\odot 2} + \vec{v}]^{\odot -1}. \quad (15)$$

101 Sequence enhancement

102 The feature enhancement is to map a sequence of corrupt features \mathbf{X} .

$$\mathbf{X} = (\vec{x}[0], \vec{x}[1], \dots, \vec{x}[M-1]) \quad (16a)$$

$$\mathbf{Y} = (\vec{y}[0], \vec{y}[1], \dots, \vec{y}[M-1]) \quad (16b)$$

103 **Acknowledgments**

104 **References**

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

References

- [1] C. Kim and R. M. Stern. Power-Normalized Cepstral Coefficients (PNCC) for Robust Speech Recognition. *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, pages 1315–1329, July 2016.
- [2] C. 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.
- [3] C. Kim and R. M. Stern. Power-normalized cepstral coefficients (pncc) for robust speech recognition. In *IEEE Int. Conf. on Acoustics, Speech, and Signal Processing*, pages 4101–4104, March 2012.
- [4] Price, P, et al. *Resource Management RM1 2.0 LDC93S3B. DVD*. Linguistic Data Consortium, Philadelphia, PA, 1993.
- [5] S. B. Davis and P. Mermelstein. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Trans. Acoust., Speech, and Signal Processing*, 28(4):357–366, Aug. 1980.