

# Supplementary material for “A Purely End-to-end System for Multi-speaker Speech Recognition”

Hiroshi Seki<sup>1,2,\*</sup>, Takaaki Hori<sup>1</sup>, Shinji Watanabe<sup>3</sup>, Jonathan Le Roux<sup>1</sup>, John R. Hershey<sup>1</sup>

<sup>1</sup>Mitsubishi Electric Research Laboratories (MERL)

<sup>2</sup>Toyohashi University of Technology

<sup>3</sup>Johns Hopkins University

## 1 Architecture of the encoder-decoder network

In this section, we describe the details of the baseline encoder-decoder network which is further extended for permutation-free training. The encoder network consists of a VGG network and bi-directional long short-term memory (BLSTM) layers. The VGG network has the following 6-layer CNN architecture at the bottom of the encoder network:

- Convolution (# in = 3, # out = 64, filter =  $3 \times 3$ )
- Convolution (# in = 64, # out = 64, filter =  $3 \times 3$ )
- MaxPooling (patch =  $2 \times 2$ , stride =  $2 \times 2$ )
- Convolution (# in = 64, # out = 128, filter =  $3 \times 3$ )
- Convolution (# in = 128, # out = 128, filter =  $3 \times 3$ )
- MaxPooling (patch =  $2 \times 2$ , stride =  $2 \times 2$ )

The first 3 channels are static, delta, and delta delta features. Multiple BLSTM layers with projection layer  $\text{Lin}(\cdot)$  are stacked after the VGG network. We defined one BLSTM layer as the concatenation of a forward LSTM  $\overrightarrow{\text{LSTM}}(\cdot)$  and a backward LSTM  $\overleftarrow{\text{LSTM}}(\cdot)$ :

$$\overrightarrow{H} = \overrightarrow{\text{LSTM}}(\cdot) \quad (29)$$

$$\overleftarrow{H} = \overleftarrow{\text{LSTM}}(\cdot) \quad (30)$$

$$H = [\text{Lin}(\overrightarrow{H}); \text{Lin}(\overleftarrow{H})], \quad (31)$$

When the VGG network and the multiple BLSTM layers are represented as  $\text{VGG}(\cdot)$  and  $\text{BLSTM}(\cdot)$ , the encoder network in Eq. (2) maps the input feature vector  $O$  to internal representation  $H$  as follows:

$$H = \text{Encoder}(O) = \text{BLSTM}(\text{VGG}(O)) \quad (32)$$

\*This work was done while H. Seki, Ph.D. candidate at Toyohashi University of Technology, Japan, was an intern at MERL.

The decoder network sequentially generates the  $n$ -th label  $y_n$  by taking the context vector  $c_n$  and the label history  $y_{1:n-1}$ :

$$y_n \sim \text{Decoder}(c_n, y_{1:n-1}). \quad (33)$$

The context vector is calculated in an location based attention mechanism (Chorowski et al., 2015) which weights and sums the  $C$ -dimensional sequence of representation  $H = (h_l \in \mathbb{R}^C | l = 1, \dots, L)$  with attention weight  $a_{n,l}$ :

$$c_n = \text{Attention}(a_{n-1}, e_n, H), \quad (34)$$

$$\triangleq \sum_{l=1}^L a_{n,l} h_l. \quad (35)$$

The location based attention mechanism defines the weights  $a_{n,l}$  as follows:

$$a_{n,l} = \frac{\exp(\alpha k_{n,l})}{\sum_{l=1}^L \exp(\alpha k_{n,l})}, \quad (36)$$

$$k_{n,l} = w^T \tanh(V^E e_{n-1} + V^H h_l + V^F f_{n,l} + b), \quad (37)$$

$$f_n = F * a_{n-1}, \quad (38)$$

where  $w, V^E, V^H, V^F, b, F$  are tunable parameters,  $\alpha$  is a constant value called inverse temperature, and  $*$  is the convolution operation. We used 10 convolution filters of width 200, and set  $\alpha$  to 2. The introduction of  $f_n$  makes the attention mechanism take into account the previous alignment information. The hidden state  $e$  is updated recursively by an updating LSTM function:

$$e_n = \text{Update}(e_{n-1}, c_{n-1}, y_{n-1}), \quad (39)$$

$$\triangleq \text{LSTM}(\text{Lin}(e_{n-1}) + \text{Lin}(c_{n-1}) + \text{Emb}(y_{n-1})), \quad (40)$$

where  $\text{Emb}(\cdot)$  is an embedding function.

Table 1: Examples of recognition results. Errors are emphasized as capital letter. “\_” is a space character, and a special token “\*” is inserted to pad deletion errors.

---

**(1) Model w/ permutation-free training (CER of HYP1: 12.8%, HYP2: 0.9%)**

**HYP1:** the\_shuttle\_\*\*\*IS\_IN\_the\_first\_tHE\_lifE\_o\*f\_since\_the\_nineteen\_eight\_y\_six\_challenger\_explosion

**REF1:** the\_shuttle\_WOULD\_BE\_the\_first\_t\*O\_lifT\_oFf\_since\_the\_nineteen\_eigh\_tty\_six\_challenger\_explosion

**HYP2:** the\_expanded\_recall\_was\_disclosed\_at\_a\_meeting\_with\_n.r.c.officia\_ls\_at\_an\_agency\_office\_outside\_chicago

**REF2:** the\_expanded\_recall\_was\_disclosed\_at\_a\_meeting\_with\_n.r.c.officia\_ls\_at\_an\_agency\_office\_outside\_chicago

---

**(2) Model w/ permutation-free training (CER of HYP1: 91.7%, HYP2: 38.9%)**

**HYP1:** IT\_WAS\_Last\_r\*AISe\*D\_IN\_JUNE\_NINeTeeN\_e\*IGHtY\_fIVe\_TO\_\*THIRTY

**REF1:** \*\*\*\*\*ast\*ronOMeRS\_SAY\_THAT\_\*\*\*\*tHe\*\_eARTh'S\_fATe\_IS\_SEALED

**HYP2:** \*\*\*\*aND\_\*st\*rongeRS\_SAY\_THAT\_\*\*\*\*tHe\*\_e\*ARtH\_fATe\_IS\_to\_fo\_rty\_five\_dollars\_from\_thirty\_five\_dollars

**REF2:** IT\_Wa\*S\_LAst\_rAISe\*D\_IN\_JUNE\_NINeTeeN\_eIGHtY\_fIVe\*\*\*\_to\_fort\_y\_five\_dollars\_from\_thirty\_five\_dollars

---



---

**Algorithm 1** Generation of multi speaker speech dataset

---

$n_{\text{reuse}} \leftarrow$  maximum number of times same utterance can be used.

$U \leftarrow$  utterance set of the corpora.

$C_k \leftarrow n_{\text{reuse}}$  for each utterance  $U_k \in U$

**for**  $U_k \in U$  **do**

$P(U_k) = C_k / \sum_l C_l$

**end for**

**for**  $U_i$  in  $U$  **do**

Sample utterance  $U_j$  from  $P(U)$  while ensuring speakers of  $U_i$  and  $U_j$  are different.

Mix utterances  $U_i$  and  $U_j$

**if**  $C_j > 0$  **then**

$C_j = C_j - 1$

**for**  $U_k \in U$  **do**

$P(U_k) = C_k / \sum_l C_l$

**end for**

**end if**

**end for**

---

## 2 Generation of mixed speech

Each utterance of the corpus is mixed with a randomly selected utterance with the probability,  $P(U_k)$ , that moderates over-selection of specific utterances.  $P(U_k)$  is calculated in the first for-loop as a uniform probability. All utterances are used as one side of the mixture, and another side is sam-

pled from the distribution  $P(U_k)$  in the second for-loop. The selected pairs of utterances are mixed at various signal-to-noise ratios (SNR) between 0 dB and 5 dB. We randomized the starting point of the overlap by padding the shorter utterance with silence whose duration is sampled from the uniform distribution within the length difference between the two utterances. Therefore, the duration of the mixed utterance is equal to that of the longer utterance among the unmixed speech. After the generation of the mixed speech, the count of selected utterances  $C_j$  is decremented to prevent of over-selection. All counts  $C$  are set to  $n_{\text{reuse}}$ , and we used  $n_{\text{reuse}} = 3$ .

## 3 Examples of recognition results and error analysis

Table 1 shows examples of recognition result. The first example (1) is one which accounts for a large portion of the evaluation set. The SNR of the HYP1 is -1.55 db and that of HYP2 is 1.55 dB. The network generates multiple hypotheses with a few substitution and deletion errors, but without any overlapped and swapped words. The second example (2) is one which leads to performance reduction. We can see that the network makes errors when there is a large difference in length between the two sequences. The word “thirty” of HYP2 is injected in HYP1, and there are deletion errors in

HYP2. We added a negative KL divergence loss to ease such kind of errors. However, there is further room to reduce error by making unshared modules more cooperative.

## References

Jan K Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, and Yoshua Bengio. 2015. Attention-based models for speech recognition. In *Advances in Neural Information Processing Systems (NIPS)*, pages 577–585.