

## A Block diagram of hierarchical character-level language model

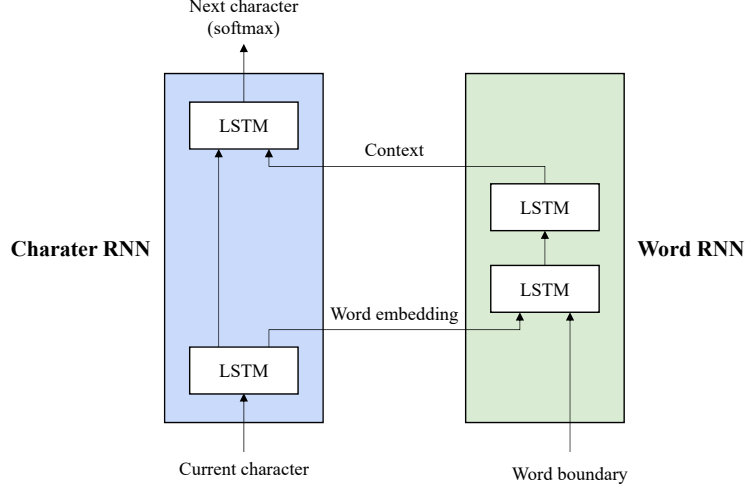


Figure 1: The architecture of the hierarchical recurrent neural network model used for language modeling.

## B Formulation of recurrent units

The recurrent units used in this work are described as follows:

LSTM:

$$\begin{aligned}
 \mathbf{f}_t &= \sigma(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f), \\
 \mathbf{i}_t &= \sigma(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i), \\
 \mathbf{o}_t &= \sigma(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o), \\
 \hat{\mathbf{c}}_t &= \tanh(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c), \\
 \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \hat{\mathbf{c}}_t, \\
 \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\mathbf{c}_t).
 \end{aligned} \tag{1}$$

GRU:

$$\begin{aligned}
 \mathbf{z}_t &= \sigma(\mathbf{W}_z \mathbf{x}_t + \mathbf{U}_z \mathbf{h}_{t-1} + \mathbf{b}_z), \\
 \mathbf{r}_t &= \sigma(\mathbf{W}_r \mathbf{x}_t + \mathbf{U}_r \mathbf{h}_{t-1} + \mathbf{b}_r), \\
 \mathbf{h}_t &= (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tanh(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h (\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{b}_h).
 \end{aligned} \tag{2}$$

GILR-LSTM:

$$\begin{aligned}
 \mathbf{g}_t &= \sigma(\mathbf{V}_g \mathbf{x}_t + \mathbf{b}_g), \\
 \mathbf{j}_t &= \tanh(\mathbf{V}_j \mathbf{x}_t + \mathbf{b}_j), \\
 \tilde{\mathbf{h}}_t &= \mathbf{g}_t \odot \tilde{\mathbf{h}}_{t-1} + (1 - \mathbf{g}_t) \odot \mathbf{j}_t, \\
 [\mathbf{f}_t, \mathbf{i}_t, \mathbf{o}_t] &= \sigma\left([\mathbf{U}_f, \mathbf{U}_i, \mathbf{U}_o] \tilde{\mathbf{h}}_{t-1} + [\mathbf{V}_f, \mathbf{V}_i, \mathbf{V}_o] \mathbf{x}_t + [\mathbf{b}_f, \mathbf{b}_i, \mathbf{b}_o]\right), \\
 \mathbf{z}_t &= \tanh\left(\mathbf{U}_z \tilde{\mathbf{h}}_{t-1} + \mathbf{V}_z \mathbf{x}_t + \mathbf{b}_z\right), \\
 \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{z}_t, \\
 \mathbf{h}_t &= \mathbf{o}_t \odot \mathbf{c}_t.
 \end{aligned} \tag{3}$$

## C Decoding algorithm

**Input:** AM output probability matrix  $P_{\text{CTC}} \in \mathbb{R}^{|\Sigma| \times T}$ , beam width  $B$ , number of search candidates  $k$ , language model weight  $\alpha$ , insertion bonus  $\gamma$ , vocabulary  $\Sigma$

```

1  $\mathbf{A}_{prev} = \{\phi\}$ 
2  $P_{\text{CTC}}(\text{blank}|\mathbf{x}_{1:0}) = 1$ 
3 for  $t = 1$  to  $T$  do
4   if  $P_{\text{CTC}}(\text{blank}|\mathbf{x}_{1:t-1}) > 0.95$  and  $P_{\text{CTC}}(\text{blank}|\mathbf{x}_{1:t}) > 0.95$  then
5     continue
6   end
7    $\mathbf{A}_{next} = \{\}$ 
8    $\mathbf{K} = \text{top-}k \text{ labels in } \Sigma \text{ according to value of } P_{\text{CTC}}(\mathbf{c}|\mathbf{x}_{1:t})$ 
9   for  $l$  in  $\mathbf{A}_{prev}$  do
10    for  $c$  in  $\mathbf{K}$  do
11      if  $c = \text{blank}$  then
12         $p_{nb}(l) = p_{nb}(l)P_{\text{CTC}}(\text{blank}|\mathbf{x}_{1:t})$ 
13         $p_b(l) = (p_{nb}(l) + p_b(l))P_{\text{CTC}}(\text{blank}|\mathbf{x}_{1:t})$ 
14      else
15         $l^+ = \text{concat}(l, c)$ 
16        if  $c = l_{end}$  then
17           $p_{nb}(l^+) = p_b(l)P_{\text{CTC}}(c|\mathbf{x}_{1:t})\gamma P_{LM}(c|l)^\alpha$ 
18           $p_{nb}(l) = p_b(l)P_{\text{CTC}}(c|\mathbf{x}_{1:t})$ 
19        else
20           $p_{nb}(l^+) = (p_b(l) + p_{nb}(l))P_{\text{CTC}}(c|\mathbf{x}_{1:t})\gamma P_{LM}(c|l)^\alpha$ 
21        end
22      end
23      add  $l^+$  to  $\mathbf{A}_{next}$ 
24    end
25  end
26  assign top- $B$  of  $\mathbf{A}_{next}$  to  $\mathbf{A}_{prev}$ 
27 end

```

**Algorithm 1:** Prefix beam search in proposed system.

The most time-consuming part in the decoding is computing the probability of  $P_{LM}(c|l)$ , which is required in the line 17 and 20. Time complexity of the LM computation is  $O(B \times T \times |\Sigma|)$ , but the actual computation complexity can be reduced to  $O(B \times |l| \times |\Sigma|)$  by reusing the result of LM for same inputs. Skipping consecutive blanks and candidate pruning are applied in line number 5 and 8, respectively. Table 1 shows the ratio of skipped repeated blank frames. The number of operations of LM per frame according to the number of candidates are shown in Table 2.

Table 1: The ratio of frames whose decoding stages are skipped due to high CTC blank output.

Acoustic model	Downsampling ratio	Percentage of skipping
WSJ - Character	$\times 2$	33.8%
WSJ - Word piece	$\times 4$	44.33%
Libri - Word piece	$\times 8$	20.23%

Table 2: The number of LM operations with the varying number of candidates.

Number of candidates	20	30	40	100
LM operations / frame	4.365	4.488	4.593	4.820

## D Training curves of the acoustic models

We compared the train and valid loss curves of some selected models. The peaks in training curves are due to curriculum-like learning scheduling. SRU without 1-D convolution was not trained well. LSTM was converged faster than other models but it reached local minimum quickly. Training loss of i-SRU was reduced faster than SRU, while they showed the similar valid loss in the end of training.

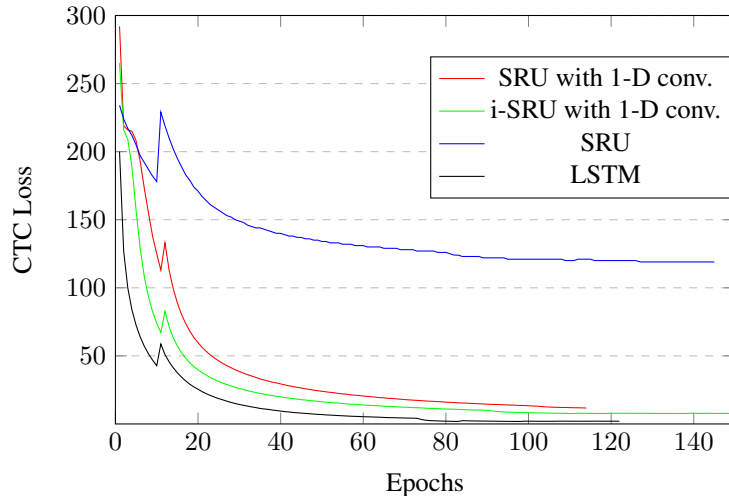


Figure 2: Training loss of acoustic models when trained on WSJ SI-284.

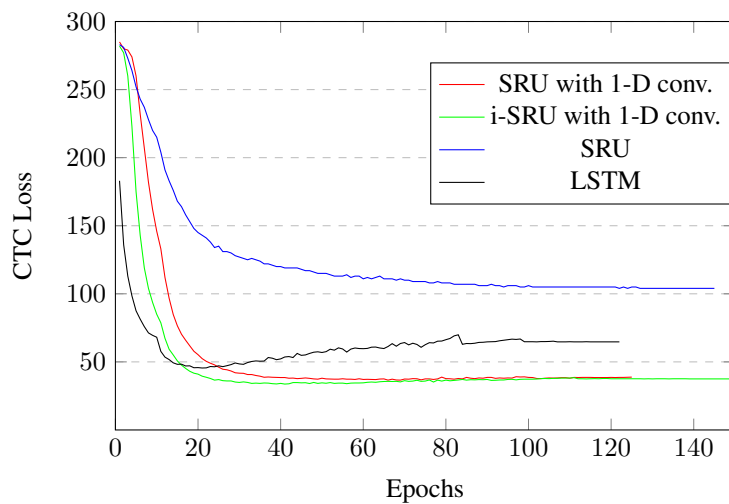


Figure 3: Validation loss of acoustic models when trained on WSJ SI-284

## E Example results of word piece and character level speech recognition

Table 3: Examples which are correct in wordpiece-level model but wrong in character-level model.

Label	THE PINK SHEET A DRUG INDUSTRY TRADE LETTER REPORTED
char-greedy char-LM	THE PANKSHET A DRUG INDUSTRY TRADE LETTER REPORTED THE BANK HIT A DRUG INDUSTRY TRADE LETTER REPORTED
wp-greedy wp-LM	THE PINK SHEET A DRUG INDUSTRY TRADE LETTER REPORTED THE PINK SHEET A DRUG INDUSTRY TRADE LETTER REPORTED
Label	THIS WEEK LOCAL GOVERNMENTS HAVE APPEARED IN THE...
char-greedy char-LM	THIS WEEK LOCAL GOVERNMENTS HAVE THE PEAR IN THE... THIS WEEK LOCAL GOVERNMENTS HAVE THE PART IN THE...
wp-greedy wp-LM	THIS WEEK LOCAL GOVERNMENTS HAVE APPEARED IN THE... THIS WEEK LOCAL GOVERNMENTS HAVE APPEARED IN THE...

Table 4: Examples that are correct in character-level model but wrong in word piece-level model.

Label	SO MR. WANG TELLS PEOPLE HE IS FIFTY
char-greedy char-LM	SO MR. WEANGTELL'S PEOPLE HE IS FIFTY SO MR. WANG TELLS PEOPLE HE IS FIFTY
wp-greedy wp-LM	SO MR. WGTELLES PEOPLE HE IS FIFTY SO MR. WANGTELL'S PEOPLE HE IS FIFTY
Label	THE MARKETS TEND TO MAGNIFY THE NEWS
char-greedy char-LM	THE MARKETS TEND DOMAGNIFY THE NEWS THE MARKETS TEND TO MAGNIFY THE NEWS
wp-greedy wp-LM	THE MARKET'S TEND TO MAGNIFY THE NEWS THE MARKET'S TEND TO MAGNIFY THE NEW

## F Example of results that are corrected when HCLM is used

Table 5: Example sentences that are corrected when HCLM is used.

Label	RESPONSES TEND TO BE MUTE
greedy LSTM-CLM HCLM	RESPONSEES TEND TO BE MUNED RESPONSE IS TEND TO BE MUTED RESPONSES TEND TO BE MUTED
Label	I'M NOT AT ALL UNHAPPY WITH WHAT I'M SEEING
greedy LSTM-CLM HCLM	I'M NOT IT ALL AND HAPPY WITH WHAT I'M SEE I'M NOT AT ALL AND HAPPY WITH WHAT I'M SEEN I'M NOT AT ALL UNHAPPY WITH WHAT I'M SEEING
Label	THE BIG SHOE IS GOING TO DROP WHEN WE SEE THE TRADE NUMBER
greedy LSTM-CLM HCLM	THE BIG SHE WAS GOING TO DROP INLY SEE THE TRADE NUMBER THE BIG SHE WAS GOING TO DROP IN WE SEE THE TRADE NUMBER THE BIG SHOE IS GOING TO DROP WHEN WE SEE THE TRADE NUMBER
Label	AND THEY BALK AT THE APPROACH USED IN MEAT AND POULTRY PLANTS...
greedy LSTM-CLM HCLM	AND THEY BALCK AT THE APPROACH USED IN MEET AND PULTRY PLANTS... AND THEY BALK AT THE APPROACH USED IN MEET AND POULTRY PLANTS... AND THEY BALK AT THE APPROACH USED IN MEAT AND POULTRY PLANTS...