

Towards End-to-End Speech Recognition with Recurrent Neural Networks

Alex Graves, Navdeep Jaitly

Introduction

In previous speech recognition, the neural networks are at present only a single component in a complex pipeline

Shortage of previous speech recognition system:

The frame-level training targets must be inferred from the alignments determined by the HMM

A pronunciation dictionary is necessary to map from words to phoneme sequences

Create pronunciation dictionary costs a lot of labor

The quality of the dictionary affects the result of speech recognition system dramatically



Alignment



Pronunciation dictionary

Improvement of this paper

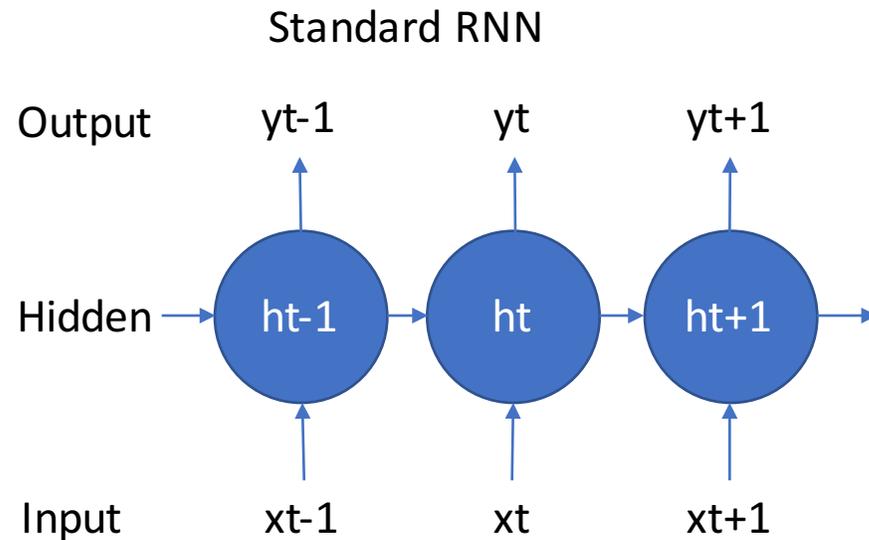
Using a deep bidirectional LSTM network with a Connectionist Temporal Classification output layer

No pronunciation dictionary is necessary

CTC integrates out over all possible input-output alignments, no forced alignment is required to provide training targets

Bidirectional Recurrent Neural Networks

In the speech recognition, we need consider the information from both past and future
 Add backward Layer to Standard RNN to construct Bidirectional RNN



$$h_t = \mathcal{H}(W_{ih}x_t + W_{hh}h_{t-1} + b_h) \quad (1)$$

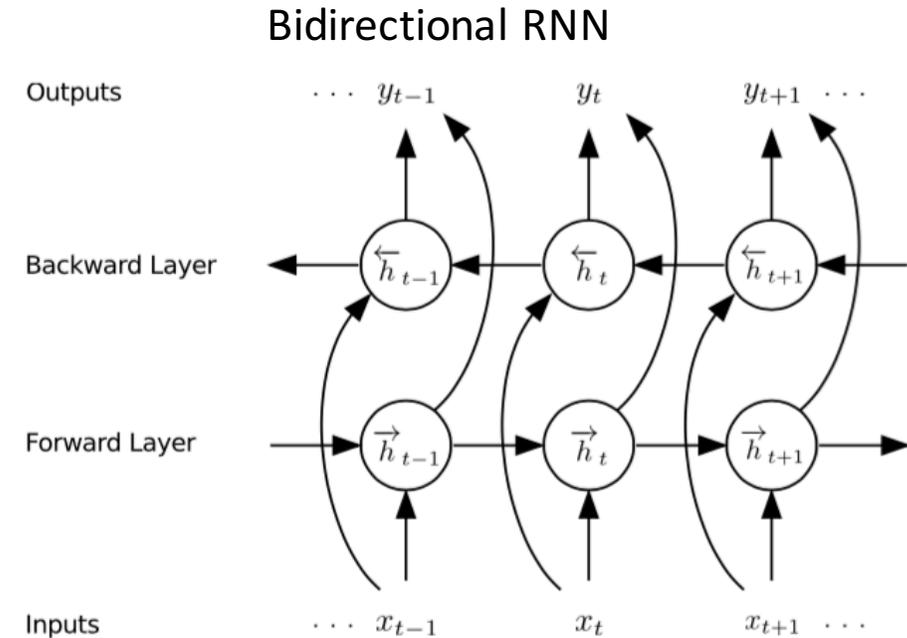
$$y_t = W_{ho}h_t + b_o \quad (2)$$

Input vector $x=(x_1,\dots,x_T)$

Hidden vector $h=(h_1,\dots,h_T)$

Output vector $y=(y_1,\dots,y_T)$

Optimize W_{ih} W_{hh} and W_{ho}



$$\overrightarrow{h}_t = \mathcal{H}(W_{x\overrightarrow{h}}x_t + W_{\overrightarrow{h}\overrightarrow{h}}\overrightarrow{h}_{t-1} + b_{\overrightarrow{h}}) \quad (8)$$

$$\overleftarrow{h}_t = \mathcal{H}(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t+1} + b_{\overleftarrow{h}}) \quad (9)$$

$$y_t = W_{\overrightarrow{h}y}\overrightarrow{h}_t + W_{\overleftarrow{h}y}\overleftarrow{h}_t + b_o \quad (10)$$

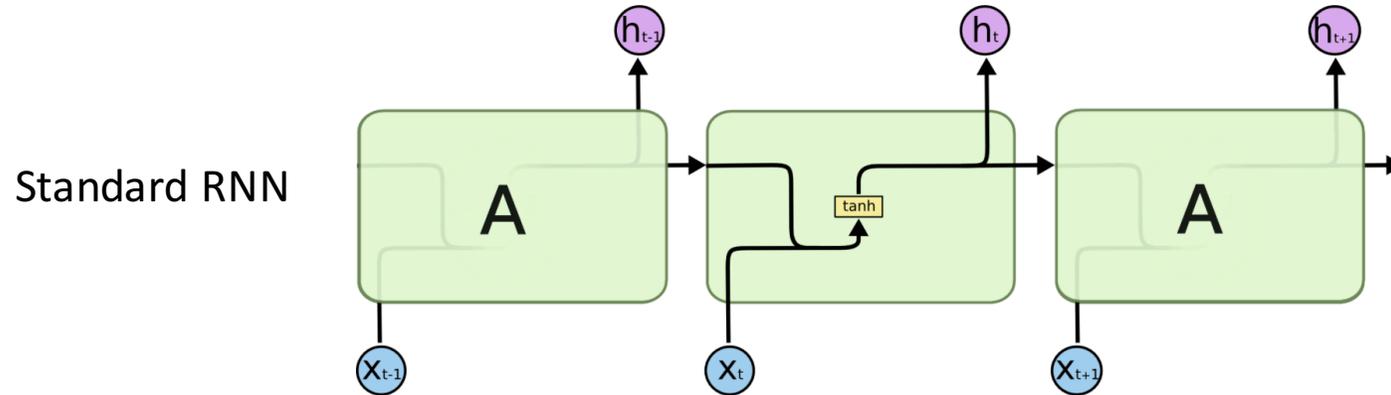
[Graves Towards End-to-End Speech Recognition with Recurrent Neural Networks](http://x-algo.cn/index.php/2016/04/25/rnn-recurrent-neural-networks-derivation-and-implementation/)

<http://x-algo.cn/index.php/2016/04/25/rnn-recurrent-neural-networks-derivation-and-implementation/>

<http://www.shareditor.com/blogshow/?blogId=116>

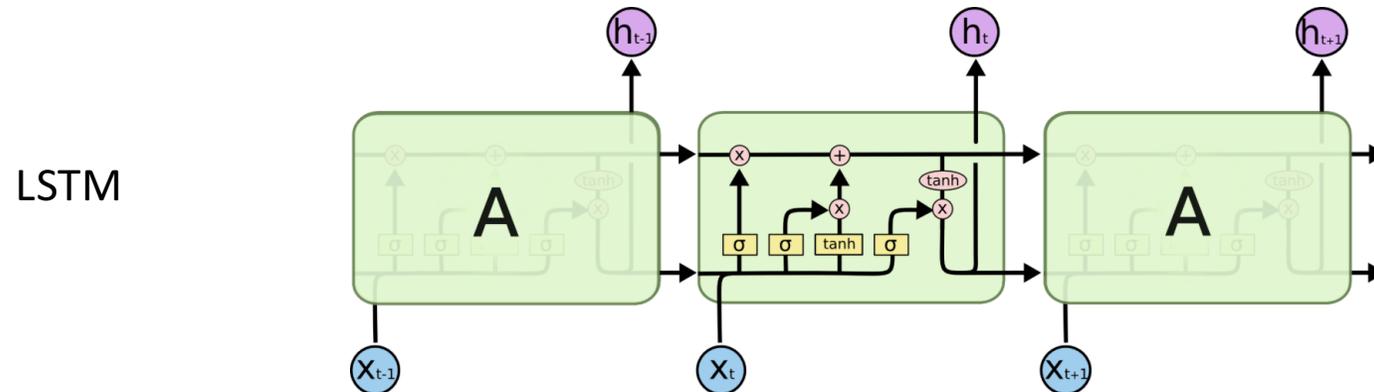
Long Short Term Memory

Standard RNN performs bad for long-term dependencies because gradients propagation over many stages tend to either vanish or explode
In practice, Long Short Term Memory is widely used to preserve long-term dependencies

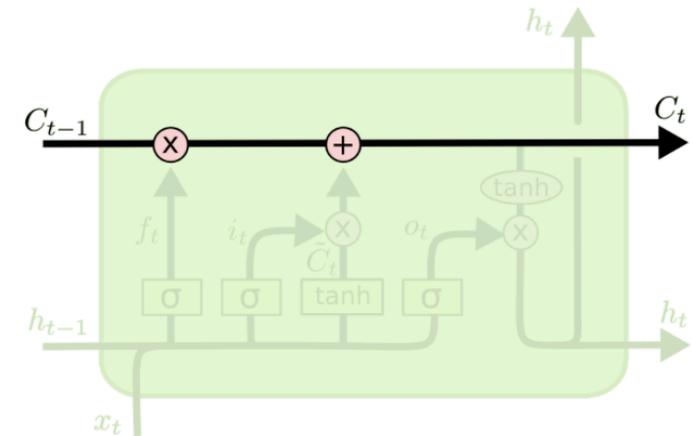


The repeating module in a standard RNN contains a single layer.

The key to LSTMs is the cell state, the horizontal line running through the top of the diagram. It's very easy for information to just flow along it unchanged.

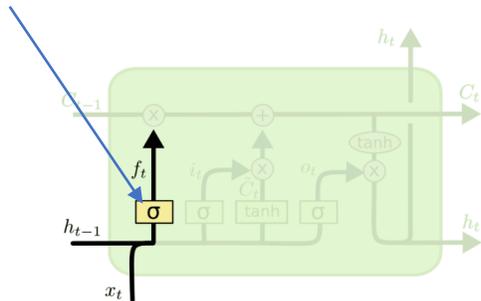


The repeating module in an LSTM contains four interacting layers.



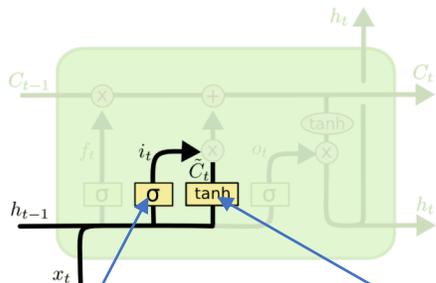
Long Short Term Memory

Forget Gate



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

The first step in our LSTM is to decide what information we're going to throw away from the cell state.



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

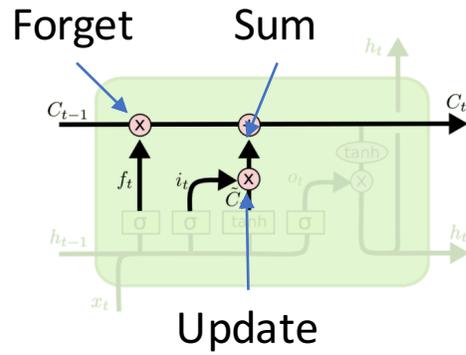
8

The next step is to decide what new information we're going to store in the cell state.

Input gate : which values we'll update

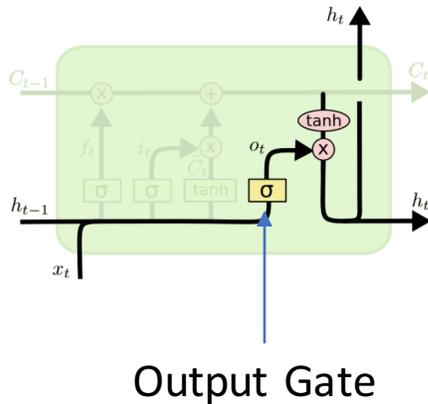
New candidate values, added to the state

Long Short Term Memory



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

Update the old cell state C_{t-1} , into the new cell state C_t



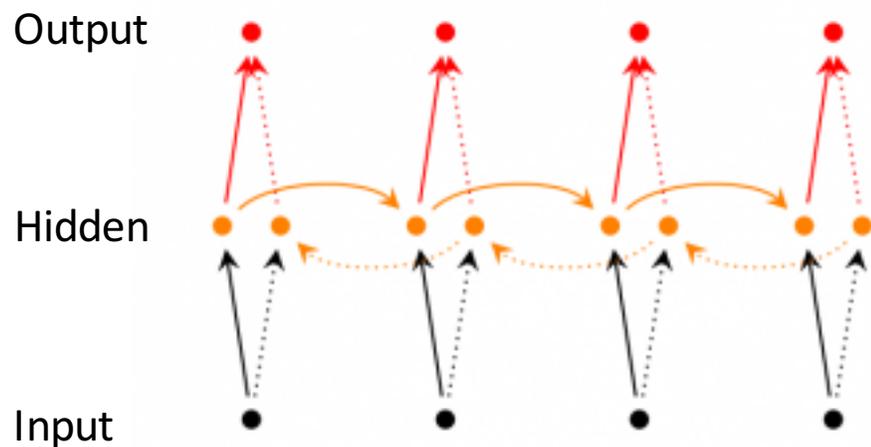
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

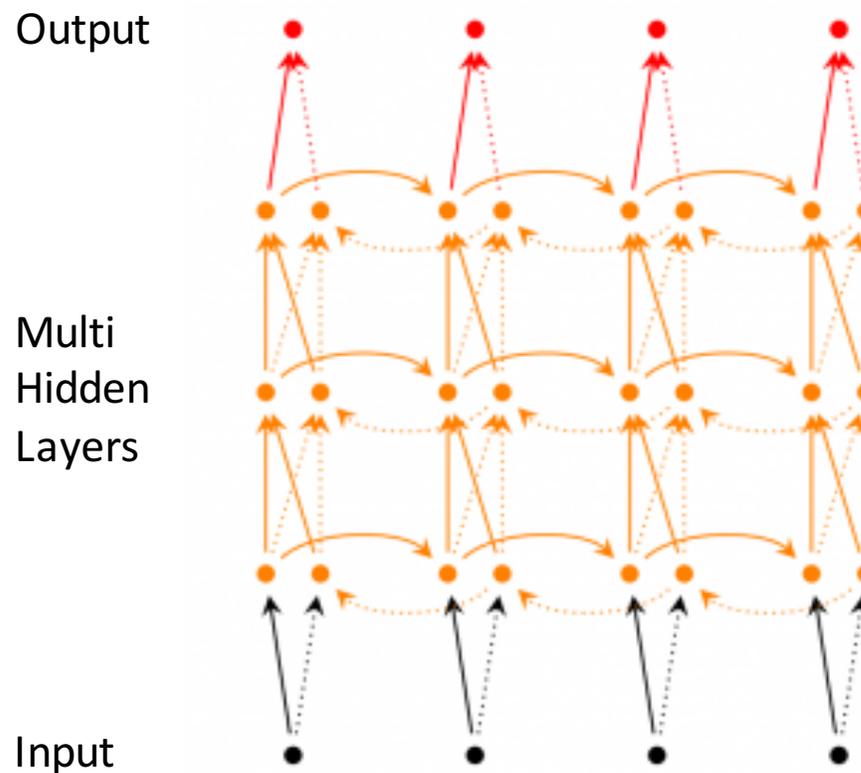
Finally, we need to decide what we're going to output. This output will be based on our cell state, but will be a filtered version.

Deep Recurrent Neural Network

Bidirectional RNN



Deep Bidirectional RNN



More ability to learning

Connectionist Temporal Classification (CTC)

Typical frame-level classifier:

Requires a training target for every frame in the audio

(1) circular dependency:

good alignment \leftrightarrow good classifier

(2) no need for character-level alignment

Using CTC as objective function:

Allows an RNN to be trained without any prior alignment between the input and the transcription

Connectionist Temporal Classification (CTC)

\mathbf{x} : Input sequence with length T

\mathbf{a} : output sequence of blank and label indices (CTC alignment)

\mathbf{y} : output transcription

\mathbf{B} : matching CTC alignment to transcription

Example:

$a_1 = (a, -, b, c, -, -)$ $a_2 = (-, -, a, -, b, c)$

$a_3 = (a, -, b, -, c, c)$ $a_4 = (a, a, b, b, c, c)$

$\mathbf{y} = (a, b, c)$

$B(a_1) = B(a_2) = B(a_3) = B(a_4) = \mathbf{y}$

$$\Pr(\mathbf{y}|\mathbf{x}) = \sum_{\mathbf{a} \in \mathcal{B}^{-1}(\mathbf{y})} \Pr(\mathbf{a}|\mathbf{x})$$

$$CTC(\mathbf{x}) = -\log \Pr(\mathbf{y}^*|\mathbf{x})$$

Expected Transcription Loss

$L(\mathbf{x}, \mathbf{y})$: transcription loss function (e.g. WER)

$L(\mathbf{x})$: expected transcription loss

$$\mathcal{L}(\mathbf{x}) = \sum_{\mathbf{y}} \Pr(\mathbf{y}|\mathbf{x}) \mathcal{L}(\mathbf{x}, \mathbf{y})$$

Approximate with Monte-Carlo sampling

$$\mathcal{L}(\mathbf{x}) \approx \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\mathbf{x}, \mathcal{B}(\mathbf{a}^i)), \quad \mathbf{a}^i \sim \Pr(\mathbf{a}|\mathbf{x})$$

Differentiate $L(\mathbf{x})$ with respect to RNN outputs

$$\frac{\partial \mathcal{L}(\mathbf{x})}{\partial \Pr(k, t|\mathbf{x})} \approx \frac{1}{N} \sum_{i=1}^N \mathcal{L}(\mathbf{x}, \mathcal{B}(\mathbf{a}^{i,t,k}))$$

Expected Transcription Loss

$$\frac{\partial \mathcal{L}(\mathbf{x})}{\partial y_t^k} \approx \frac{\Pr(k, t | \mathbf{x})}{N} \sum_{i=1}^N \mathcal{L}(\mathbf{x}, \mathcal{B}(\mathbf{a}^i)) - \mathcal{Z}(\mathbf{a}^i, t)$$

where

$$\mathcal{Z}(\mathbf{a}^i, t) = \sum \Pr(k', t | \mathbf{x}) \mathcal{L}(\mathbf{x}, \mathcal{B}(\mathbf{a}^{i,t,k'}))$$



The derivative added to y_t^k equals to the difference between the loss with $a_t^i = k$ and the expected loss with a_t^i sampled from $\Pr(k', t | \mathbf{x})$



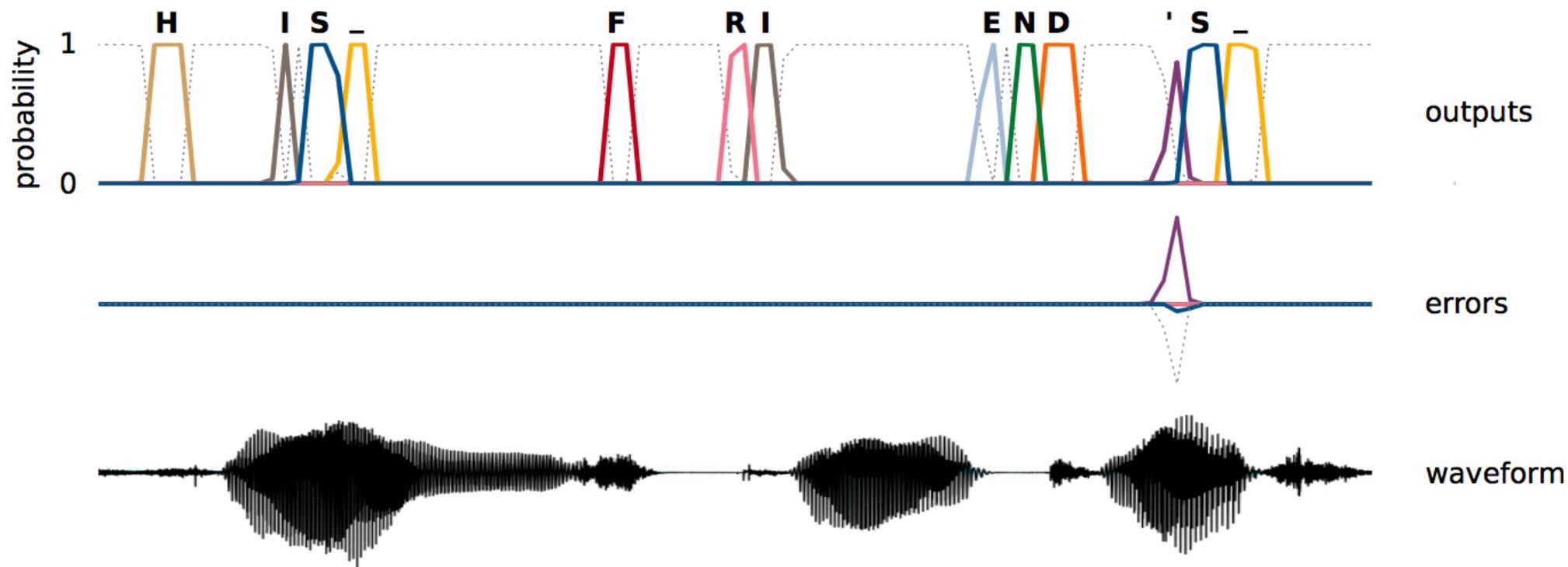
Encourage outputs changing that alter the loss functions

Expected Transcription Loss

Target: HIS_FRIENDS_

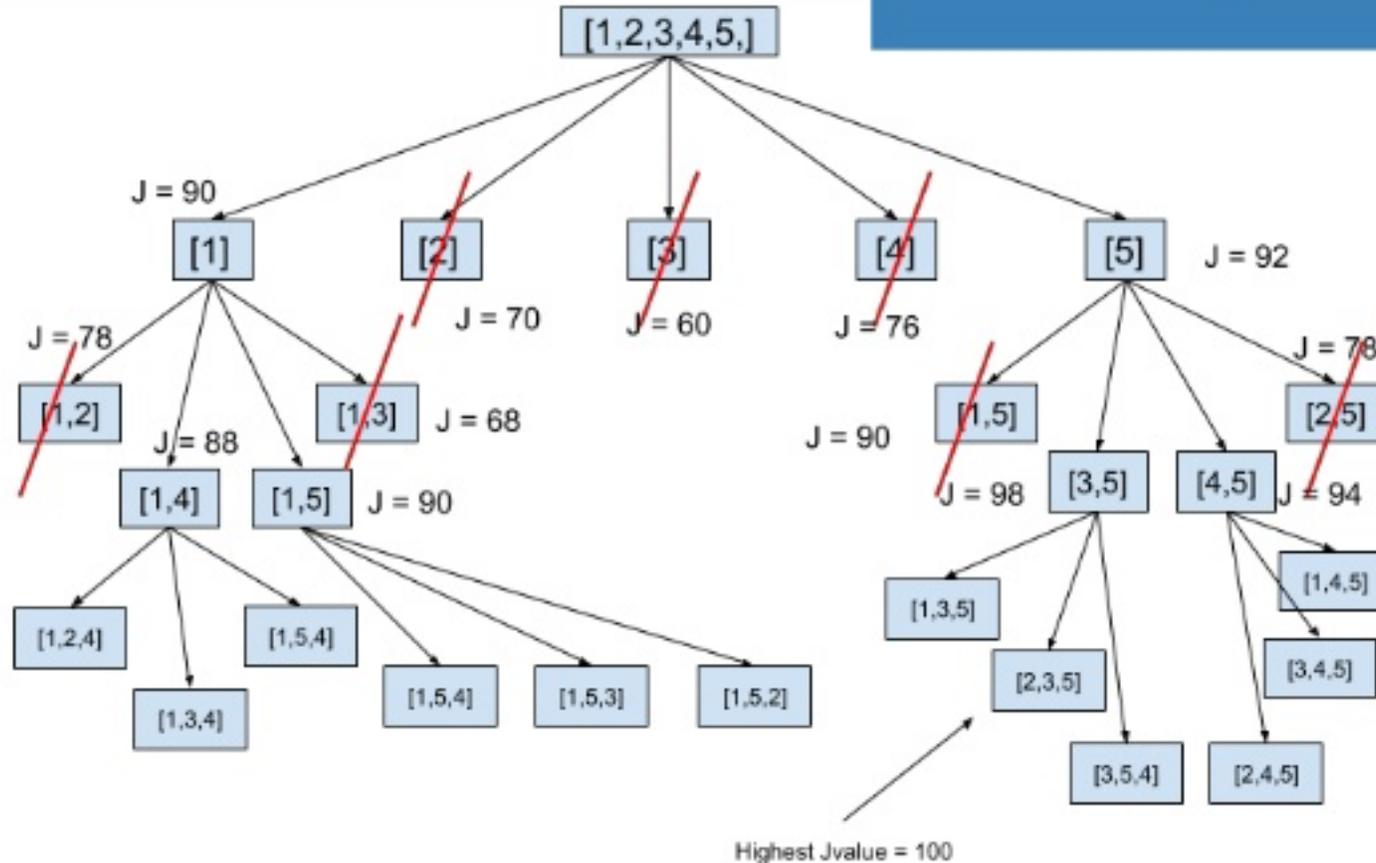
Output: HIS_FRIEND'S_

Changing of the apostrophe is encouraged



Decoding

Decode the output of CTC to a sequence of result: beam search algorithms



Doing pruning during the search and cut the nodes with less probability and only keep top n nodes

Experiments

Table 1. Wall Street Journal Results. All scores are word error rate/character error rate (where known) on the evaluation set. 'LM' is the Language model used for decoding. '14 Hr' and '81 Hr' refer to the amount of data used for training.

SYSTEM	LM	14 HR	81 HR
RNN-CTC	NONE	74.2/30.9	30.1/9.2
RNN-CTC	DICTIONARY	69.2/30.0	24.0/8.0
RNN-CTC	MONOGRAM	25.8	15.8
RNN-CTC	BIGRAM	15.5	10.4
RNN-CTC	TRIGRAM	13.5	8.7
RNN-WER	NONE	74.5/31.3	27.3/8.4
RNN-WER	DICTIONARY	69.7/31.0	21.9/7.3
RNN-WER	MONOGRAM	26.0	15.2
RNN-WER	BIGRAM	15.3	9.8
RNN-WER	TRIGRAM	13.5	8.2
BASELINE	NONE	—	—
BASELINE	DICTIONARY	56.1	51.1
BASELINE	MONOGRAM	23.4	19.9
BASELINE	BIGRAM	11.6	9.4
BASELINE	TRIGRAM	9.4	7.8
COMBINATION	TRIGRAM	—	6.7

Combine RNN
with DNN

Performance is close to the baseline, while baseline has a lot of priority knowledge.
Improve 1% when combine the model

Discussion

Examples: (NO dictionary or language model used for decoding)

target: **TO ILLUSTRATE** THE POINT A PROMINENT MIDDLE EAST ANALYST IN WASHINGTON RECOUNTS
A CALL FROM ONE CAMPAIGN

output: **TWO ALSTRAIT** THE POINT A PROMINENT MIDILLE EAST ANALYST IM WASHINGTON
RECOUNCACALL FROM ONE CAMPAIGN

target: T. W. A. ALSO PLANS TO HANG ITS **BOUTIQUE SHINGLE** IN AIRPORTS AT LAMBERT SAINT

output: T. W. A. ALSO PLANS TOHING ITS **BOOTIK SINGLE** IN AIRPORTS AT LAMBERT SAINT

target: ALL THE EQUITY RAISING IN MILAN GAVE THAT STOCK MARKET INDIGESTION LAST YEAR

output: ALL THE EQUITY RAISING IN MULONG GAVE THAT STACRK MARKET IN TO JUSTIAN LAST YEAR

target: THERE'S UNREST BUT WE'RE NOT GOING TO LOSE THEM TO DUKAKIS

output: THERE'S UNREST BUT WERE NOT GOING TO LOSE THEM TO DEKAKIS

Thanks!