

Soltau, H.
Liao, H.
Sak, H.

Google

Neural Speech Recognizer

Acoustic-to-Word LSTM
Model for Large
Vocabulary Speech
Recognition

arXiv:1610.09975

- Previous Study (Sak et al. 2015)
 - LSTM RNN + CTC
 - Learn an alignment between acoustic input and label sequences
 - Can recognize **whole words**
 - Vocabulary of 90k words
 - Fast and accurate, without decoding, but still far from the sub-word phone-based models
- This Paper
 - Applied the techniques on a **larger** dataset
 - Data sparsity can be alleviated

Labels	Initialization			+sMBR	
	Method	Uni	Bi	Uni	Bi
CD state	CE	15.6	14.0	14.0	12.9
CI phone	CTC	15.5	14.1	14.2	12.7
CD phone	CTC	14.3	13.6	12.9	12.2

Table 2: *WERs (%) for sequence-trained LSTM RNN models.*

Vocabulary	OOV	WER (%)	In vocab. WER (%)
25k Word	4.8	19.5	14.5
7k Word	13	26.8	11.8

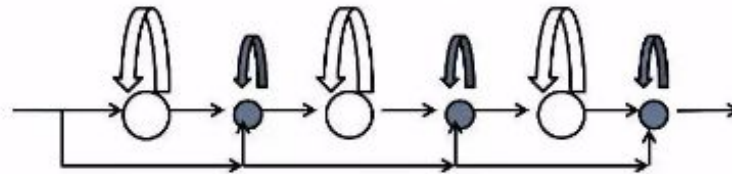
Table 3: *LSTM RNN CTC word acoustic models. The WERs and out of vocabulary (OOV) rates for word models are on held-out data with no decoding or language model. WERs in the last column are computed ignoring utterances containing OOVs.*

- Connectionist Temporal Classification
 - A sequence alignment/labeling technique
 - An additional unit for the **blank** label used to represent outputting no label at a given time

a b c = blank a a b blank c c c blank

= blank a blank b b blank c blank

= blank a a a a blank b b b c c blank



- Relieves the network from having to label each frame by introducing the blank label, enables the use of **longer duration modeling units**

- Loss Function of CTC

$$\mathcal{L}_{CTC} = - \sum_{(\mathbf{x}, \mathbf{l})} \ln p(\mathbf{z}^{\mathbf{l}} | \mathbf{x}) = - \sum_{(\mathbf{x}, \mathbf{l})} \mathcal{L}(\mathbf{x}, \mathbf{z}^{\mathbf{l}})$$

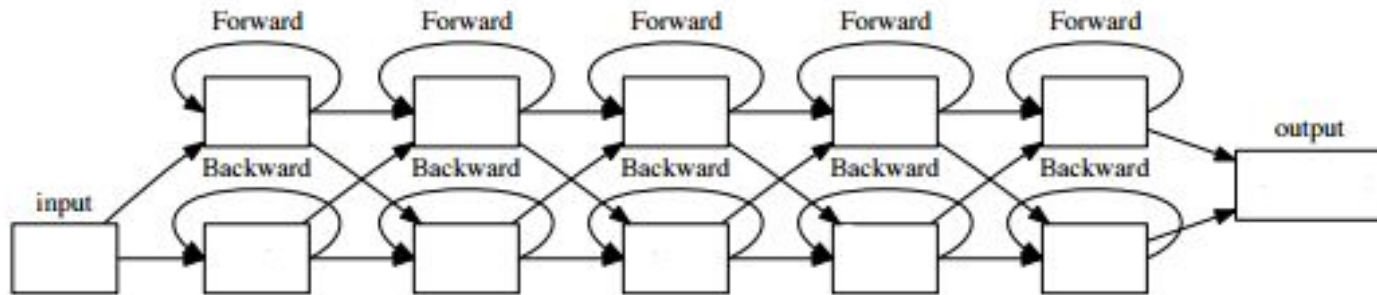
- \mathbf{x} : input sequence of acoustic frames
- \mathbf{l} : input label sequence
- $\mathbf{z}^{\mathbf{l}}$: lattice encoding all possible alignments of \mathbf{x} with \mathbf{l}
- $p(\mathbf{z}^{\mathbf{l}} | \mathbf{x})$: probability for correct labelings

- Gradient of CTC Loss

$$\frac{\partial \mathcal{L}(x, z^l)}{\partial a_l^t} = y_l^t - \frac{1}{p(z^l|x)} \sum_{u \in \{u: z_u^l = l\}} \alpha_{x, z^l}(t, u) \beta_{x, z^l}(t, u)$$

- y_l^t : softmax activation for a label l at time step t
- u : lattice states aligned with label l at time t
- $\alpha_{x, z^l}(t, u)$: forward variable, the summed probability of all paths in the lattice z^l starting in the initial state at time 0 and ending in state u at time t
- $\beta_{x, z^l}(t, u)$: backward variable

- Bidirectional LSTM RNN
 - 5x600
 - 7x1000
- Layer Connections in Bidirectional LSTM



- Input: mel-spaced log filterbank features
- Output: word posterior probabilities
- Distributed Training
 - Asynchronous SGD
 - Optimized Native TensorFlow CPU kernel

- Youtube
 - Test Set
 - Videos from Google Preferred channels
 - 296 videos from 13 categories (avg. 5 min)
 - ~25 hours, 250k words
 - Training Set (semi-supervised)
 - Leverage user-uploaded captions for labels
 - select only audio segments in a video where the user uploaded caption matches the transcript produced by an ASR system
 - ~125k hours, 1.2B words, vocabulary of 1.7M
 - Spoken Vocabulary
 - >100 times, 82k words, OOV 0.63%
 - Written Vocabulary
 - > 80 times, 98k words, OOV 0.7%

- Conventional State/Phone based Models
 - CD triphone states
 - CD single-state phone units

Table 1: Bidirectional-LSTM acoustic models trained on data sets of varying sizes.

Model	Training Criterion	Size	Data (hrs)	WER(%)
CD states	CE	5x600	650	29.0
	CE	5x600	5000	21.2
CD phones	CE	5x600	5000	20.3
	CE	5x600	50000	17.7
	CE	5x600	125000	16.7
	CTC	5x600	125000	16.5
	CTC, multi_lstm_op ¹	5x600	125000	15.5
	CTC, multi_lstm_op ¹	7x1000	125000	14.2

- There is **little difference** between CE and CTC training criteria.
 - Asynchronous SGD gives better results with faster parameter updates

State vs. Phone

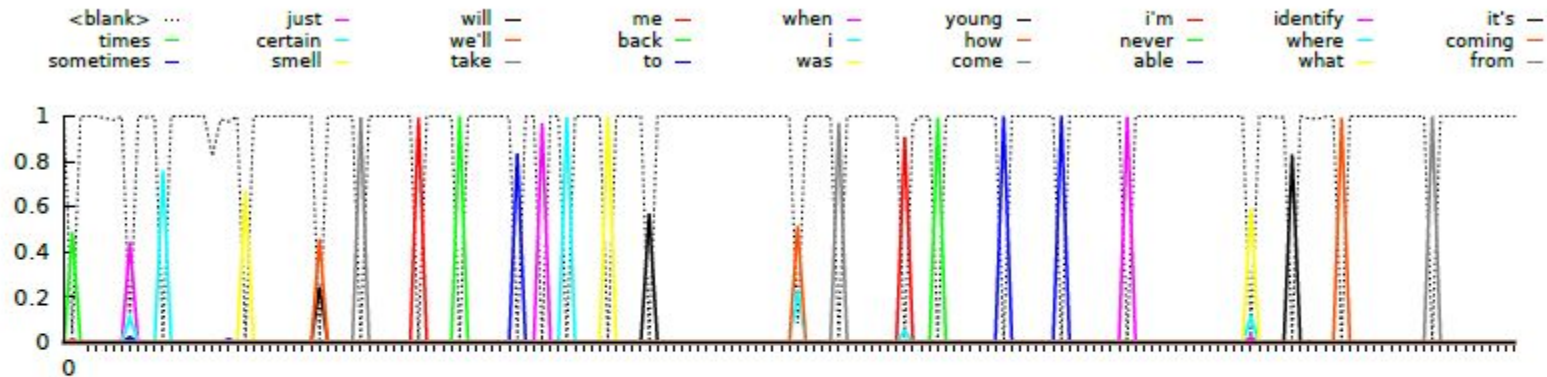


Figure 1: The word posterior probabilities as predicted by the NSR model at each time-frame (30 msec) for a segment of music video 'Stressed Out' by Twenty One Pilots. We only plot the word with highest posterior and the missing words from the correct transcription: '*Sometimes a certain smell will take me back to when I was young, how come I'm never able to identify where it's coming from*'.

- Word Models Compared with Phone Models
 - Word model can be used **without decoding or language model** → end-to-end recognizer

Table 2: CTC CD phone models compared with CTC word models.

Model	Layers	Outputs	Params	Vocab	OOV(%)	WER(%)	
						w/ LM	w/o LM
CTC CD phone	5x600	6400	14m	500000	0.24	15.5	—
	7x1000	6400	43m	500000	0.24	14.2	—
	7x1000	35326	75m	500000	0.24	14.5	—
	7x1000	6400	43m	82473	0.63	14.7	—
CTC spoken words	5x600	82473	57m	82473	0.63	14.5	15.8
	7x1000	82473	116m	82473	0.63	13.5	14.8
CTC written words	7x1000	97827	137m	97827	0.70	13.4	13.9

- Capable of accurate speech recognition with no LM or decoding involved

- Error Rate Correction for Spoken Word Model
 - References are in written domain while model output is in spoken domain
 - errors like “three” vs. “3”
 - Force align the utterances with a graph
 - $C * L * \text{project} (V * T)$
 - C: context transducer
 - L: lexicon transducer
 - V: spoken-to-written transducer
 - project: map the input symbols to the output symbols
 - $\text{project} (V * G)$
 - convert written language model G to a spoken form
 - use the spoken LM to build the decoding graph

Table 3: Comparison of CD phone with spoken word models in spoken domain.

Model	Layers	Outputs	Params	Vocab	OOV(%)	Spoken WER(%)	
						w/ LM	w/o LM
CTC CD phone	7x1000	6400	43m	500000	0.24	12.3	—
CTC spoken words	7x1000	82473	116m	82473	0.63	11.6	12.0

- Word models without use of any language model or decoding performs at 12.0% WER, slightly better than the CD phone model that uses an LVCSR decoder and incorporates a 30m 5-gram language model.
- Adding LM for the CTC spoken word model improves the error rate from 12.0% to 11.6%, not too much.

- The final system performs better than a well-trained, conventional CD phone-based system on a difficult YouTube video transcription task
 - word model of bidirectional LSTM plus CTC loss having 7x1000 layers with 116 parameters and 82k vocabulary size
 - 13.4% WER for written domain with LM
 - 11.6% spoken WER for spoken domain with LM