# RNN-T Based ASR Systems

Mahaveer Jain, Facebook

#### ASR And Machine Translation As Sequence To **Sequence Problem** Speller (eas)



#### Machine Translation

Image from Bhiksha Raj's lecture on attention models

LAS attends to all audio embeddings and uses text history produced so far to generate probability distribution over output units

$$f(y_u | y_{u-1}, y_{u-2}, ..., y_0, h_0, h_2, h_3, ..., h_v$$



#### **ASR As Sequence To Sequence Problem**

- Alignment not known during training
- At each audio frame index, neural network generates probability of producing each output symbol. e. g. Probability of generating output symbol "k" at 1<sup>st</sup> frame.

 $P_K^1 = prob(s_1 = K \mid f_1)$ 

- #audio frames != # letters. We also don't know which portion of audio was emitted when.
- Long short-term memory (LSTM), B-LSTM (bi-directional LSTM) and, Transformer are some of the commonly used audio encoders.



### ASR As Sequence To Sequence Problem

- Input Sequence: Audio frame features
- Output Sequence : Words from True Transcript Text written as letters
- Example Training Instance: Input Sequence : 7 audio frames and True Transcript Text: "BE"
- At each audio frame index, network\* generates probability of producing each output symbol. E. g. Probability of generating output symbol "K" at 1st frame.
- #au  $P_K^1 = prob(s_1 = K \mid f_1)$  lent i.e. which portion of audio aligns to what letter in true transcript.
- Collapse continuous occurrences of same letter into one:
  - First alignment choice (BBEEEEE -> BE)
  - Second alignment choice (BBBEEEE -> BE)
  - and so on.....
- How to deal with repeated letters (as in word BEE)? (next slide)



Audio



Neural Network (i.e. LSTM)

Softmax Layer (Toy example: If we assume there are just three letters in the language (B, E, C))

# **Connectionist Temporal Classification (CTC)**

- Input Sequence: Audio frame features
- Output Sequence: Letters
- At each audio frame index, network generates probability of producing each output symbol and blank symbol (Ø)
- How to deal with #audio frame != # letters -> Collapse letters
- How to deal with repeated letters (as in word BEE)? -> blank symbol (Ø). Ø also indicates "emit nothing" (silence). ØBEØEEØ
- Conditional Independence: Probability output at each frame does not depend on history of transcript produced so far, this non auto-regression on text history makes it less effective on learning LM information and, you may still need language model -> RNN-T.



#### **Conditional Independence Assumption**



#### CTC:

 Only uses audio embedding at time "t" (x<sub>t</sub>) to generate probability distribution over output units. CTC assumes that each output unit is conditionally independent of others.



#### LAS:

Attends to all audio embeddings and uses text history produced so far to generate probability distribution over output units.

$$f(y_u | y_{u-1}, y_{u-2}, ..., y_0, h_0, h_2, h_3, ..., h_v)$$

Streaming

# **Conditional Independence Assumption**



#### RNN-T:

Uses both audio embedding at time "t"  $(x_t)$  and text history produced so far  $(y_{\mu-1})$  to generate probability distribution over output units.

Attends to all audio embeddings and uses text history produced so far to generate probability distribution over output units.

 $f(y_u|y_{u-1}, y_{u-2}, ..., y_0, h_0, h_2, h_3, ..., h_v)$ 

Streaming

Only uses audio embedding at time "t"  $(x_t)$  to

CTC assumes that each output unit is

conditionally independent of others.

generate probability distribution over output units.

CTC:

•

### Why RNN-T ASR From Production Point Of View

- Single deployable non modularized neural model, all components of ASR in one model. ٠
- Allows compact on-device streaming ASR. Does not need a decoder graph which can be ٠ large. Unlimited words in vocab. Standard on-device ASR choice across industry.



An All-Neural On-Device Speech Recognizer Tuesday, March 12, 2019 Posted by Johan Schalkwyk, Google Fellow, Speech Team

Image from https://ai.googleblog.com/2019/03/an-all-neural-on-devicespeech.html

Achieves comparable accuracy and compute with much smaller size model compared to ٠ modularized (hybrid) systems for production when training data is the same.

Table 3. Comparison of Hybrid model with RNN-T				
Test Set	System	WER	Throughput	rtf@40
vid-clean	hybrid	14.0	55	.70
vid-clean	RNN-T	14.0	63	.60
vid-noisy	hybrid	20.7	55	.71
vid-noisy	RNN-T	21.0	65	.60

Table 2 Communication of Hadrid and John MALT

Example Study

Image from Jain et al., "RNN-T For Latency Controlled ASR With Imoroved Beam Search'

# Components Of RNN-T ASR System

- Audio Encoder: To encode sequence of audio features into audio embeddings. Long short-term memory (LSTM), B-LSTM (bi-directional LSTM), Transformer are commonly used audio encoders. Acts as acoustic model.
- Text Predictor: To encode transcript produced so far into text embedding. Typically a LSTM. Acts as language model.
- Joiner combines output of audio encoder and text predictor.
- A Linear Layer followed by Softmax produces probability distribution over output units.
   Pr(k|t,u)(P<sup>k</sup><sub>(t,u</sub>) is probability of emitting "k" from (t, u).
- How is Ø different from other symbols: If emitted output symbols is blank (Ø) then move to next time frame else stay in the same time frame. Ø also indicates "emit nothing".



# Training



# **Training Objective**

- Training utterance:
  "Input Audio" -> "BEE" (transcript)
- #output units (u): 3 BEE -> [[B], [E], [E]]
- #audio frames (t) : 4
- We don't know alignment i.e. which portion of audio aligns to what output unit (A path taken in lattice)
- Training Objective

$$\begin{split} P(BEE|X) &= \sum_{\substack{alignment}} P(alignment, BEE|X) \\ P(BEE|alignment, X) &= 1 & \sum_{\substack{alignment\\alignment}} P(BEE|alignment, X) P(alignment|X) \\ \sum_{\substack{alignment}} P(alignment|X) \end{split}$$



Audio















# **RNN-T Lattice And Training**

- Training utterance:
  "Input Audio" -> "BEE" (transcript)
- #output units (u): 3 BEE -> [[B], [E], [E]]
- #audio frames (t) : 4
- We don't know alignment i.e. which portion of audio aligns to what output unit
- Probability of alignment is multiplication of probabilities assigned along the path of alignment

$$\begin{aligned} & P(BEE|X) = \sum_{\substack{alignment \\ alignment}} P(alignment, BEE|X) \\ P(BEE|alignment, X) = 1 & \sum_{\substack{alignment \\ alignment}} P(BEE|alignment, X) P(alignment|X) \end{aligned}$$

- Lattice contains all valid alignment paths(traversals). During training, we change (optimize) neural network parameters to maximize sum of probabilities of all alignment paths
- Computation is done efficiently through dynamic programming (slide 54 to 58)







#### Graves et al., "Sequence transduction with recurrent neural networks"

# Inference: Find a Transcript given an audio and Trained RNN-T model

Goal: Find candidate transcript of a given audio during test time. Challenge: There are infinitely many alignments that can be assigned to a test audio.



# **Operation: Extend Hypothesis**

- Hypothesis is defined as a candidate output sequence(including Ø) during search
- Example Hypothesis: "EB" at time frame "t"
- Output symbols: "E", "B", "C", "Ø"

Time = t + 1

- Extend Hypothesis: Append hypothesis with each of the output symbols (k) and  $\phi$
- The Ø extension goes to next time frame (t + 1) and non-blank extensions remain in the same time frame (t)
- Every extended hypothesis has lower probability compared to the hypothesis it was extended from

EB

EBE

EB

**p** \* **P**<sup>Ø</sup><sub>(t, "BE")</sub>

Time = t

**p** \* **P**<sup>B</sup><sub>(t, "BE")</sub>

EB

EB

**p** \* **P**<sup>E</sup><sub>(t, "BE")</sub>

**p** \*

Р<sup>С</sup>(t, "BE")





• Goal: Get a candidate final transcript by making local optimal choice at each distinct "t" and "u"



Е 0.2 Lets make a local optimal • В 0.3 choice: If highest С 0.1 h<sup>t,u</sup> probability is given to Linear Layer blank (Ø), move to the  $z^{t,u}$ next audio frame else stay in the same audio frame Do it until all audio frames ٠ ft  $g_u$ [< *BOS* are processed Audio Encoder (E Ø [ < BOS > ] $x_1$ Mm Mmmm MM  $X_1$  $X_2$ **RNN-T model** features helper helper ? Transcrip Audio Audio

 $\Pr(k|t,u)$ 

probability

t = 1

В

 $X_3$ 

С

 $X_4$ 

0.4

k

Ø

Е 0.04 Lets make a local optimal ٠ В 0.3 choice: If highest С 0.01 h<sup>t,u</sup> probability is given to Linear Layer blank (Ø), move to the  $z^{t,u}$ next audio frame else stay in the same audio frame Do it until all audio frames ٠ ft  $g_u$ are processed [ < BOS ]Audio Encoder (E Ø [ < BOS > ] $x_2$ 1mhhmmm/mlm MM  $X_3$  $X_1$ **X**<sub>2</sub> features **RNN-T model** helper helper ? Transcrip Audio Audio

 $\Pr(k|t,u)$ 

probability

t = 2

В

С

 $X_4$ 

0.04

k

Ø



k





probability

С

k



probability

С

k



## Subset Of Hypotheses At t = 1 During Decoding





- We don't want to stay in time frame "t" forever. But we do not know when to move to "t+1" as there are infinitely many candidates that can be explored in time frame "t".
- Goal: We would like to obtain *n* candidate hypotheses at time frame "t+1" before exiting to decode at "t" which would be better than all possible future extensions at "t".
- *n* is hyper parameter

## Beam Search: Expand the best hypothesis among candidate hypotheses



Is beam size (n) number of extensions enough? No.



MWMMWWW



 $X_1$ 







 $X_4$ 

P<sup>®</sup>IW PE(LU)

P<sup>B</sup><sub>0</sub>... P<sup>C</sup><sub>tab</sub>

Output symbols

 $\Pr(k|t, u) \rightarrow \Pr_{(t,u)}^{k}$ 

ht,u

 $y_{u-1}$ 

## Beam Search: When to exit search at time frame t



#### Thanks

 Many thanks to Rohit Prabhavalkar, Mark Tygert, Michael Picheny, Nayan Singhal, Kritika Singh, Xiaohui Zhang, Duc Le, Yuan Shangguan, Paco Guzmán, Yatharth Saraf and Mike Seltzer for brainstorming with ideas to improve content of this presentation.