

Deep Learning Audio

Lecture 4

Pavel Severilov

Moscow Institute of Physics and Technology

2022

Outline

1. RNN-Transducer (RNN-T)
2. Language models for ASR
3. Byte-pair encoding (BPE)

ASR: SOTA models

- ▶ RNN-T (2018, Google)
- ▶ MoChA (2018, Google) (Tricks to make LAS online)
- ▶ wav2vec (2019, Facebook AI Research) (use WAV not spectrograms)
- ▶ Jasper (2019, Nvidia) (Encoder – CNN ; Loss: CTC)
- ▶ QuartzNet (2019, Nvidia) (Encoder – TDS CNN ; Loss – CTC)
- ▶ ContextNet (2020, Google) (Encoders – CNN and LSTM ; Loss – RNN-T)
- ▶ Whisper (2022, OpenAi) (Encoder - Transformer; Loss - multitask)

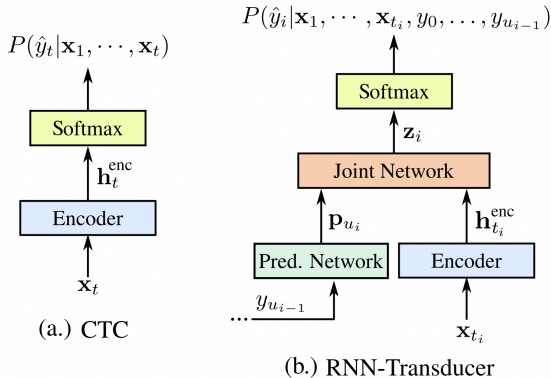
Outline

1. RNN-Transducer (RNN-T)
2. Language models for ASR
3. Byte-pair encoding (BPE)

CTC and Attention models: recap

	CTC	Listen, Attend and Spell: LAS	?
Summary	Maximize probability of all possible CTC-paths leading to target.	Encoder-decoder architecture with attention.	???
Online	+	-	+
Context dependent	-	+	+
Multiple outputs for each input	-	-	+

RNN-T: idea



- ▶ Predictor is autoregressive: takes as input the previous outputs.
- ▶ Joiner – feedforward network, combines the encoder vector \mathbf{h}_t and predictor vector \mathbf{p}_u

He et al. Streaming End-to-end Speech Recognition for Mobile Devices / 2019, Google, Inc.

RNN-T: model

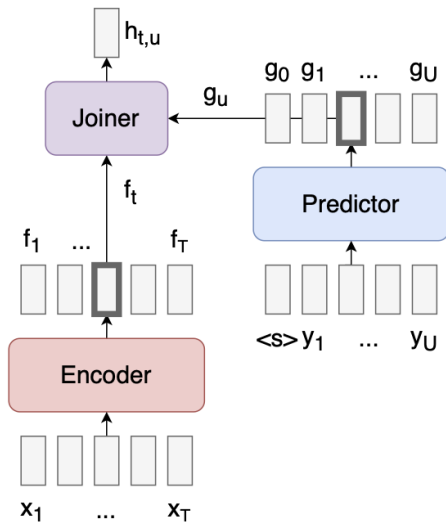


Figure: RNN-T architecture

RNN-T: model

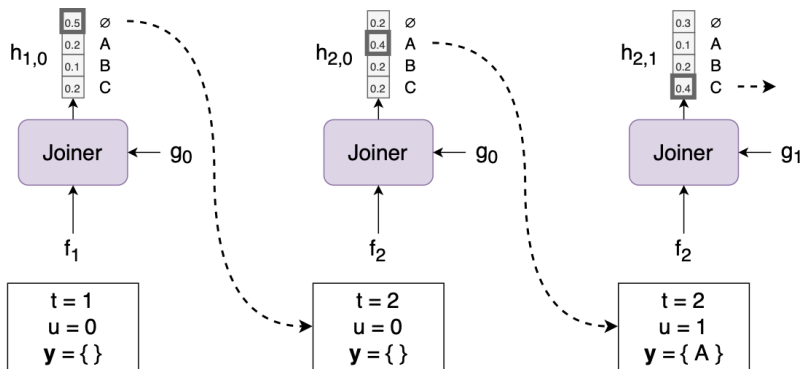


Figure: Steps example of RNN-T inference: t – audio-encoder timestamp, u – Predictor (char network) step

RNN-T: training

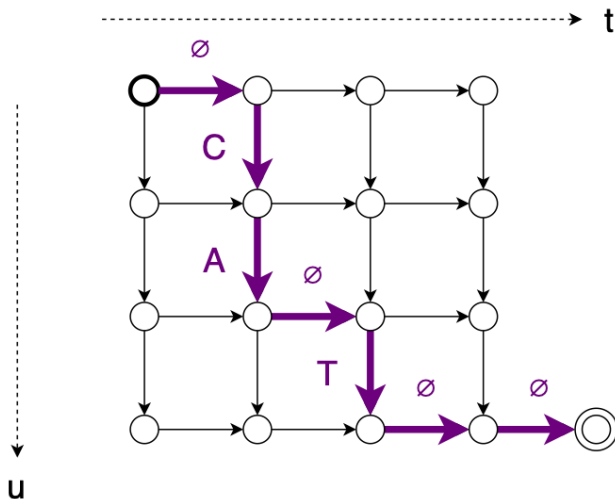


Figure: Alignment $\{\emptyset, C, A, \emptyset, T, \emptyset, \emptyset\}$ for input sequence of length $T = 4$ and an output sequence "CAT" of length $U = 3$

RNN-T: training

We need to get $p(\mathbf{y}|\mathbf{x})$ as the sum of the probabilities of all possible alignments between \mathbf{x} and \mathbf{y}

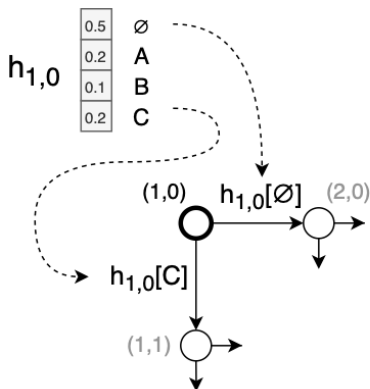


Figure:

$\mathbf{z} = \emptyset, C, A, \emptyset, T, \emptyset, \emptyset$

$$p(\mathbf{z} | \mathbf{x}) = h_{1,0}[\emptyset] \cdot h_{2,0}[C] \cdot h_{2,1}[A] \cdot h_{2,2}[\emptyset] \cdot h_{3,2}[T] \cdot h_{3,3}[\emptyset] \cdot h_{4,3}[\emptyset]$$

Sequence-to-sequence learning with Transducers blog post

RNN-T: training

To compute the sum efficiently, compute $\alpha_{t,u}$, for $1 \leq t \leq T$ and $0 \leq u \leq U$

$$\begin{aligned}\alpha_{t,u} &= \alpha_{t-1,u} \cdot h_{t-1,u}[\emptyset] \\ &\quad + \alpha_{t,u-1} \cdot h_{t,u-1}[y_{u-1}]\end{aligned}$$

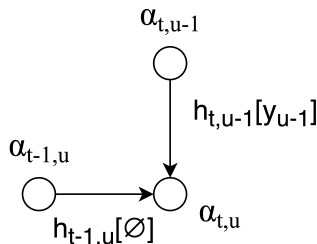


Figure: Computing $\alpha_{t,u}$ to get $p(\mathbf{y} \mid \mathbf{x}) = \alpha_{T,U} \cdot h_{T,U}[\emptyset]$

Outline

1. RNN-Transducer (RNN-T)
2. Language models for ASR
3. Byte-pair encoding (BPE)

Language models (LM): why need in ASR?

- ▶ Language models recap: a model that estimates the probability of a text.
 - ▶ N-gramms
 - ▶ Neural networks (BERT, GPT-3, ...)
 - ▶ Example:
 $P(\text{let's go two a movie}) = 0.01$
 $P(\text{let's go to a movie}) = 0.6$
- ▶ ASR problem:
 - ▶ Spelling of a word heavily depends on its context
 - ▶ Labeled audio data is difficult to obtain
- ▶ How LM helps:
 - ▶ Improves final WER
 - ▶ Improves performance for small audio datasets
 - ▶ Can be used to adapt model to new domain

LM: how to integrate in ASR?

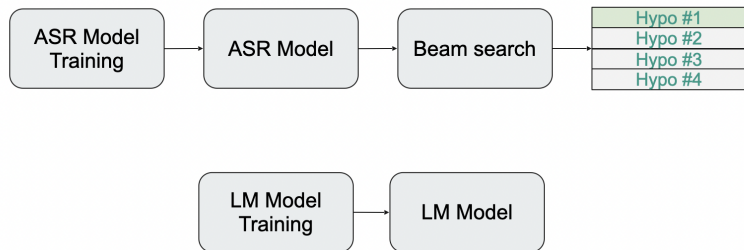


Figure: ASR pipeline VS Language models pipeline

LM: final hypothesis rescoring

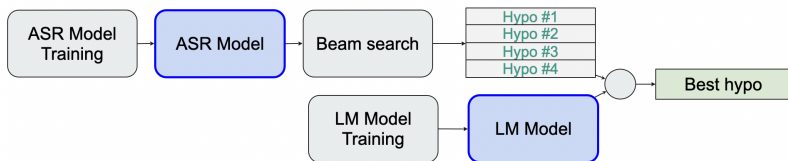


Figure: Final hypothesis rescoring: rescore beam-search output with LM probs

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \log p(\mathbf{y} \mid \mathbf{x}) + \lambda \log p_{LM}(\mathbf{y}) + \beta \cdot \text{len}(\mathbf{y})$$

$\text{len}(\mathbf{y})$ – function of word length, anti-penalty for long words

LM: shallow fusion

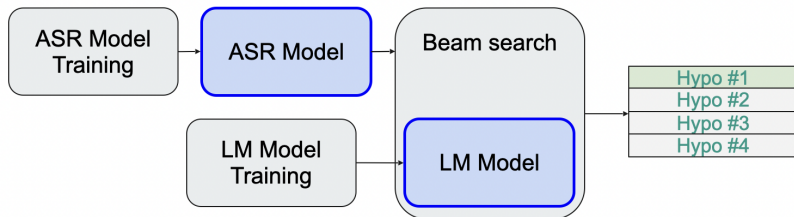


Figure: Shallow fusion: use LM rescoring after each beam search step

Practice:

- ▶ requires much more LM runs
- ▶ use light LM for shallow fusion
- ▶ use heavy LM for second-pass rescoring

LM: Deep fusion

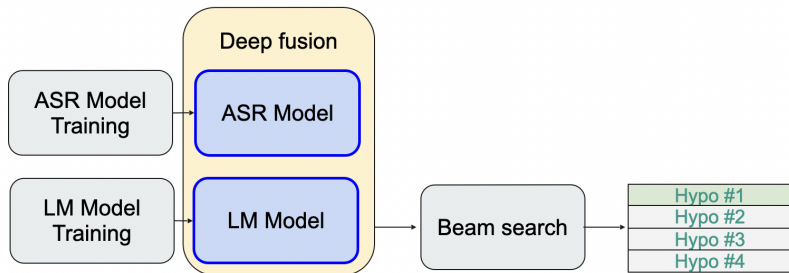


Figure: Deep fusion: integrates the external LM into the encoder-decoder model (ASR) by fusing together the hidden states of the external LM and the decoder

LM: Cold fusion

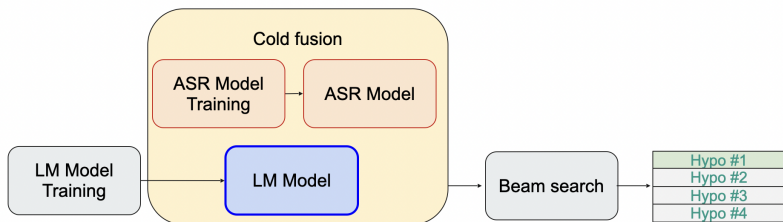


Figure: Cold fusion: train like in Deep fusion, but jointly with ASR model

LM in ASR: comparison of approaches

Model	SWB	CH	Full
LAS	17.1	27.9	22.6
Shallow Fusion	15.6	26.6	21.1
Deep Fusion	16.3	27.2	21.7
Cold Fusion	16.3	27.3	21.8

Table: Word error rates (%) on Eval2000 for the LAS baseline model and fusion approaches. SWB=Switchboard, CH=CallHome, Full=Eval2000.

Outline

1. RNN-Transducer (RNN-T)
2. Language models for ASR
3. Byte-pair encoding (BPE)

BPE: motivation & idea

- ▶ Motivation: a lot of characters have different pronunciation in different contexts
- ▶ Idea: let's use n-gramms as tokens in addition
- ▶ Advantages:
 - ▶ Less decoder steps → faster training and inference
 - ▶ Better generalization → better WER

BPE: algorithm

1. Each character – token
2. Most popular n-gram: add new token
3. Replace n-gram with a new token
4. Restrict maximum length of tokens
5. New vocabulary = all characters + new tokens

Iteration	Sequence	Vocabulary
0	a b a b c a b c	{a, b, c}
1	ab ab c ab c	{a, b, c, ab}
2	ab abc abc	{a, b, c, ab, abc}
3	ababc abc	{a, b, c, ab, abc, ababc}
4	ababcabc	{a, b, c, ab, abc, ababc, ababcabc}

Table: BPE: example for sequence {ababcabc}