# End-to-end systems 1: CTC
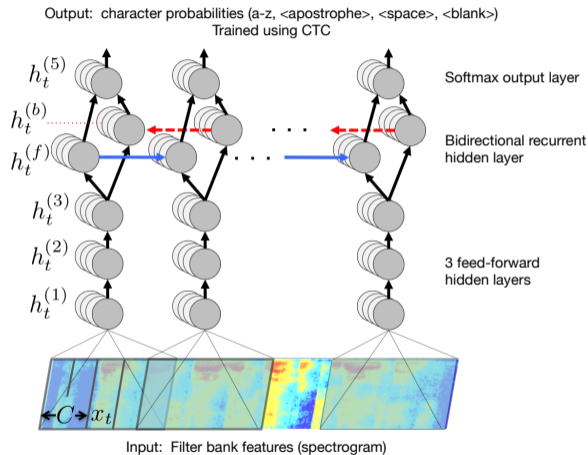## (Connectionist Temporal Classification)

Steve Renals

Automatic Speech Recognition – ASR Lecture 15
11 March 2019

## End-to-end systems

- End-to-end systems are systems which learn to directly map from an input sequence $X$ to an output sequence $Y$, estimating $P(Y|X)$
  - $Y$ can be a sequence of words or subwords
- ML trained HMMs are kind of end-to-end system – the HMM estimates $P(X|Y)$, and when combined with a language model gives an estimate of $P(Y|X)$
- Sequence discriminative training of HMMs (using GMMs or DNNs) can be regarded as end-to-end
  - But training is quite complicated – need to estimate the denominator (total likelihood) using lattices, first train conventionally (ML for GMMs, CE for NNs) then finetune using sequence discriminative training
  - Lattice-free MMI is one way to address these issues
- Other approaches based on recurrent networks which directly map input to output sequences
  - CTC – Connectionist Temporal Classification
  - Encoder-decoder approaches

# End-to-end systems

- End-to-end systems are systems which learn to directly map from an input sequence $X$ to an output sequence $Y$, estimating $P(Y|X)$
  - $Y$ can be a sequence of words or subwords
- ML trained HMMs are kind of end-to-end system – the HMM estimates $P(X|Y)$, and when combined with a language model gives an estimate of $P(Y|X)$
- Sequence discriminative training of HMMs (using GMMs or DNNs) can be regarded as end-to-end
  - But training is quite complicated – need to estimate the denominator (total likelihood) using lattices, first train conventionally (ML for GMMs, CE for NNs) then finetune using sequence discriminative training
  - Lattice-free MMI is one way to address these issues
- Other approaches based on recurrent networks which directly map input to output sequences
  - **CTC – Connectionist Temporal Classification**
  - Encoder-decoder approaches (*next lecture*)

# Deep Speech



Output: character probabilities (a-z, <apostrophe>, <space>, <blank>)
Trained using CTC

$h_t^{(5)}$ — Softmax output layer

$h_t^{(b)}$
$h_t^{(f)}$ — Bidirectional recurrent hidden layer

$h_t^{(3)}$
$h_t^{(2)}$ — 3 feed-forward hidden layers
$h_t^{(1)}$

$C$ $x_t$

Input: Filter bank features (spectrogram)

Hannun et al (2014), "Deep Speech: Scaling up end-to-end speech recognition",

https://arxiv.org/abs/1412.5567.

# Deep Speech: Results

| Model | SWB | CH | Full |
|---|---|---|---|
| Vesely et al. (GMM-HMM BMMI) [44] | 18.6 | 33.0 | 25.8 |
| Vesely et al. (DNN-HMM sMBR) [44] | 12.6 | 24.1 | 18.4 |
| Maas et al. (DNN-HMM SWB) [28] | 14.6 | 26.3 | 20.5 |
| Maas et al. (DNN-HMM FSH) [28] | 16.0 | 23.7 | 19.9 |
| Seide et al. (CD-DNN) [39] | 16.1 | n/a | n/a |
| Kingsbury et al. (DNN-HMM sMBR HF) [22] | 13.3 | n/a | n/a |
| Sainath et al. (CNN-HMM) [36] | 11.5 | n/a | n/a |
| Soltau et al. (MLP/CNN+I-Vector) [40] | **10.4** | n/a | n/a |
| **Deep Speech SWB** | 20.0 | 31.8 | 25.9 |
| **Deep Speech SWB + FSH** | 12.6 | **19.3** | **16.0** |

Table 3: Published error rates (%WER) on Switchboard dataset splits. The columns labeled "SWB" and "CH" are respectively the easy and hard subsets of Hub5'00.

# Deep Speech Training

- Maps from acoustic frames $X$ to subword sequences $S$, where $S$ is a sequence of characters (in some other CTC approaches, $S$ can be a sequence of phones)
- CTC loss function
- Makes good use of large training data
  - Synthetic additional training data by jittering the signal and adding noise
- Many computational optimisations
- n-gram language model to impose word-level constraints
- Competitive results on standard tasks

# Deep Speech Training

- Maps from acoustic frames $X$ to subword sequences $S$, where $S$ is a sequence of characters (in some other CTC approaches, $S$ can be a sequence of phones)
- **CTC loss function**
- Makes good use of large training data
  - Synthetic additional training data by jittering the signal and adding noise
- Many computational optimisations
- n-gram language model to impose word-level constraints
- Competitive results on standard tasks

# Connectionist Temporal Classification (CTC)

- Train a recurrent network to map from input sequence $X$ to output sequence $S$
  - sequences can be different lengths – for speech, input sequence $X$ (acoustic frames) is much longer than output sequence $S$ (characters or phonemes)
  - CTC does not require frame-level alignment (matching each input frame to an output token)
- CTC sums over all possible alignments (similar to forward-backward algorithm) – "alignment free"
- Possible to back-propagate gradients through CTC

Gopod overview of CTC: Awni Hannun, "Sequence Modeling with CTC", *Distill*. https://distill.pub/2017/ctc

# CTC: Alignment

- Imagine mapping $(x_1, x_2, x_3, x_4, x_5, x_6)$ to $[a, b, c]$
  - Possible alignments: *aaabbc*, *aabbcc*, *abbbbc*,...
- However
  - Don't always want to map every input frame to an output symbol (e.g. if there is "inter-symbol silence")
  - Want to be able to have two identical symbols adjacent to each other – keep the difference between
- Solve this using an additional *blank* symbol ($\epsilon$)
- CTC output compression
  1. Merge repeating characters
  2. Remove blanks

  Thus to model the same character successively, separate with a blank
- Some possible alignments for $[h, e, l, l, o]$ and $[h, e, l, o]$ given a 10-element input sequence
  - $[h, e, l, l, o]$: *h$\epsilon$$\epsilon$e$\epsilon$ll$\epsilon$lo*; *he$\epsilon$ll$\epsilon$l$\epsilon$oo*
  - $[h, e, l, o]$: *h$\epsilon$$\epsilon$e$\epsilon$lllo*; *hh$\epsilon$e$\epsilon$l$\epsilon$$\epsilon$o$\epsilon$*

# CTC: Alignment example

| h | h | e | $\epsilon$ | $\epsilon$ | l | l | l | $\epsilon$ | l | l | o |
|---|---|---|---|---|---|---|---|---|---|---|---|

First, merge repeat characters.

| h | e | $\epsilon$ | l | $\epsilon$ | l | o |
|---|---|---|---|---|---|---|

Then, remove any $\epsilon$ tokens.

| h | e | | l | | l | o |
|---|---|---|---|---|---|---|

The remaining characters are the output.

| h | e | l | l | o |
|---|---|---|---|---|

Consider an output [c, a, t] with an input of length six

**Valid Alignments**

| $\epsilon$ | c | c | $\epsilon$ | a | t |

| c | c | a | a | t | t |

| c | a | $\epsilon$ | $\epsilon$ | $\epsilon$ | t |

**Invalid Alignments**

| c | $\epsilon$ | c | $\epsilon$ | a | t |

corresponds to $Y = [c, c, a, t]$

| c | c | a | a | t | |

has length 5

| c | $\epsilon$ | $\epsilon$ | $\epsilon$ | t | t |

missing the 'a'

# CTC: Alignment properties

- Monotonic – Alignments are monotonic (left-to-right model); no re-ordering (unlike neural machine translation)
- Many-to-one – Alignments are many-to-one; many inputs can map to the same output (however a single input cannot map to many outputs)
- CTC doesn't find a single alignment: it sums over all possible alignments

# CTC: Loss function (1)

- Let $\boldsymbol{C}$ be an output label sequence, including blanks and repetitions – same length as input sequence $\boldsymbol{X}$
- Posterior probability of output labels $\boldsymbol{C} = (c_1, \ldots c_t, \ldots c_T)$ given the input sequence $\boldsymbol{X} = (x_1, \ldots x_t, \ldots x_T)$:

$$P(\boldsymbol{C}|\boldsymbol{X}) = \prod_{t=1}^{T} y(c_t, t)$$

where $y(c_t, t)$ is the output for label $c_t$ at time $t$
- This is the probability of a single alignment

# CTC: Loss function (2)

- Let $S$ be the target output sequence after compression
- Compute the posterior probability of the target sequence $S = (s_1, \ldots s_m, \ldots s_M)$ ($M \leq T$) given $X$ by summing over the possible CTC alignments:

$$P(S|X) = \sum_{c \in A(S)} P(C|X)$$

where $A$ is the set of possible output label sequences $c$ that can be mapped to $S$ using the CTC compression rules (merge repeated labels, then remove blanks)

- The CTC loss function $\mathcal{L}_{CTC}$ is given by the negative log likelihood of the sum of CTC alignments:

$$\mathcal{L}_{CTC} = -\log P(S|X)$$

- Perform the sum over alignments using dynamic programming – similar structure as used in forward-backward algorithm and Viterbi (see Hannun for details)
- Various NN architectures can be used for CTC – usually use a deep bidirectional LSTM RNN

# CTC: Distribution over alignments

We start with an input sequence, like a spectrogram of audio.

The input is fed into an RNN, for example.

The network gives $p_t(a \mid X)$, a distribution over the outputs {h, e, l, o, $\epsilon$} for each input step.
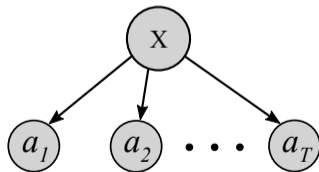
With the per time-step output distribution, we compute the probability of different sequences

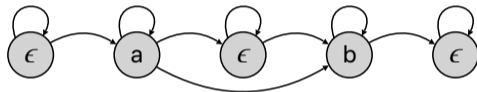By marginalizing over alignments, we get a distribution over outputs.

- Each output is dependent on the entire input sequence (in Deep Speech this is achieved using a bidirectional recurrent layer)
- Given the inputs, each output is independent of the other outputs (conditional independence)
- CTC does not learn a language model over the outputs, although a language model can be applied later
- Graphical model showing dependences in CTC:

Left-to-right HMM          CTC HMM

- CTC can be interpreted as an HMM with additional (skippable) blank states, trained discriminatively

- Direct interpolation of a language model with the CTC acoustic model:

$$\hat{\boldsymbol{W}} = \arg \max_W (\alpha \log P(\boldsymbol{S}|\boldsymbol{X}) + \log P(W))$$

  Only consider word sequences $W$ which correspond to the subword sequence $\boldsymbol{S}$ (using a lexicon)

- $\alpha$ is an empirically determined scale factor to match the acoustic model to the language model

- Lexicon-free CTC: use a "subword language model" $P(\boldsymbol{S})$ (Maas et al, 2015)

- WFST implementation: create an FST $T$ which transforms a framewise label sequence $\boldsymbol{c}$ into the subword sequence $\boldsymbol{S}$, then compose with $L$ and $G$: $T \circ \min(\det(L \circ G))$ (Miao et al, 2015)

# Mozilla Deep Speech

- Mozilla have released an Open Source TensorFlow implementation of the Deep Speech architecture:

- https://hacks.mozilla.org/2017/11/a-journey-to-10-word-error-rate/

- https://github.com/mozilla/DeepSpeech

- Close to state-of-the-art results on librispeech

- Mozilla Common Voice project: https://voice.mozilla.org/en

# Summary and reading

- CTC is an alternative approach to sequence discriminative training, typically applied to RNN systems
- Used in "Deep Speech" architecture for end-to-end speech recognition
- Reading
  - A Hannun et al (2014), "Deep Speech: Scaling up end-to-end speech recognition", ArXiV:1412.5567. https://arxiv.org/abs/1412.5567
  - A Hannun (2017), "Sequence Modeling with CTC", *Distill*. https://distill.pub/2017/ctc
- Background reading
  - Y Miao et al (2015), "EESEN: End-to-end speech recognition using deep RNN models and WFST-based decoding", ASRU-2105. https://ieeexplore.ieee.org/abstract/document/7404790
  - A Maas et al (2015). "Lexicon-free conversational speech recognition with neural networks", NAACL HLT 2015, http://www.aclweb.org/anthology/N15-1038