

Final Project: Listen, Attend and Spell (LAS)

Nguyen Huu Kim 2020.06.24

Contents

- Introduction
- Implementation and Results

Introduction

- Automatic Speech Recognition Model (ASR)
- Input: Acoustic features $x = (x_1, ..., x_T)$

Output: Characters or phonemes $y = (\langle sos \rangle, y_1, ..., y_S, \langle eos \rangle)$

Objective: maximize likelihood

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i} P(y_i|\mathbf{x}, y_{< i})$$

Introduction

- Speller:
 - Attention-based LSTM transducer

$$c_i = \text{AttentionContext}(s_i, \mathbf{h})$$

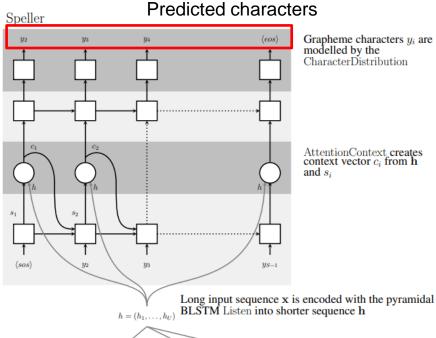
$$s_i = \text{RNN}(s_{i-1}, y_{i-1}, c_{i-1})$$

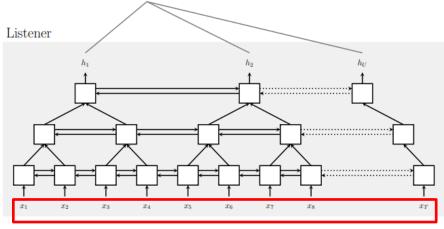
$$P(y_i | \mathbf{x}, y_{< i}) = \text{CharacterDistribution}(s_i, c_i)$$

- Listener
 - Pyramid BLSTM structure

$$h_i^j = \text{BLSTM}(h_{i-1}^j, h_i^{j-1})$$

$$h_i^j = \text{pBLSTM}(h_{i-1}^j, \left[h_{2i}^{j-1}, h_{2i+1}^{j-1}\right])$$





Acoustic feature



• Experiment settings:

• Dataset: VCTK Corpus

• Sampling rate: 16kHz

• Frame length: 50 ms

• Hop length: 12.5 ms

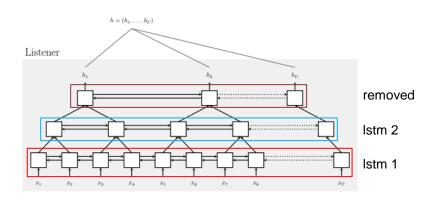
• Window: Hann

Hyperparameter setup

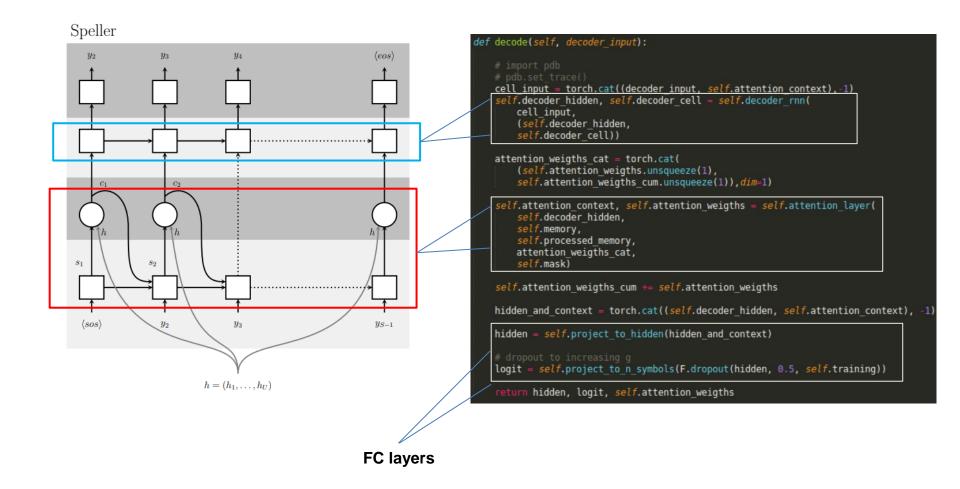
Transcription features	Phonemes
Acoustic features	80-dim log mel- spectrogram
Optimizer	Adam
Learning rate	10 ⁻³



- Different with original paper:
 - 2 layer pyramid BLSTM vs. 3 layer pyramid BLSTM
 - VCTK corpus vs. private Google dataset

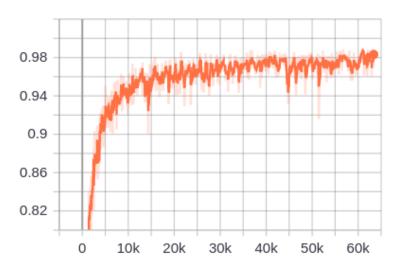


```
class Listener(nn.Module):
   def __init__(self, hparams):
    super(Listener, self).__init__()
       if hparams.spemb input:
           input dim * hparams.n_mel_channels * hparams.speaker_embedding_dim
           input dim = hparams.n mel channels
       self.lstm1 = nn.LSTM(input dim, int(hparams.audio encoder hidden dim / 2),
                            num layers=1, batch first=True, bidirectional=True)
        self.lstm2 = nn.LSTM(hparams.audio encoder hidden dim hparams.n frames per step encoder,
                            int(hparams.audio_encoder_hidden_dim / 2),
       self.concat hidden dim = hparams.audio encoder hidden dim hparams.n frames per step encoder
       self.n frames per step = hparams.n frames per step encoder
   def forward(self, x, input lengths):
       x sorted, sorted lengths, initial index = sort batch(x, input lengths)
       x_packed = nn.utils.rnn.pack_padded_sequence(
           x sorted, sorted lengths.cpu().numpy(), batch first=True)
       self.lstml.flatten parameters()
       outputs, = self.lstml(x packed)
       outputs, = nn.utils.rnn.pad packed sequence(
           outputs, batch first=True, total length=x.size(1)) # use total length make sure the recovered
       outputs = outputs.reshape(x.size(0), -1, self.concat hidden dim)
       output lengths = torch.ceil(sorted lengths.float() / self.n frames per step).long()
       outputs = nn.utils.rnn.pack_padded_sequence(
           outputs, output lengths.cpu().numpy() , batch first=True)
       self.lstm2.flatten_parameters()
       outputs, = self.lstm2(outputs)
       outputs, _ = nn.utils.rnn.pad_packed sequence(
           outputs, batch first=True)
       return outputs[initial_index], output_lengths[initial_index]
```

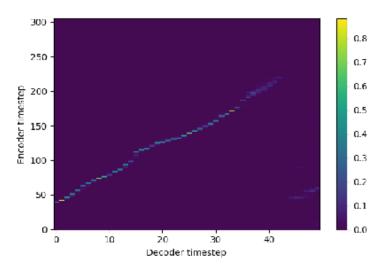


Results

Identification task: Accuracy 98%



Training accuracy



Attention plot

