

2019 음성인식 해커톤

Chapter 2. Deep Learning for Speech Recognition

Ready for your next

AI Hackathon?

참가 신청하기

https://campaign.naver.com/aihackathon_speech/

세번째 네이버 AI 해커톤, AI Speech

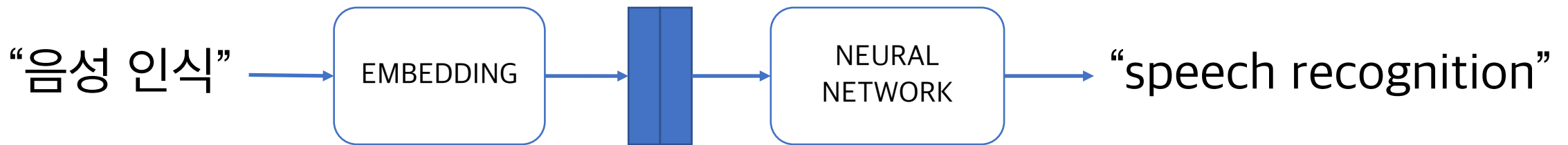
여러분들의 뜨거운 관심과 치열한 고민을 이어갈 수 있도록 이번에는 AI의 핵심 기술중의 하나인 음성인식 주제를 가지고 찾아왔습니다.

Overview

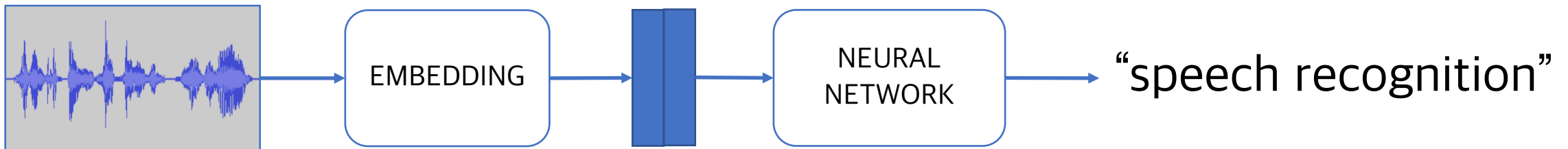
End-to-end Speech Recognition

Deep learning for 음성인식

- 기계 번역



- 음성 인식

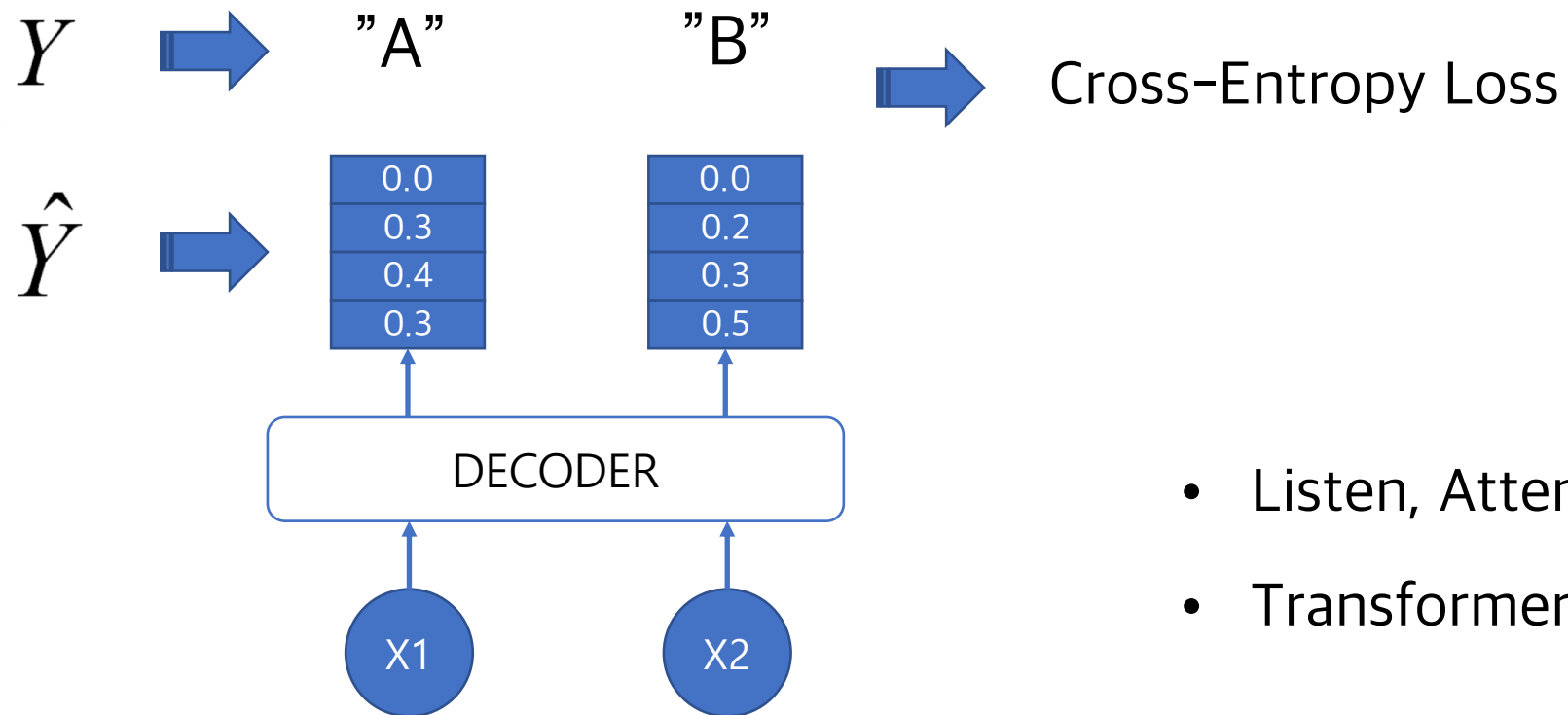


음성인식의 특징

- Long input sequence, short output sequence
 - 48,000 (3초 X 16,000 samples) => 20 (글자)
- Input and output sequences in same order
- Loss Function
 - CrossEntropy, MSE, ...
 - CTCLoss
- Performance
 - Levenshtein Distance (Edit Distance)

LOSS Function

Sequence to sequence model



LOSS Function

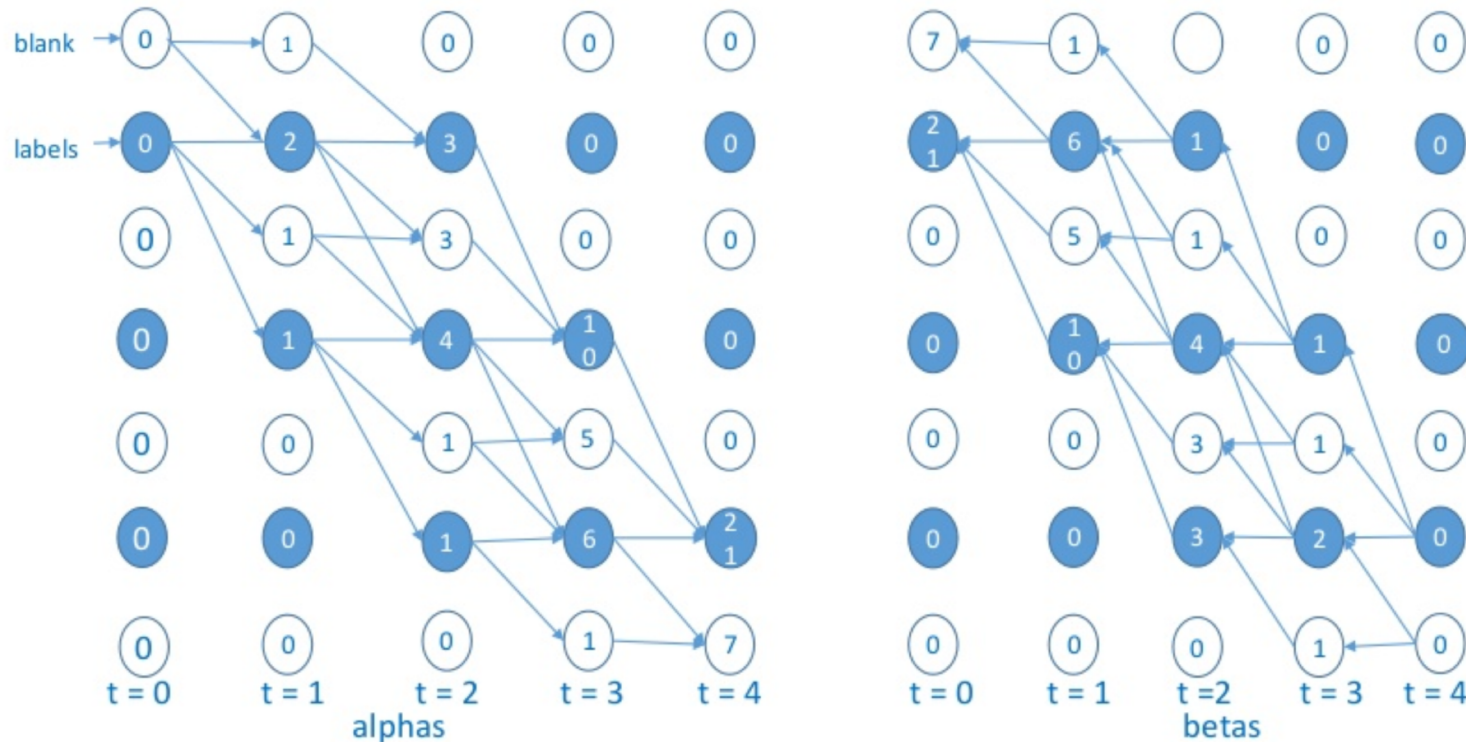
Sequence to sequence model

11 lines (7 sloc) | 306 Bytes

```
1 criterion = nn.CrossEntropyLoss(reduction='sum', ignore_index=PAD_token)
2
3 def train()
4     # ...
5     y_hat = model(features, features_lengths, scripts, teacher_forcing_ratio)
6     y_hat = torch.stack(y_hat, dim=1)
7     y = scripts[:, 1:]
8
9     loss = criterion(y_hat.view(-1, y_hat.size(-1)), y.view(-1))
10
```

LOSS Function

Connectionist Temporal Classification



$$p(\mathbf{l}|\mathbf{x}) = \sum_{t=1}^T \sum_{s=1}^{|\mathbf{l}|} \frac{\alpha_t(s)\beta_t(s)}{y_{l_s}^t}.$$

- Suhas Pillai, "Intelligent Handwriting Recognition_MIL_presentation_v3_final", Slideshare
- Alex Graves et al. "Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks", ICML, 2006

LOSS Function

Decoding Connectionist Temporal Classification

h h e ϵ ϵ | | | ϵ | | o

h e ϵ | ϵ | o

h e | | o

h e l l o

First, merge repeat characters.

Then, remove any ϵ tokens.

The remaining characters are the output.

- Deep Speech 2
- Awni Hannun, "Sequence Modeling With CTC", 10.23915/distill.00008

Performance

Edit Distance

정답	오	늘	은	날	씨	가	어	때	
인식	오	는		날	시	가	어	때	요

substitution = 2

deletion = 1

insertion = 1

$$ED = \frac{2 + 1 + 1}{\text{length of Ref.}} = \frac{4}{8} = 0.5$$

Performance

Edit Distance

4 lines (2 sloc) | 54 Bytes

```
1 import Levenshtein as Lev
```

```
2
```

```
3 CER = Lev.ratio(ref, hyp)
```

Baseline system

End-to-end Speech Recognition

Baseline system

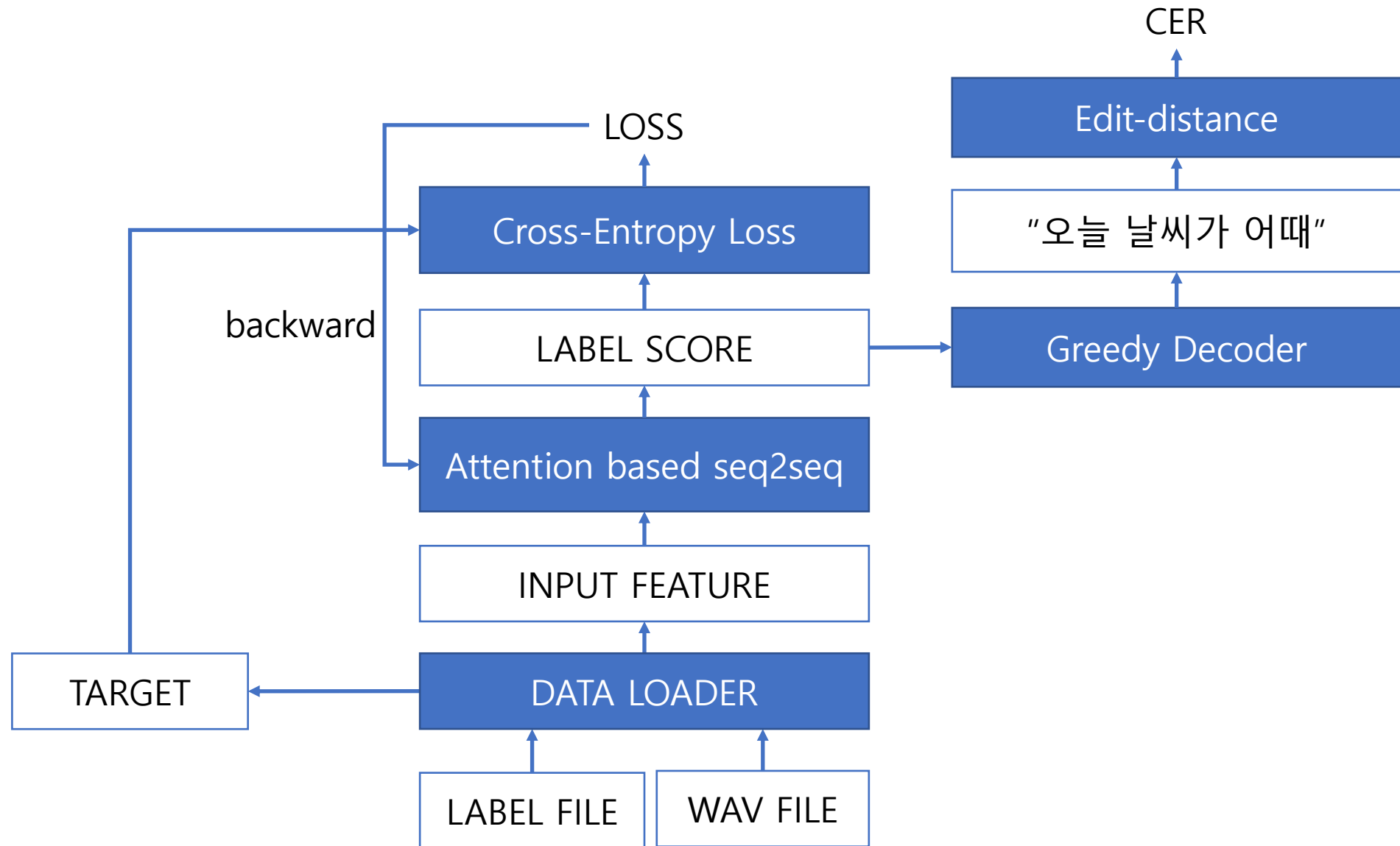
Dataset

약 3천 명의 사람들이 전화망을 통해 녹음한 3만개의 음성 파일과 정답 스크립트

FILE	CONTENTS
20190508192212_01074147469_378_0.wav	RIFF (little-endian) data, WAVE audio, Microsoft PCM, 16 bit, mono 16000 Hz
20190508192212_01074147469_378_0.script	크리스마스날 영업하나요?
20190508192212_01074147469_378_0.label	533 94 384 623 384 706 662 211 532 690 62 661 123 662

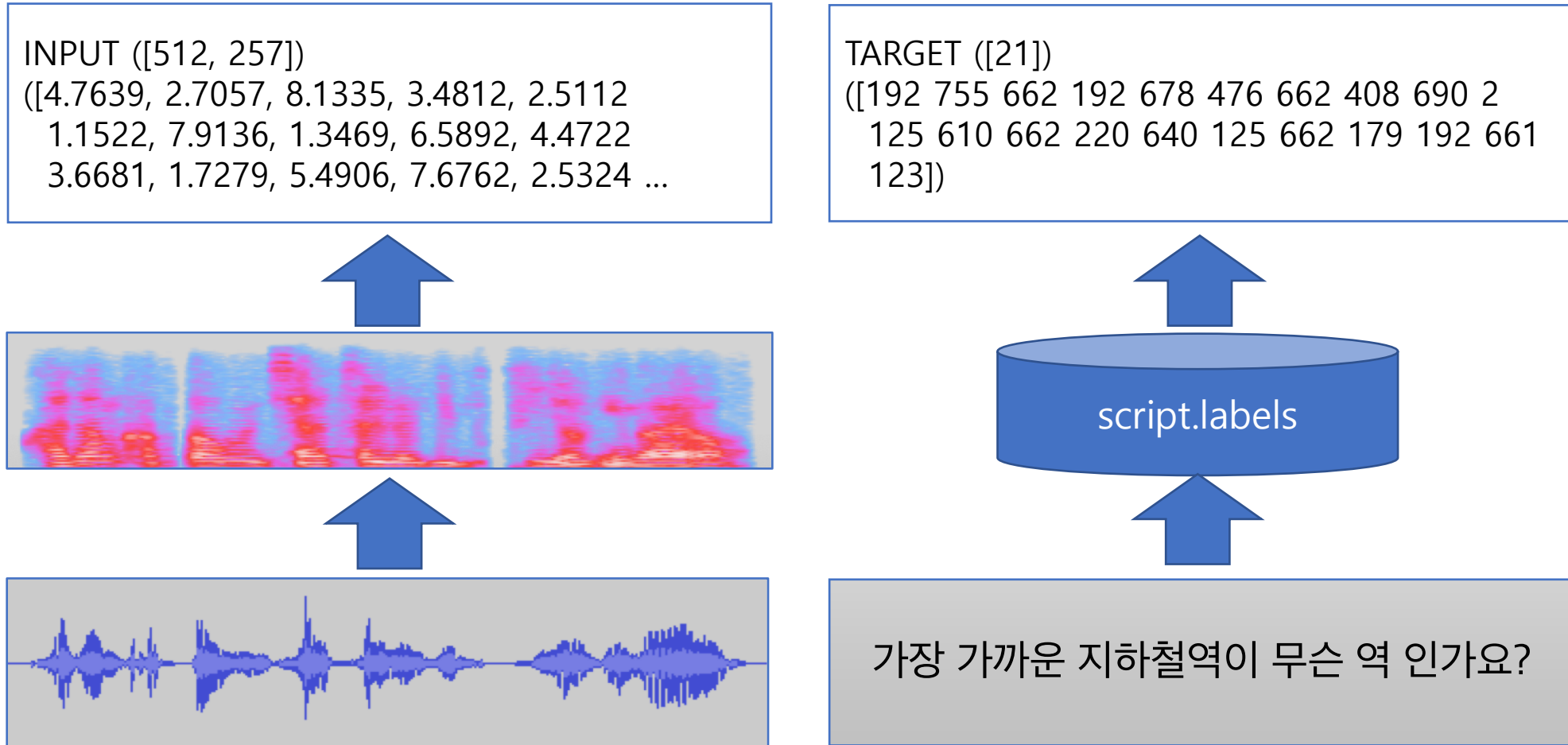
META-FILE	CONTENTS
train/train_data/data_list.csv	3만개의 파일 목록
script.labels	Character와 label 간의 매핑 테이블 정보

Baseline system



Baseline system

Data Loader



Baseline system

Feature extractor

20 lines (16 sloc) | 681 Bytes

```
1  def get_feature(filepath, feature_size):
2      (rate, width, sig) = wavio.readwav(filepath)
3
4      sig = sig.ravel()
5
6      stft = torch.stft(torch.FloatTensor(sig),
7                        512,
8                        hop_length=int(0.01*16000),
9                        win_length=int(0.030*16000),
10                       window=torch.hamming_window(int(0.030*16000)),
11                       center=False,
12                       normalized=False,
13                       onesided=True)
14     stft = (stft[:, :, 0].pow(2) + stft[:, :, 1].pow(2)).pow(0.5);
15     spect = stft.numpy();
16     spect = torch.FloatTensor(spect)
17
18     spect = torch.FloatTensor(spect).transpose(0, 1)
19     return spect
```


Baseline system

DataLoader를 python thread로 구현해서 속도 향상

23 lines (18 sloc) | 674 Bytes

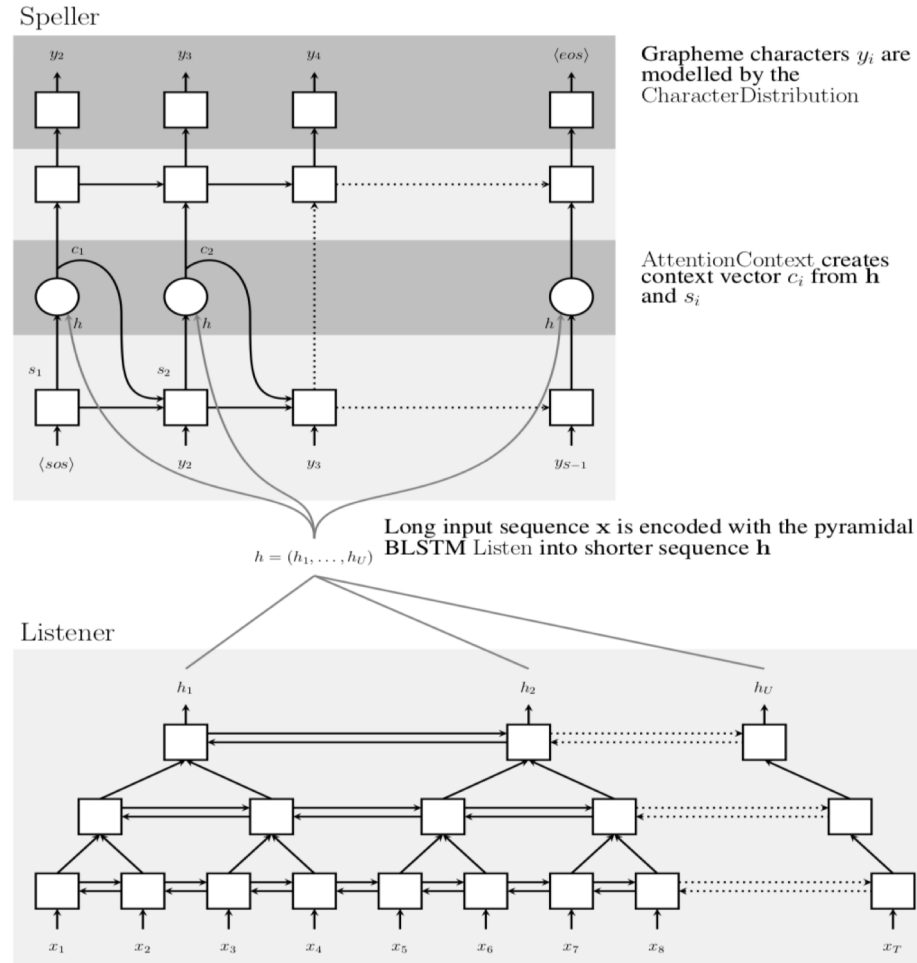
```
1 class BaseDataLoader(threading.Thread):
2     def run(self):
3         while True:
4             items = self.read_batches()
5             if len(items) == 0:
6                 break
7
8             random.shuffle(items)
9             batch = self.collate_fn(items)
10            self.queue.put(batch)
11
12 class MultiLoader():
13     def __init__(self, dataset_list, queue, batch_size, worker_size):
14         self.worker_size = worker_size
15         self.loader = list()
16         for i in range(self.worker_size):
17             self.loader.append(BaseDataLoader(dataset_list[i], queue, batch_size, i))
18
19     def start(self):
20         for i in range(self.worker_size):
21             self.loader[i].start()
```

End-to-end models

End-to-end Speech Recognition

Listen, Attend and Spell

Sequence to sequence with Attention



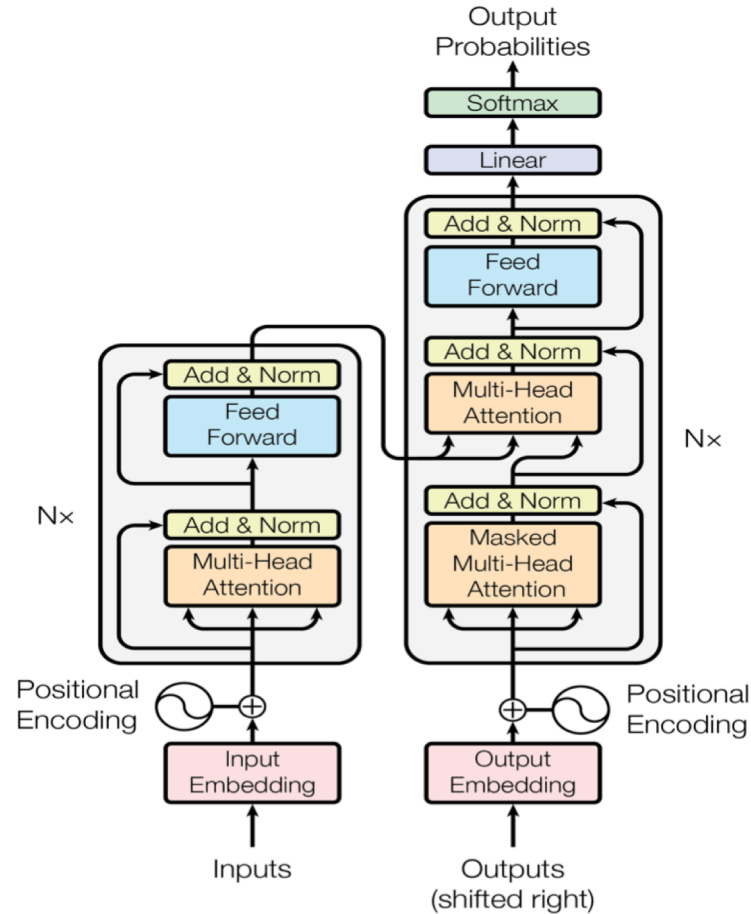
cross-entropy loss

input context에 대해서
global attention 사용

pyramidal LSTM을 사용해서
input sequence를 ¼로 축소

Transformer

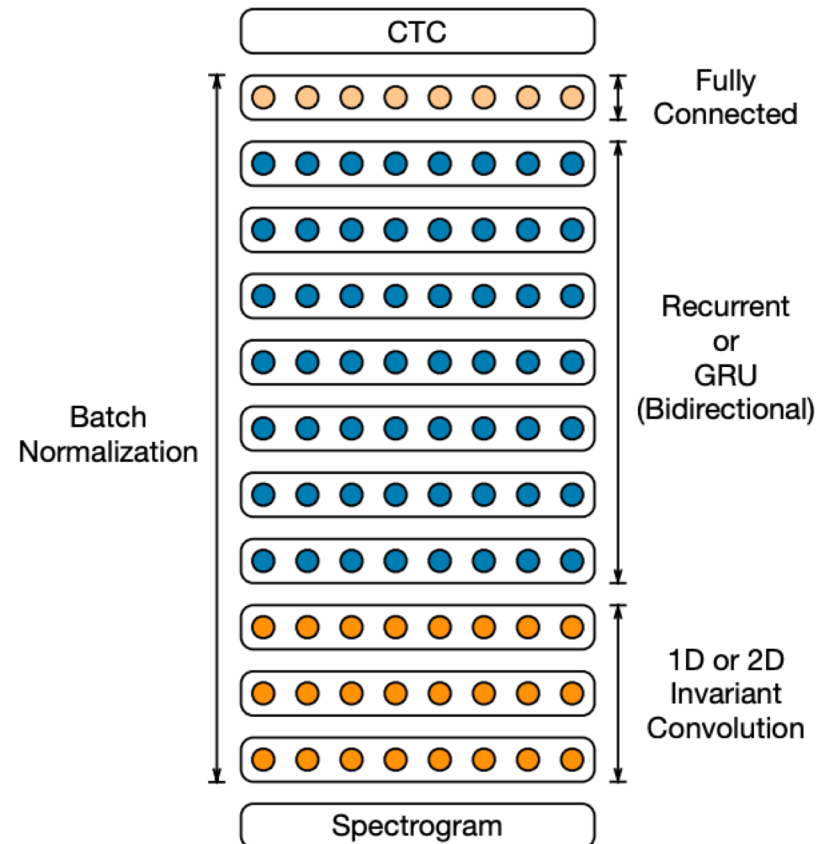
Sequence to sequence with Attention



RNN 대신에, Positional Encoding과 Feed-Forward layer를 이용해서 Parallel Encoding & Decoding, 학습 속도 향상

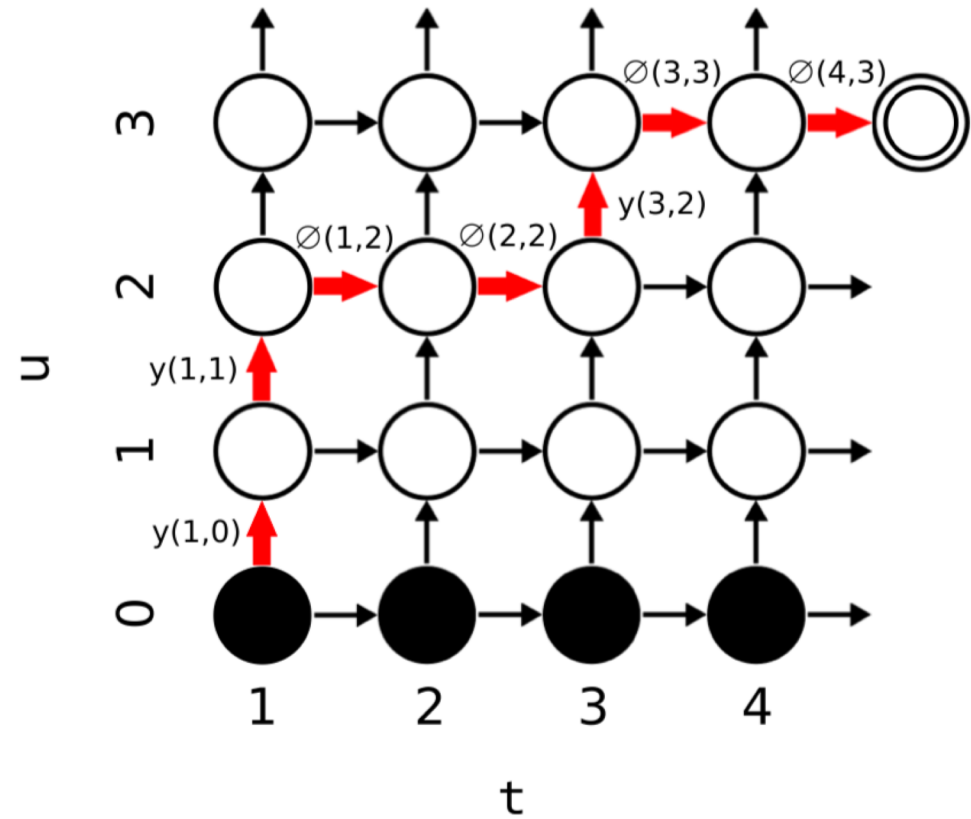
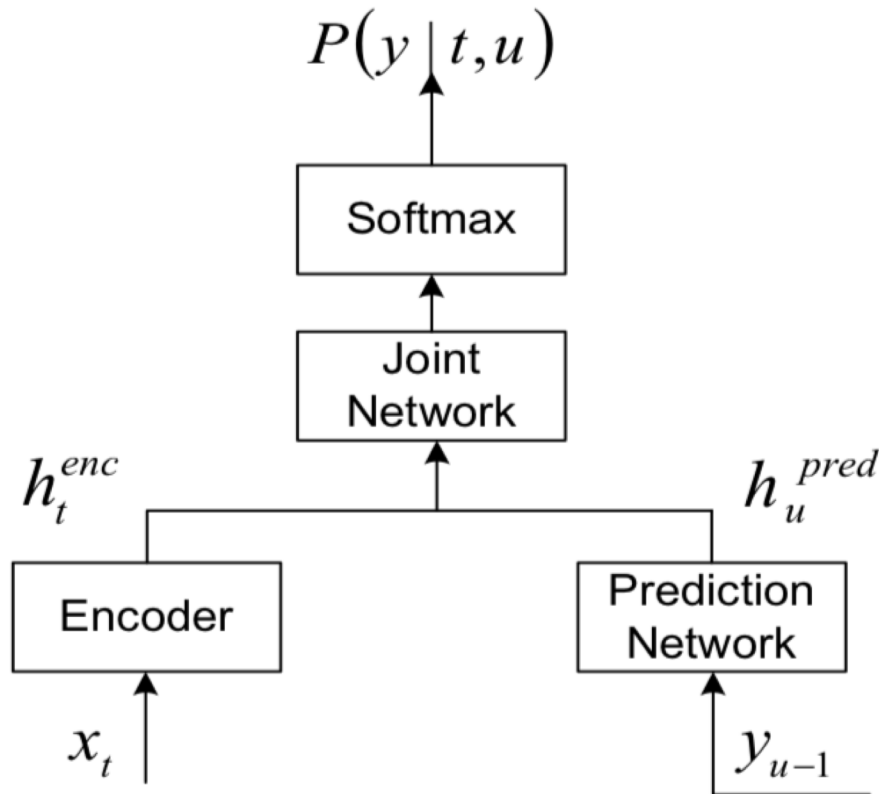
Encoder with Forward-backward algorithm

Deep Speech2



Encoder with Forward-backward algorithm

RNN Transducer



- Alex Graves, “Sequence Transduction with Recurrent Neural Networks”, ICML, 2012
- Senmao Wang et al, “Exploring RNN-Transducer for Chinese Speech Recognition”, arXiv:1811.05097

End of Document