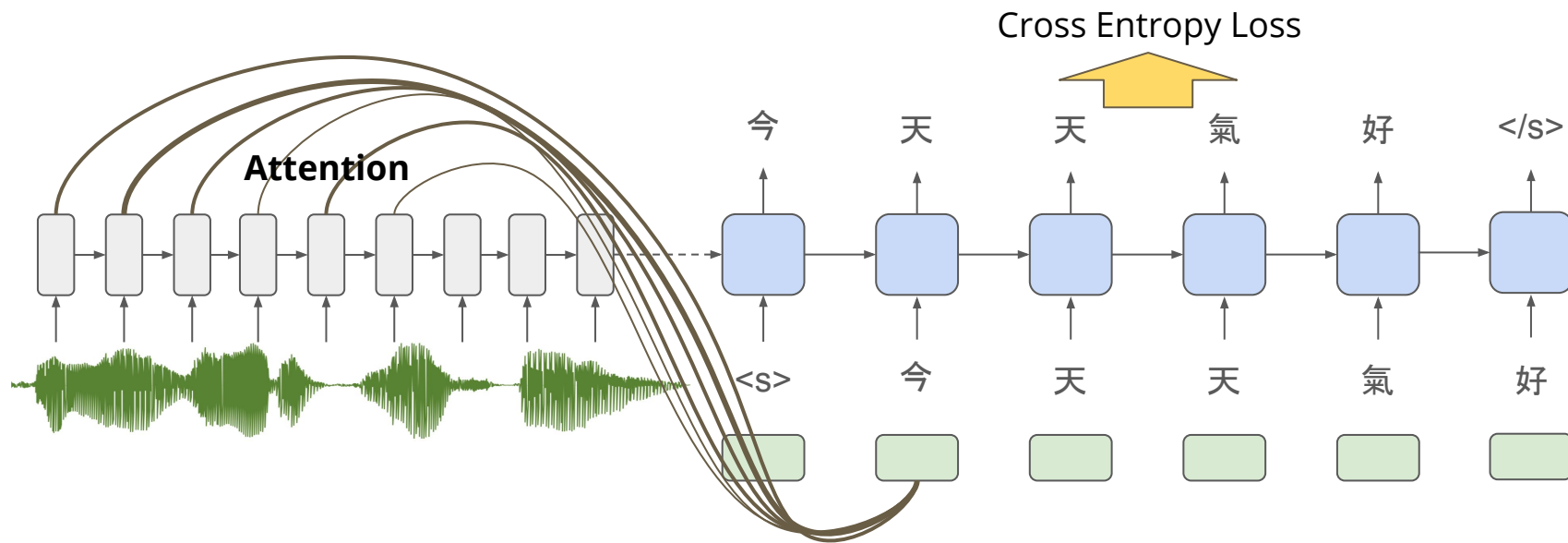# DLHLP - HW1
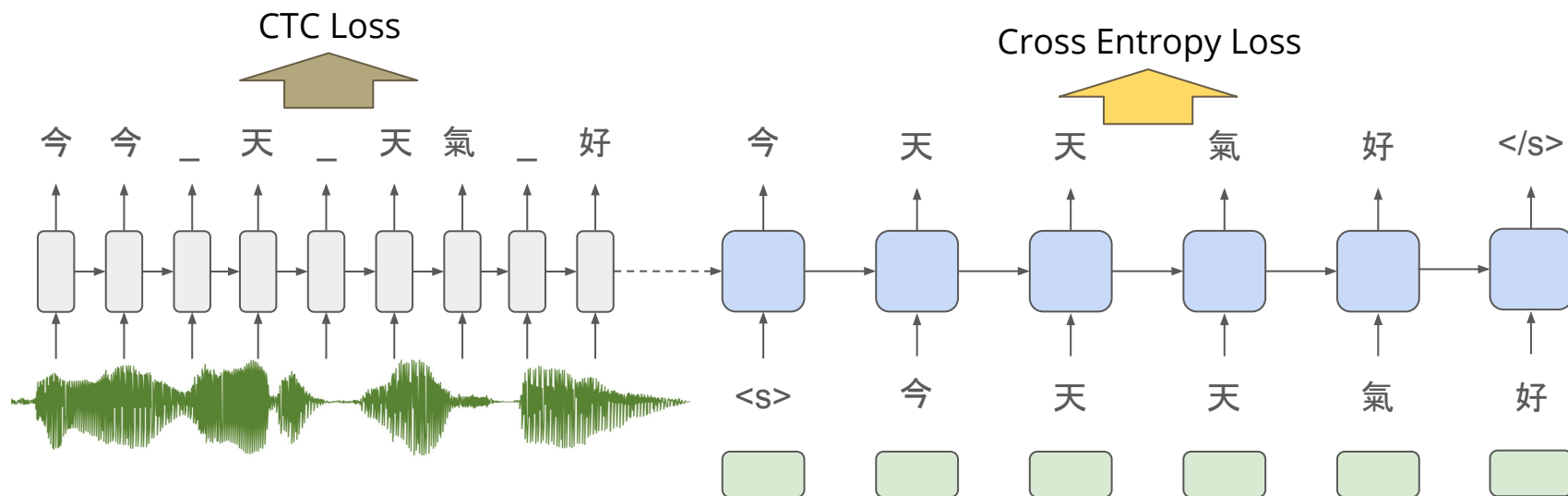# End-to-end Speech Recognition

TA: 莊永松、柯上優

dlhlp.ta@gmail.com
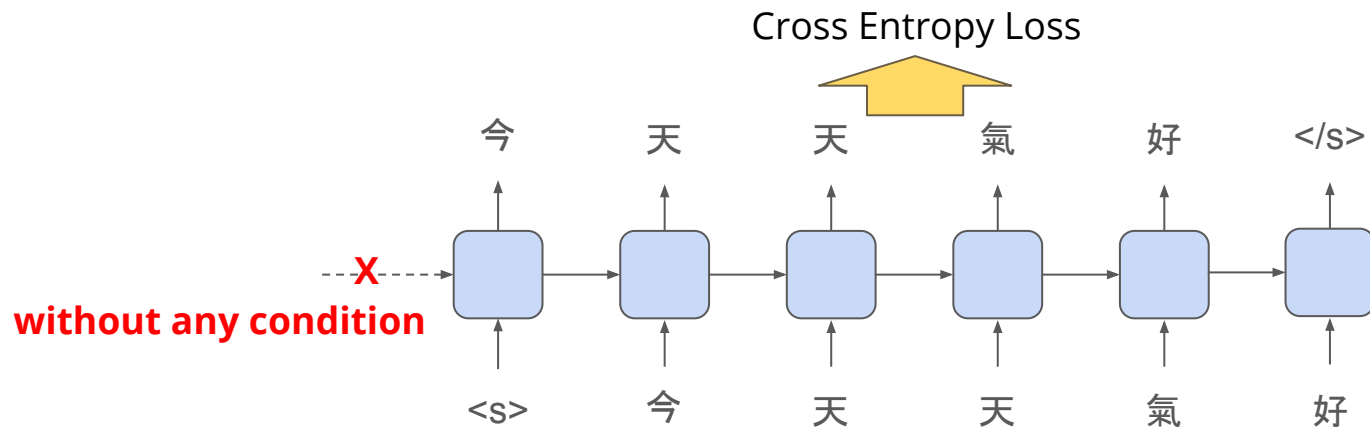
# Seq2Seq attention-based ASR

# Joint CTC-attention based ASR

# RNN-LM (for re-scoring)

Cross Entropy Loss



X
**without any condition**

今　天　天　氣　好　</s>

<s>　今　天　天　氣　好

# Dataset

To make HW1 easier…

- We use a small 10-hours Mandarin Chinese audio corpus
- All from a single speaker (no speaker varient problem)
- We transform text to ㄅㄆㄇㄈ to make vocabulary size smaller and help the model to converge faster
- You would not need to train a "real-world" ASR model
- About **2~4 hours** to converge on K80, while real-world ASR may need weeks to converge

# Download Data

- Download Data(2.88G)
  - https://drive.google.com/file/d/1daFU8tPPUyhN7Fc6JUTohEfHXIn6ZDgq/view?usp=sharing
  - train: 8000, dev: 1000, test: 1000
  - the transcript file in testing set is not removed for convenience, but all the answer is replaced with ㄅㄆㄇㄈ…

Hint: In Linux, use `bash get_dataset.sh`
https://github.com/DLHLP2020/hw1-speech-recognition/blob/master/get_dataset.sh

```
/DLHLP
├ /train
│ ├ 000001.wav
│ ├ ......
│ ├ 008000.wav
│ └ bopomo.trans.txt
├ /dev
│ ├ 008001.wav
│ ├ ......
│ ├ 009000.wav
│ └ bopomo.trans.txt
├ /test
│ ├ 009001.wav
│ ├ ......
│ ├ 010000.wav
│ └ bopomo.trans.txt
└ text-data.txt
```
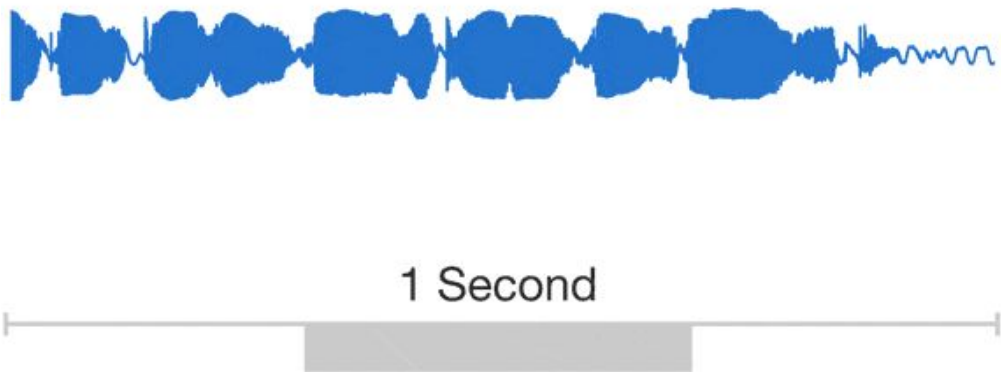
# Extract features from audio file w/ torchaudio

**Step 1.** **Load WAV File**

```
waveform, sample_freq = torchaudio.load("000001.wav")
waveform.shape: (1, 116400)
sample freq: 48000
```



1 Second
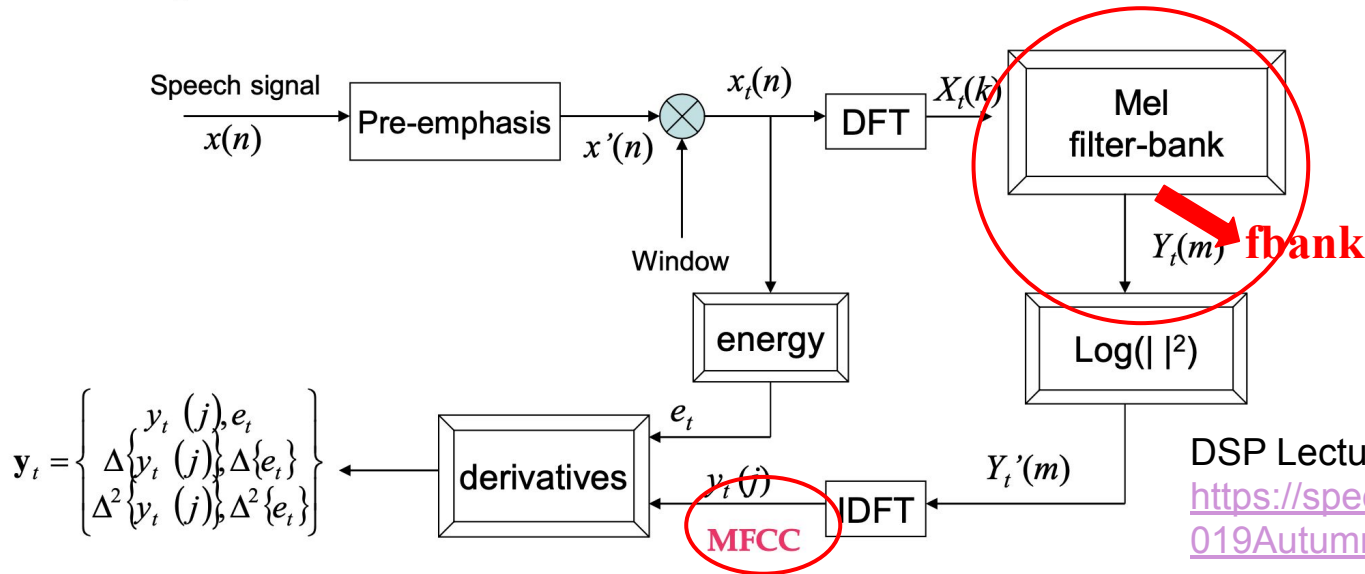
# Step 2. Get fbank or MFCC from waveform

- **Mel-Frequency Cepstral Coefficients (MFCC)**
  - Most widely used in the speech recognition
  - Has generally obtained a better accuracy at relatively low computational complexity
  - The process of MFCC extraction :



DSP Lecture 7:
https://speech.ee.ntu.edu.tw/DSP2019Autumn/Slides/7.0.pptx

## Step 2. Get Fbank or MFCC from waveform

**MFCC**
```
feature = torchaudio.compliance.kaldi.mfcc(waveform,
                            sample_frequency=sample_freq)
```
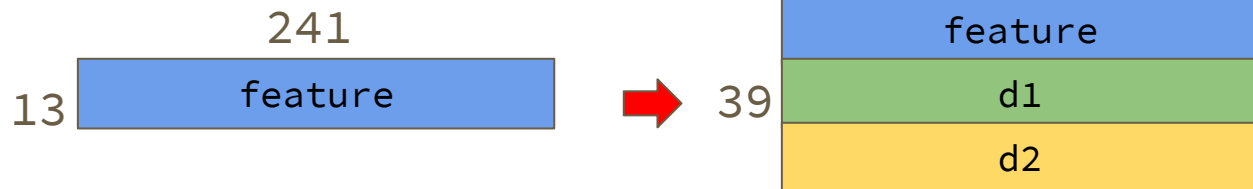
feature.shape: (241, 13)

**Fbank**
```
feature = torchaudio.compliance.kaldi.fbank(waveform,
                            sample_frequency=sample_freq,
                            num_mel_bins=40)
```
feature.shape: (241, 40)

**Step 3. Add Deltas**

241
13 | feature

241
39 | feature / d1 / d2

```
d1 = torchaudio.functional.compute_deltas(feature)
d2 = torchaudio.functional.compute_deltas(d1)
feature = torch.cat([feature, d1, d2], dim=-1)
feature.shape: (241, 39) or (241, 120)
```
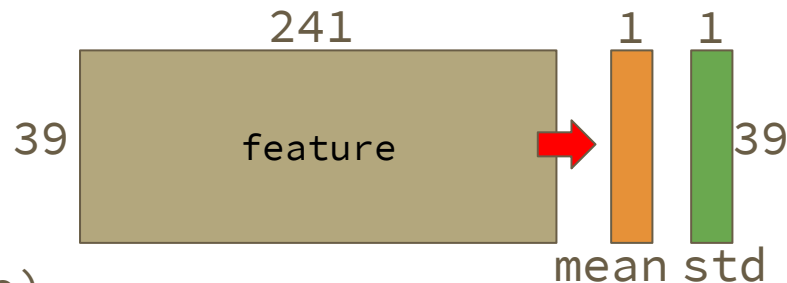
**Step 4. CMVN (Normalization)**

241
39 | feature → 1 1 / mean std / 39

```
eps = 1e-10
mean = feature.mean(0, keepdim=True)
std = feature.std(0, keepdim=True)
feature = (feature - mean) / (std + eps)
```

# Sample code (strongly recommended)

https://github.com/Alexander-H-Liu/End-to-end-ASR-Pytorch

A completed & stable ASR implementation (Thanks to 劉浩然學長!)

*Please read README before using it!*

# TODO: Add new Dataset

- See the example corpus/librispeech.py
- Create a copy corpus/dlhlp.py. Replace all "Libri" to "Dlhlp" in the file.
  1. open each **.flac** files → open each **.wav** files
  2. def read_text()
     - origin: src_file = '-'.join(file.split('-')[:-1])+'.trans.txt'
     - ours:   src_file = file.rsplit('/', 1)[0]+'/bopomo.trans.txt'
     - Base on the format of our dataset

```
28  class DlhlpDataset(Dataset):
29      def __init__(self, path, split, tokenizer, bucket_size, ascending=False):
30          # Setup
31          self.path = path
32          self.bucket_size = bucket_size
33
34          # List all wave files
35          file_list = []
36          for s in split:
37              split_list = list(Path(join(path, s)).rglob("*.wav"))
38              assert len(split_list) > 0, "No data found @ {}".format(join(path,s))
39              file_list += split_list
```

```
15  def read_text(file):
16      '''Get transcription of target wave file,
17          it's somewhat redundant for accessing each txt multiplt times,
18          but it works fine with multi-thread'''
19      src_file = file.rsplit('/', 1)[0]+'/bopomo.trans.txt'
20      idx = file.split('/')[-1].split('.')[0]
21
22      with open(src_file, 'r') as fp:
23          for line in fp:
24              if idx == line.split(' ')[0]:
25                  return line[:-1].split(' ', 1)[1]
26
```

# TODO: Import new Dataset

- In the [src/data.py](src/data.py)
  - def create_dataset
  - def create_textset

```python
64  def create_dataset(tokenizer, ascending, name, path, bucketing, batch_size,
65                      train_split=None, dev_split=None, test_split=None):
66      ''' Interface for creating all kinds of dataset'''
67
68      # Recognize corpus
69      if name.lower() == "librispeech":
70          from corpus.librispeech import LibriDataset as Dataset
71      elif name.lower() == "dlhlp":
72          from corpus.dlhlp import DlhlpDataset as Dataset
73      else:
74          raise NotImplementedError
75
```

```python
111  def create_textset(tokenizer, train_split, dev_split, name, path, buck
112      ''' Interface for creating all kinds of text dataset'''
113      msg_list = []
114
115      # Recognize corpus
116      if name.lower() == "librispeech":
117          from corpus.librispeech import LibriTextDataset as Dataset
118      elif name.lower() == "dlhlp":
119          from corpus.dlhlp import DlhlpTextDataset as Dataset
120      else:
121          raise NotImplementedError
```

# TODO: Prepare vocab-count file

```
$ python util/generate_vocab_file.py
                --input_file <input_text_file>
                --mode character
                --output_file <output_vocab_file>
```

we provide "text-data.txt" in the dataset as input_text file.

# TODO: Write your own config file (1/3)

- See example config/libri/asr_example.yaml
- Create copys as you want to train a new model (e.g. asr_dlhlp.yaml)
  - `path`: where you unzip the data (e.g. 'data/DLHLP')
  - `train/dev_split`: dir name for train/dev under 'path'
  - audio feature setting: follow the original setting
  - mode: 'character'
  - `vocab_file`: the vocab file your have prepared

```yaml
 1  data:
 2    corpus:
 3      name: 'Dlhlp'
 4      path: 'data/DLHLP'
 5      train_split: ['train']
 6      dev_split: ['dev']
 7      bucketing: True
 8      batch_size: 16
 9    audio:
10      feat_type: 'fbank'
11      feat_dim:  40
12      frame_length: 25
13      frame_shift: 10
14      dither: 0
15      apply_cmvn: True
16      delta_order: 2
17      delta_window_size: 2
18    text:
19      mode: 'character'
20      vocab_file: 'bopomo_vocab_file'
```

# TODO: Write your own config file (2/3)

- Training Hyperparams
  - set valid_step to 500 (step for one epoch)
  - max_step: 12k step is enough actually
  - teather forcing: always use it
    => tf_start: 1.0, tf_end:1.0
  - optimizer and lr and eps:
    just follow the original settings

```
22  hparas:
23   valid_step: 500
24   max_step: 12001
25   tf_start: 1.0
26   tf_end: 1.0
27   tf_step: 500000
28   optimizer: 'Adadelta'
29   lr: 1.0
30   eps: 0.00000001
31   lr_scheduler: 'fixed'
32   curriculum: 0
```

# TODO: Write your own config file (3/3)

- Model Architecture:
  - A thinner model (1~2 layer LSTM) is enough
  - ctc_weight:
    - set to 0.0 to train seq2seq without CTC
    - set between 0.0~1.0 to jointly optimize for CTC + seq2seq
  - Other settings: just follow the original settings

```yaml
34  model:
35    ctc_weight: 0.0
36    encoder:
37      prenet: 'vgg'
38      # vgg: True
39      module: 'LSTM'
40      bidirection: True
41      dim: [512,512]
42      dropout: [0,0]
43      layer_norm: [False,False]
44      proj: [True,True]
45      sample_rate: [1,1]
46      sample_style: 'drop'
47    attention:
48      mode: 'loc'
49      dim: 300
50      num_head: 1
51      v_proj: False
52      temperature: 0.5
53      loc_kernel_size: 100
54      loc_kernel_num: 10
55    decoder:
56      module: 'LSTM'
57      dim: 512
58      layer: 1
59      dropout: 0
```

# TODO: Write your own config file for LM

- See example config/libri/lm_example.yaml
- Need to modify:
  - path
  - train/dev_spit
  - vocab_file
- You can also prepare your own ㄅㄆㄇㄈcorpus larger much larger than training set!
  - just replace [`train`] with [`xxx.txt`] where `xxx.txt` is your collected corpus

```yaml
1  data:
2    corpus:
3      # The following depends on corpus
4      name: 'dlhlp'
5      path: 'data/DLHLP'
6      train_split: ['train']
7      dev_split: ['dev']
8      bucketing: True
9      batch_size: 32
10   text:
11     mode: 'character'
12     vocab_file: 'bopomo_vocab_file'
13
```

# Train model!

```
ASR:
$ python3 main.py --config config/dlhlp/asr_dlhlp.yaml

LM:
$ python3 main.py --config config/dlhlp/lm_dlhlp.yaml --lm
```

# Tensorboard

```
$ tensorboard --logdir log/ --port <port_you_want>
```

Training curve
(loss &WER / train&dev)

# Alignment

# Output example

SCALARS  IMAGES  **TEXT**

# Test your model - seq2seq

- See example config/libri/decode_example.yaml
- set max_len_ratio: 0.30
- set ctc_weight: 0.0
- set beam size 2~20 (1 not support)

If you want to use external LM:

- set lm_config&lm_path
- set lm_weight>0.0

```yaml
1   # Most of the parameters will be imported from t
2   src:
3     ckpt: 'ckpt/asr_dlhlp_bopomo_sd0/best_att.pth'
4     config: 'config/dlhlp/asr_dlnlp_bopomo.yaml'
5   data:
6     corpus:
7       name: 'Dlhlp'
8       dev_split: ['dev']
9       test_split: ['test']
10  decode:
11    beam_size: 2
12    min_len_ratio: 0.01
13    max_len_ratio: 0.30
14    lm_path: 'lm_dlhlp_sd0/best_ppx.pth'
15    lm_config: 'config/dlhlp/lm_dlhlp.yaml'
16    lm_weight: 0.3
17    ctc_weight: 0.0
```

# Test your model - CTC

- Original code has some problems for jointly decode by seq2seq&CTC
- Add few lines of code to support decode only by CTC (works for ctc_weight=1.0)
- src/decode.py line 101

```python
 95            # CTC decoding
 96        if self.apply_ctc:
 97            ctc_output = F.log_softmax(
 98                self.asr.ctc_layer(encode_feature), dim=-1)
 99            ctc_prefix = CTCPrefixScore(ctc_output)
100            ctc_state = ctc_prefix.init_state()
101            if self.ctc_w == 1.0:
102                output_seq = ctc_output[0].argmax(dim=-1)
103                hypothesis = [Hypothesis(decoder_state=dec_state, output_seq=output_seq,
104                                         output_scores=[0]*len(output_seq), lm_state=None, ctc_prob=0,
105                                         ctc_state=ctc_state, att_map=None)]
106            return hypothesis
```

# Patch for pure CTC decode (add to src/decode.py line 101)

Simply Greedy Decode!

```python
    if self.ctc_w == 1.0:
        output_seq = ctc_output[0].argmax(dim=-1)
        hypothesis = [Hypothesis(decoder_state=dec_state, output_seq=output_seq,
                        output_scores=[0]*len(output_seq), lm_state=None, ctc_prob=0,
                        ctc_state=ctc_state, att_map=None)]
        return hypothesis
```

Actually, pure CTC decode also needs beam search.
*It could be a good bonus topic for you to implement it.*
https://towardsdatascience.com/beam-search-decoding-in-ctc-trained-neural-networks-5a889a3d85a7

# Test your model - CTC

- set ctc_weight = 1.0
- external LM is not supported to jointly decode with CTC in our code (*But you can implement it as bonus!*)

```
1   # Most of the parameters will be imported from the
2   src:
3     ckpt: 'ckpt/asr_dlhlp_bopomo_sd0/best_ctc.pth'
4     config: 'config/dlhlp/asr_dlnlp_bopomo.yaml'
5   data:
6     corpus:
7       name: 'Dlhlp'
8       dev_split: ['dev']
9       test_split: ['test']
10  decode:
11    beam_size: 2
12    min_len_ratio: 0.01
13    max_len_ratio: 0.30
14    lm_path:
15    lm_config:
16    lm_weight: 0.0
17    ctc_weight: 1.0
```

# Test!

```
$ python3 main.py --config <config file> --test --njobs 8
```

- result will be produced at:
  - result/<config file>_dev_output.csv
  - result/<config file>_dev_beam.csv 用不到
  - result/<config file>_test_output.csv
  - result/<config file>_test_beam.csv 用不到
- The format of output file is `<id> <predicted seq> <truth seq>` (line by line)
- The truth sequence of testing set was replaced with ㄅㄆㄇㄈ.
- `--njobs` decides the number of threads used for decoding, very important in terms of efficiency. You can set it higher as your machine have more cores.

# Process the output from CTC

CTC has repeated tokens
Process it by yourself!



Or you can add additional code to the patch in page 25 to process the
'output_seq' directly before it is written to file.

# Eval your output (dev set)

```
$ python3 eval.py --file result/<config_file>_dev_output.csv
```
*hint: python2 is not supported here (unicode would be treat as two char)*

```
============ Result of result/decode_bopomofo_beam2_lm0.3_dev_output.csv ============
--------------------------------------------------------------------------
| Statics               |   Truth    |  Prediction  | Abs. Diff.      |
--------------------------------------------------------------------------
| Avg. # of chars       |   66.99    |   66.96      |  0.37           |
| Avg. # of words       |   17.14    |   17.12      |  0.02           |
------------------report this!---------------------------------------
| Error Rate (%)| Mean             |  Std.        | Min./Max.     |
--------------------------------------------------------------------------
| Character     | 2.1157           |  2.43        | 0.00/33.33    |
| Word          | 7.0030           |  7.48        | 0.00/100.00   |
--------------------------------------------------------------------------
Note : If the text unit is phoneme, WER = PER and CER is meaningless.
```

# Submit your result (test set)

- Extract the result in kaggle format
- Sample script:
  https://github.com/DLHLP2020/hw1-speech-recognition/blob/master/format.py
- usage:
  ```
  python3 format.py result/<config>_test_output.csv answer.csv
  ```
- upload answer.csv to kaggle

# Kaggle rules

- website: https://www.kaggle.com/c/dlhlp2020spring-asr
- Your team name should be in [team_github_id]_[any_string] e.g. daikin_大金
- 5 submission per team & per day
- Using any extra kaggle account to submit is cheating!

# Kaggle Evaluation Metric

Mean **Levenshtein Distance** calculated for each sentence in char-level.

- the minimum number of single-character edits (insertions, deletions or substitutions) required to change one word into the other.

$$LevDistance = N_{ins} + N_{del} + N_{sub}$$

Word Error Rate:

$$WER = \frac{N_{ins} + N_{del} + N_{sub}}{N_{target-length}}$$

# Baselines (5 points)

Public Simple Baseline: 1.394 (2 pt)
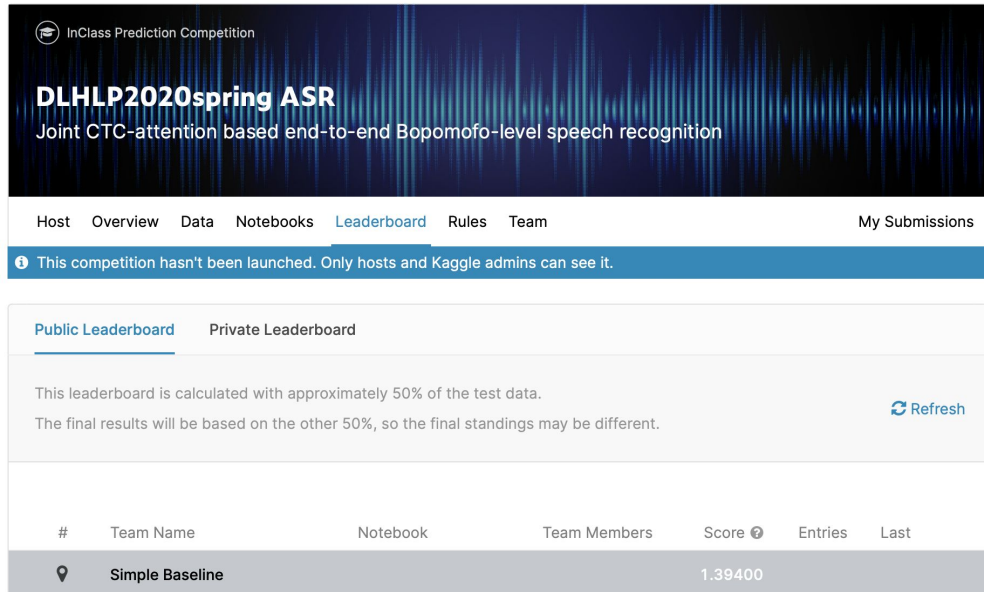Public Strong Baseline: 1.278 (1 pt)

Private Simple Baseline: (?)    (1 pt)
Private Strong Baseline: (?)    (1 pt)

p.s. private score will be shown after kaggle deadline

Both would be easy to beat~

https://www.kaggle.com/c/dlhlp2020spring-asr

InClass Prediction Competition

**DLHLP2020spring ASR**
Joint CTC-attention based end-to-end Bopomofo-level speech recognition

| Host | Overview | Data | Notebooks | Leaderboard | Rules | Team | My Submissions |

ⓘ This competition hasn't been launched. Only hosts and Kaggle admins can see it.

Public Leaderboard     Private Leaderboard

This leaderboard is calculated with approximately 50% of the test data.
The final results will be based on the other 50%, so the final standings may be different.

↻ Refresh

| # | Team Name | Notebook | Team Members | Score ❓ | Entries | Last |
|---|-----------|----------|--------------|----------|---------|------|
| 📍 | Simple Baseline | | | 1.39400 | | |

# Report Questions (1/2)

In 1. ~ 3., just decode without language model.
In 1. ~ 4., set beam size = 2 for speeding.

1. Train a **seq2seq attention-based** ASR model. Paste the learning curve and alignment plot from tensorboard. Report the CER/WER of dev set and kaggle score of testing set. (2 points)

2. Repeat 1. by training a **joint CTC-attention** ASR model (decoding with seq2seq decoder). Which model converges faster? Explain why. (2 points)

3. Use the model in 2. to **decode only in CTC** (ctc_weight=1.0). Report the CER/WER of dev set and kaggle score of testing set. Which model performs better in 1. 2. 3.? Explain why. (2 points)

# Report Questions (2/2)

4. Train an **external language model.** Use it to help the model in 1. to decode. Report the CER/WER of dev set and kaggle score of testing set. (2 points)
5. Try decoding the model in 4. with **different beam size** (e.g. 2, 5, 10, 20, 50). Which beam size is the best? (2 points)

**Bonus. Other Improvement and Innovation**

- If you don't have anything to share, you can just say 'Nothing'. It will be fine.
- Guidiance
  - Read new papers
  - Browse github
  - Apply some cool tips

# Scoring: Submission

- Kaggle Deadline: 3/22(Sun) 23:59
- Github submission Deadline: 3/25(Wed) before class
- Create a folder 'hw1' under your team Github repo
- 'hw1/' contains:
  - report.pdf
  - reproduce.sh
  - other files and directories
- Report Template: https://docs.google.com/document/d/1NyIgXrIai9Lgysplh9p742zgn2j3cE3B7epqiLDAesE/
- We restrict Python version 3.6.8 and must compatible with these package
- Scoring
  - report 10pts
  - kaggle 5pts (over baseline + successfully reproduce)

# Github maximum capacity

- within 100MB
- use Dropbox to put your model, use 'wget' to download
- Dropbox Tutorial:
  https://docs.google.com/presentation/d/1SsIeIij9ZOEN_TGdbAS1oWcI6bT1uSTI6b5__u2wdDc/edit?usp=sharing
- your shell script files should be able to download the model automatically

# Scoring: Reproduce

- `bash reproduce.sh $1 $2`
  - `` `$1` `` is the audio dataset directory (e.g. data/DLHLP)
  - `` `$2` `` is name of the output prediction csv file (e.g. ans.csv )
- This script should produce the same result as your **best submission on public leaderboard.**
- Hint: use sed to modify the data path (e.g. 'data/DLHLP') in your config
- It HAVE to automatically download EVERYTHING you want to wget
- If your code reproduce fail, your **CANNOT** get kaggle score (5%)
- Prepare it carefully!

# Scoring: Bonus

- Presentation in class
- Selection criteria
    1. Innovation
    2. Different ways compare to LAS
        - Exiting github is valid, but you have to understand and explain it.
    3. Good performance
- 1 extra pts
- The team quota and the presentation time will be announced based on the time we have.

# Late submission policy

- You can submit file until 3 days after deadline
- The score will be calculated as:

$$\text{score}_{\text{final}}(\text{hr}) = \begin{cases} \text{score}_{\text{original}} \times 0.985^{\text{hr}} & , \text{hr} \leqslant 72 \\ 0 & , \text{hr} > 72 \end{cases}$$

- Late submision form would be anounced after deadline

**FB Group:**
Deep Learning for Human
Language Processing
(2020,Spring)

# Q&A

dlhlp.ta@gmail.com