# HW1改進與調整

交大電機碩一 張詔揚

# Espnet 資料目錄

- espnet/            # Python modules
- egs2/              # The complete recipe for each corpora
-     [name]/
-        asr1/        # ASR recipe
-            downloads/
-                resource_aishell/
-                    -speaker.info
-                    -lexicon.txt
-                data_aishell        #放訓練資料
-                    …

# 訓練資料準備

- data_aishell/
-    wav/
-       train/
-          -global/       # 訓練音檔
-       test/
-          -global/       # Kaggles競賽用音檔
-       -dev/
-          -global/       # 測試音檔(不能是空的)
-    transcript/
-       aishell_transcript.txt     #放與音檔對應的文字檔

# 怎麼產生我們要的文字檔

- downloads/
- resource_aishell/
- data_aishell
- wav/
- transcript/
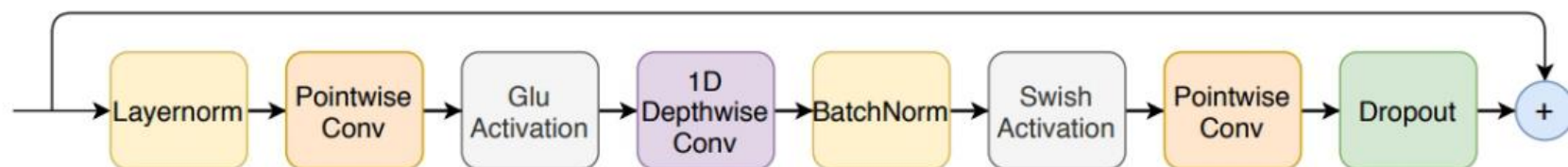- - aishell_transcript.txt

# transcript.txt的擺放方式

- 為了生成文字檔，必須把想要辨認的音檔文字改成 a e i o u 之後加到transcript.txt的末尾

```
3115 3115 hian tsai tok sin e neh
3116 3116 ho thinn khi hong bi bi tshue tioh tshai tshinn tshinn long tshuan koo niu khuai lok kue lit tsi
3117 3117 he leh hoo khah ling tsit e
3118 3118 an ne gua bo kin lai tsong tsit tai a tian si be saih
3119 3119 ah lin tse tso hue honnh
3120 test_1 a e i o u
3121 test_2 a e i o u
3122 test_3 a e i o u
3123 test_4 a e i o u
3124 test_5 a e i o u
3125 test_6 a e i o u
3126 test_7 a e i o u
3127 test_8 a e i o u
3128 test_9 a e i o u
3129 test_10 a e i o u
3130 test_11 a e i o u
3131 test_12 a e i o u
3132 test_13 a e i o u
3133 test_14 a e i o u
3134 test_15 a e i o u
```

# conformer架構

- 引入CNN 和 multi-head self attention的transformer架構

- multi-head self attention: 利用多個尺度的特徵來對長序列進行建模。這種機制可以更好地處理長序列，提高模型的性能和效率。

- separable convolution: 將卷積拆分成兩個部分，pointwise Conv & depthwise Conv ，可以大幅減少計算量
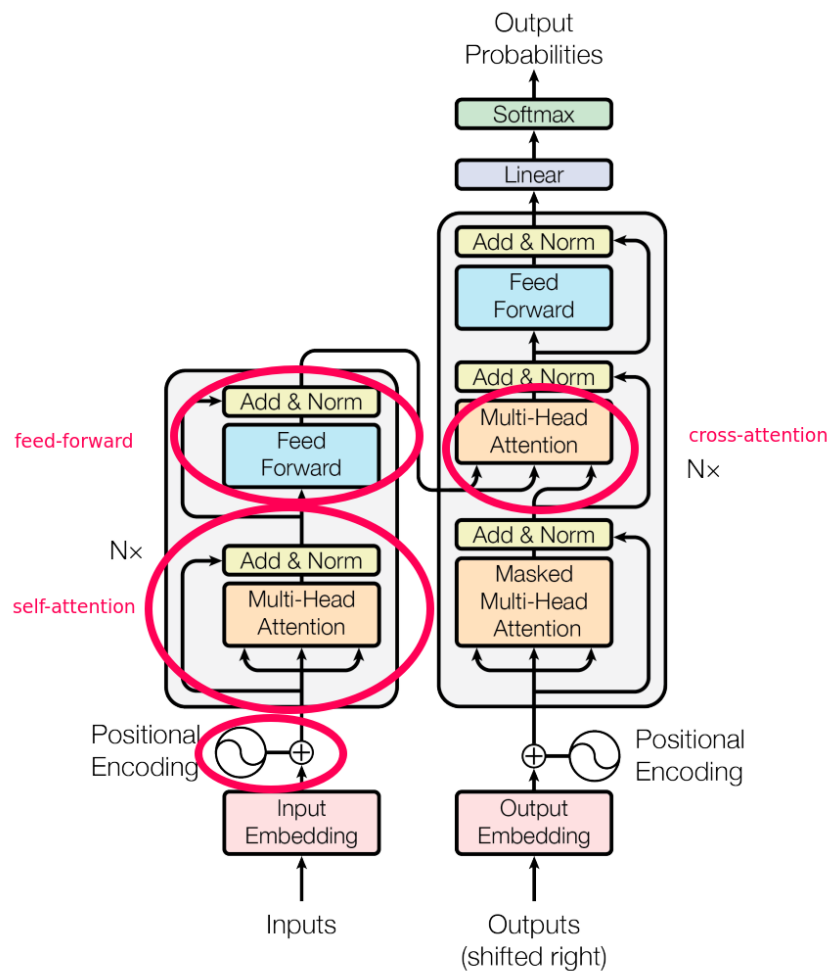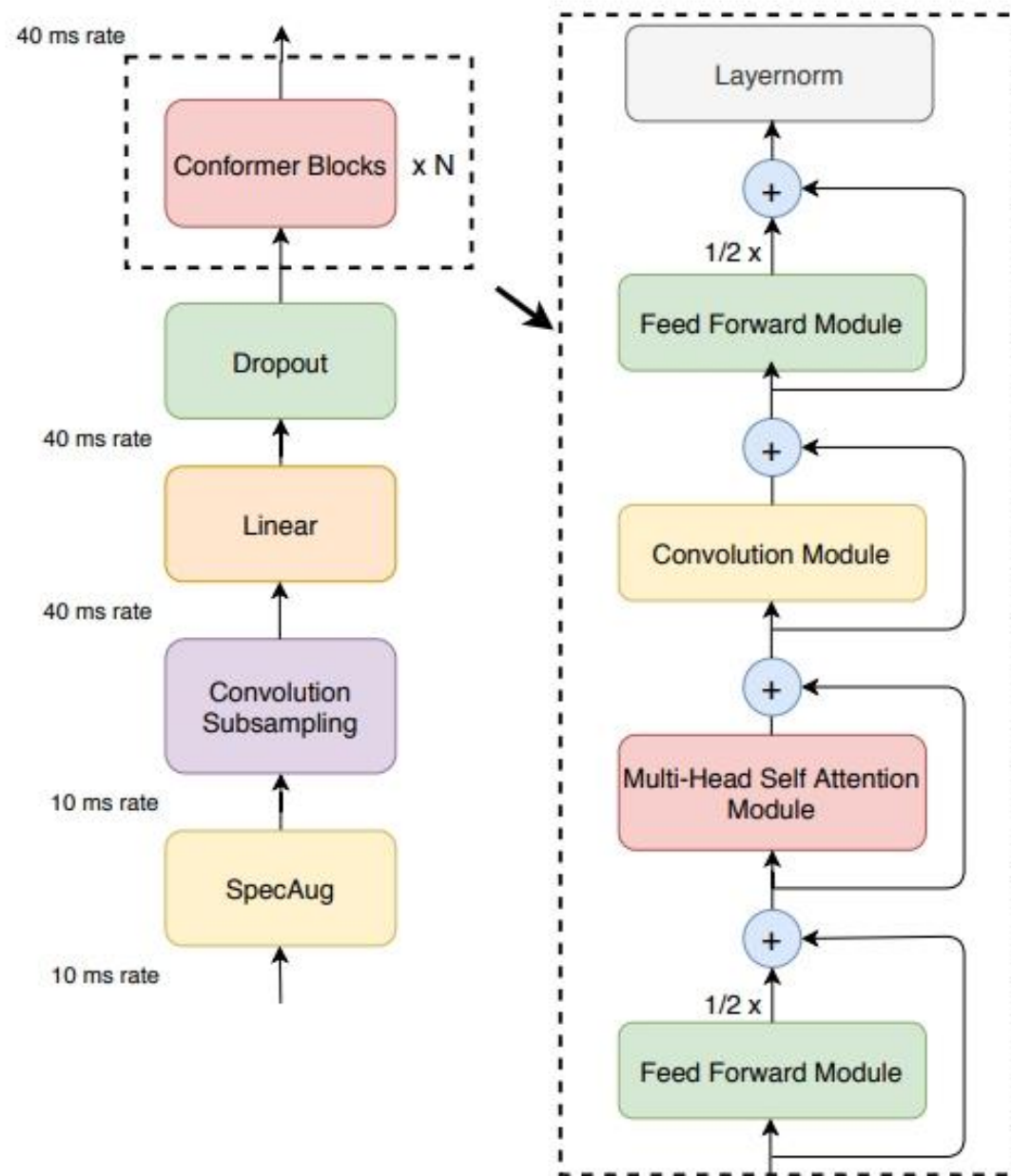
# conformer架構圖



Figure 1: The Transformer - model architecture.

# Depwise Conv & Pointwise Conv

- 參數量: N_std = 4 × 3 × 3 × 3 = 108



3 channel Input          Filters * 4          Maps * 4

知乎 @ashergaga

# Depthwise Conv

- 參數量: N_depthwise = 3 × 3 × 3 = 27



3 channel Input　　　Filters * 3　　　Maps * 3

- depthwise 只改變feature map的大小，不改變通道數。

# Pointwise Conv

- 參數量: N_pointwise = $1 \times 1 \times 3 \times 4 = 12$



Maps * 3          Filters * 4          Maps * 4

知乎 @ashergaga

- pointwise指改變通道數，不改變feature map大小。
- N_separable = N_depthwise + N_pointwise = 39

# Espnet 配置文件

- conformer 模型的配置文件 train.yaml

```
 1 # network architecture
 2 # encoder related
 3 encoder: conformer
 4 encoder_conf:
 5     output_size: 256       # dimension of attention
 6     attention_heads: 4
 7     linear_units: 2048     # the number of units of position-wise feed forward
 8     num_blocks: 12         # the number of encoder blocks
 9     dropout_rate: 0.1
10     positional_dropout_rate: 0.1
11     attention_dropout_rate: 0.1
12     input_layer: conv2d # encoder architecture type
13     normalize_before: true
14     rel_pos_type: latest
15     pos_enc_layer_type: rel_pos
16     selfattention_layer_type: rel_selfattn
17     activation_type: swish
18     macaron_style: true
19     use_cnn_module: true
20     cnn_module_kernel: 32
```

# Conformer的表現

| dataset | Snt | Wrd | Corr | Sub | Del | Ins | Err | S.Err |
|---|---|---|---|---|---|---|---|---|
| beam10_ctc0.4/dev | 14326 | 205341 | 95.8 | 4.1 | 0.1 | 0.1 | 4.3 | 33.1 |
| beam10_ctc0.4/test | 7176 | 104765 | 95.4 | 4.4 | 0.1 | 0.1 | 4.6 | 34.7 |

CER

# Branchformer 架構

- Branchformer是將global和local特徵並行提取，相互之間是獨立的，兩種特徵最後可以進行結合

- branching mechanism: 將輸入序列分解成多個子序列，並且為每個子序列分配一個獨立的分支，進行獨立的自注意力(self-attention )計算。

- Branchformer 擅長處理長文本序列(因為將輸入序列分解成多個子序列)

# Branchformer 架構圖



(a) Overall architecture of the encoder. A stack of identical Branchformer blocks are used to capture local and global dependencies.

(b) Architecture of our Branchformer encoder block. It consists of two parallel branches. One branch employs attention to capture global context, while the other branch utilizes the MLP with convolutional gating to extract local context.

# Espnet 配置文件

```
 1 # network architecture
 2 # encoder related
 3 encoder: branchformer
 4 encoder_conf:
 5     output_size: 256
 6     use_attn: true
 7     attention_heads: 4
 8     attention_layer_type: rel_selfattn
 9     pos_enc_layer_type: rel_pos
10     rel_pos_type: latest
11     use_cgmlp: true
12     cgmlp_linear_units: 2048
13     cgmlp_conv_kernel: 31
14     use_linear_after_conv: false
15     gate_activation: identity
16     merge_method: concat
17     cgmlp_weight: 0.5                # used only if merge_method is "fixed_ave"
18     attn_branch_drop_rate: 0.0      # used only if merge_method is "learned_ave"
19     num_blocks: 24
20     dropout_rate: 0.1
21     positional_dropout_rate: 0.1
22     attention_dropout_rate: 0.1
23     input_layer: conv2d
24     stochastic_depth_rate: 0.0
```

# Branchformer的表現(官方網站)

CER

| dataset | Snt | Wrd | Corr | Sub | Del | Ins | Err | S.Err |
|---|---|---|---|---|---|---|---|---|
| beam10_ctc0.4/dev | 14326 | 205341 | 96.0 | 4.0 | 0.1 | 0.1 | 4.1 | 32.7 |
| beam10_ctc0.4/test | 7176 | 104765 | 95.7 | 4.2 | 0.1 | 0.1 | 4.4 | 34.1 |

# 測試結果(自己的)

- 在解碼文件輸出目錄中，會生成一個result.txt 文件，該文件保存了最終在測試集和驗證集上的解碼結果，如下圖所示：

```
                  SYSTEM SUMMARY PERCENTAGES by SPEAKER

|exp/asr_train_asr_branchformer_raw_zh_char_sp/decode_asr_branchformer_asr_model_valid.acc.ave/dev/score_wer/hyp.trn|
|------------------------------------------------------------------------------------------------------------------|
| SPKR      |  # Snt    # Wrd  |   Corr      Sub       Del       Ins       Err        S.Err |
|-----------+------------------+------------------------------------------------------------|
| global    |   20        311  |   94.9      4.2       1.0       1.9       7.1         40.0 |
|===========+==================+============================================================|
| Sum/Avg   |   20        311  |   94.9      4.2       1.0       1.9       7.1         40.0 |
|===========+==================+============================================================|
|  Mean     |  20.0      311.0 |   94.9      4.2       1.0       1.9       7.1         40.0 |
|  S.D.     |   0.0        0.0 |    0.0      0.0       0.0       0.0       0.0          0.0 |
|  Median   |  20.0      311.0 |   94.9      4.2       1.0       1.9       7.1         40.0 |
|------------------------------------------------------------------------------------------------------------------|
```
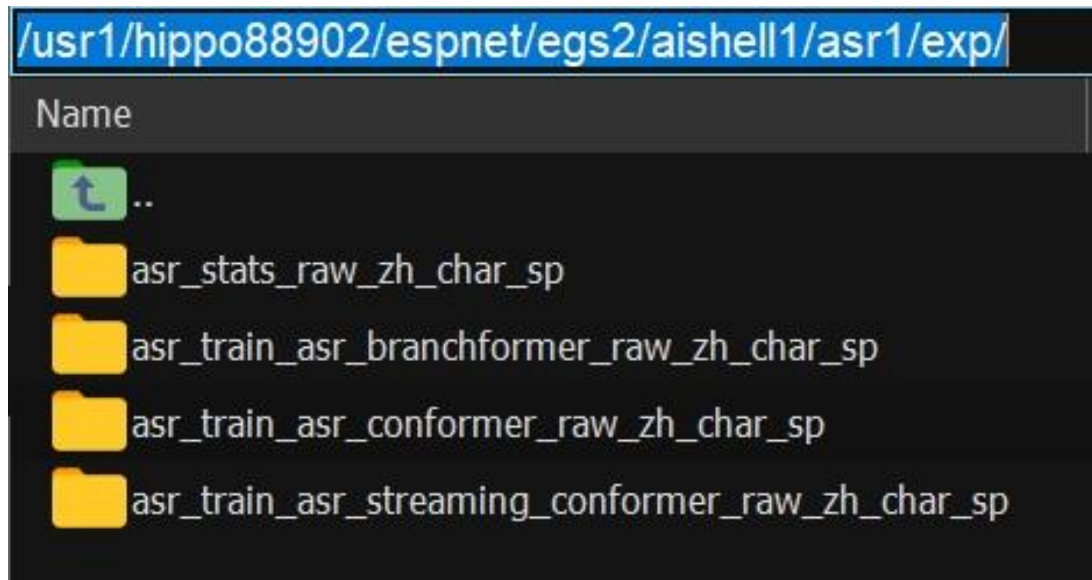
# 怎麼更換model

- 更改/[username]/espnet/egs2/aishell1/asr1/run.sh的文件

```
1 #!/usr/bin/env bash
2 # Set bash to 'debug' mode, it will exit on :
3 # -e 'error', -u 'undefined variable', -o ... 'error in pipeline', -x 'print commands',
4 set -e
5 set -u
6 set -o pipefail
7
8 train_set=train
9 valid_set=dev
10 test_sets="dev test"
11
12 asr_config=conf/train_asr_streaming_conformer.yaml
13 inference_config=conf/decode_asr_branchformer.yaml
```

# 不同model出來的結果放在哪

- /usr1/[username]/espnet/egs2/aishell1/asr1/exp/



- /decode_asr_branchformer_asr_model_valid.acc.ave/test/text 就是辨認後的文字檔