

HW 2 - GAN-based Mel-Vocoder

Name: 譚至斌

ID: d12942015

Novelty

- No novelty
 - I just reproduced Hifi-GAN in the homework dataset.....

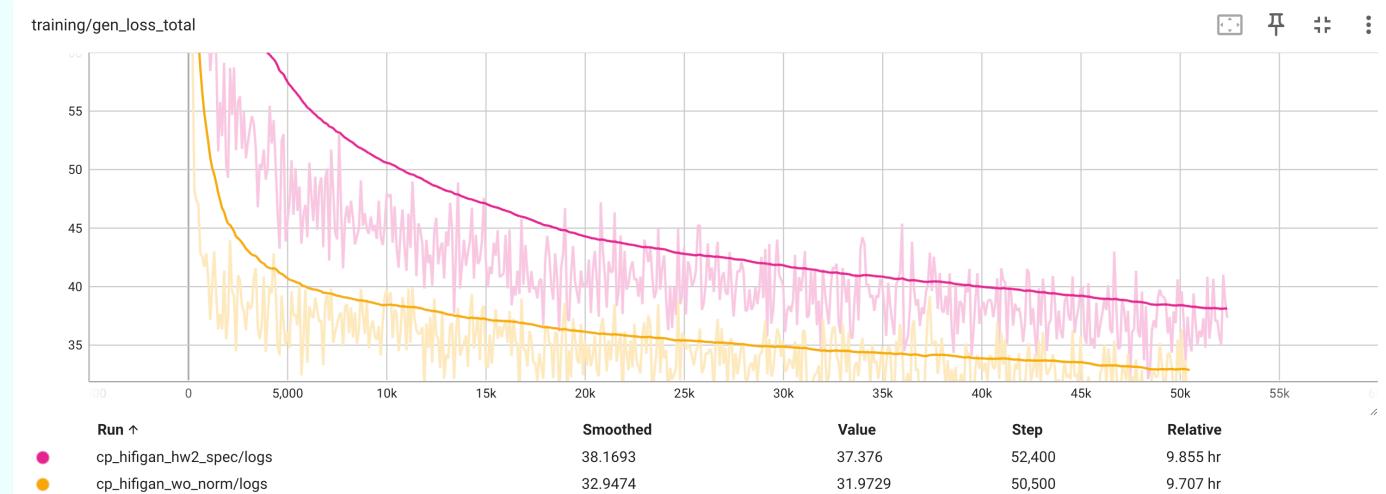
Methodology

- The methodology can be described by 4 steps:
 1. Resample all the audio to 22050
 2. Replace Hifi-GAN's `melspectrogram` with hw2's `melspectrogram`
 3. Train Hifi-GAN
 4. Inference

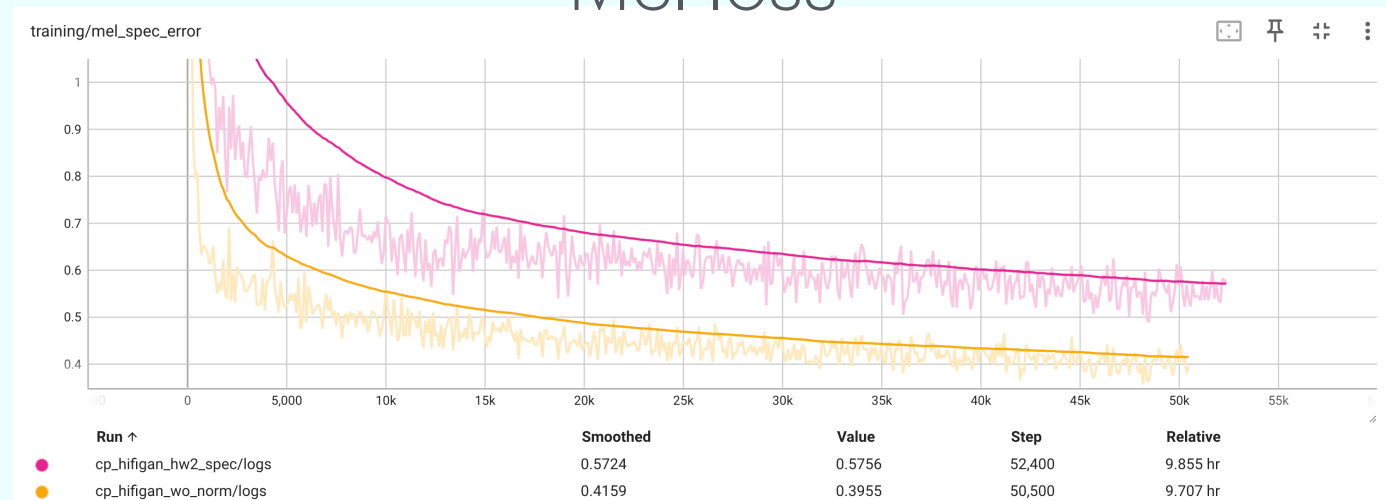
Result

- About 50k steps, the total training loss decreased to **~36** and the mel loss decreased to **~0.58**
- The validation loss of mel decreased to **~0.76**
- Note that we consider the pink line which represents the loss computed by hw2's `melspectrogram`

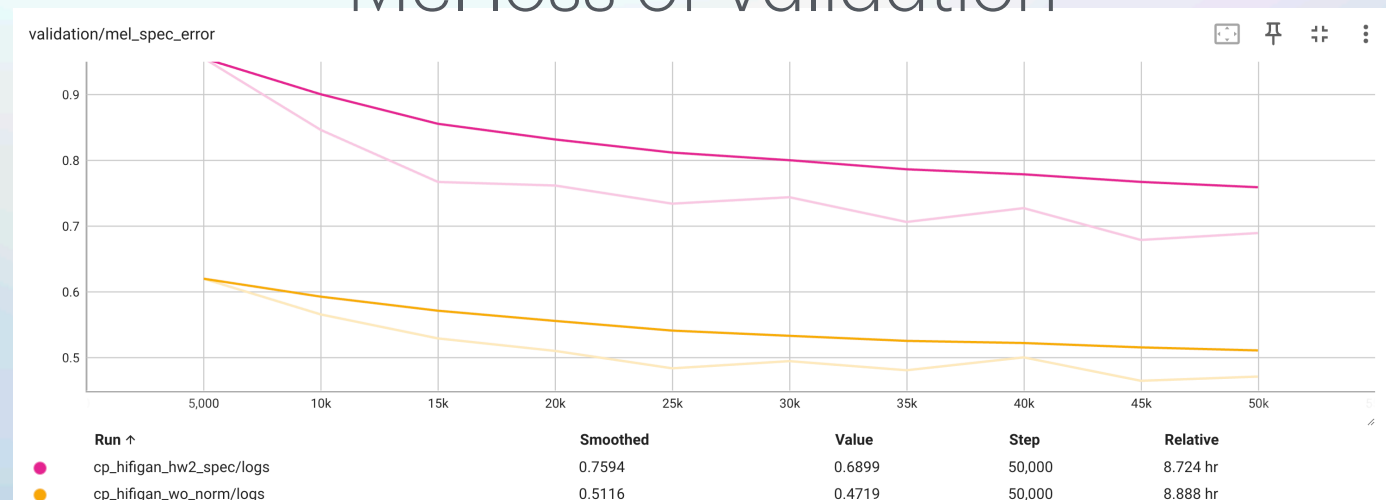
Discriminative loss + Mel loss



Mel loss



Mel loss of validation



Result

- Evaluation on the validation set:
 - Griffin & Lim algorithm gets best scores

Griffin & Lim

Hifi-GAN (hw2)

Hifi-GAN (original)

```
{
  "M-STFT": 0.3219452682323159,
  "log2f0_mean": 0.061804500390175085,
  "FAD": 0.02174284536276261
}
{
  "M-STFT": 1.1204531639814377,
  "log2f0_mean": 0.17611471862812564,
  "FAD": 0.32018343540095096
}
{
  "M-STFT": 9.274185792325785,
  "log2f0_mean": 6.465007469306043,
  "FAD": 8.585338238296353
}
```

Findings

- Griffin&Lim algorithm **outperform** on the validation dataset
- It's so surprising to me while hearing the results generated by this classical algorithm and find them nearly same as the source audios
- I find the audio generated by Hifi-GAN containing obvious artifacts which sounds like **trembling**
- I guess the reason is Hifi-GAN generate audio by frames and cannot keep the consistence among a group of frames

Findings

- Some interesting findings about implementation
 - Even if all the parameters set to the same values, ``Melspectrogram`` in ``torchaudio`` still makes **different** outputs from the hand-craft ``melspectrogram`` implemented by Hifi-GAN's authors
 - ``melspectrogram`` implemented by Hifi-GAN's authors always get **smaller** loss values in the training process

Discussion

- This homework is a good experience for me to learning lots of audio processing coding skills like:
 - How to deal with spectrogram transforming
 - How to train a GAN
- Since some mistakes on coding, I have to give up an **~100k-steps** checkpoint
- However, the generated audios still **worse than those from Griffin&Lim algorithm**

Discussion

- I choose the results generated by Hifi-GAN which is trained by myself rather than Griffin&Lim algorithm
- It makes me feel more sense of accomplishment.....

Appendix

- source code: <https://github.com/tanchihpin0517/NTU/tree/master/deepmir/hw2>
- model: https://www.dropbox.com/scl/fi/b0kdwrxfq91tp5jm8b6j/g_00050000?rlkey=20bved2g0ijev9zx0yycnlv3u&dl=0