

○ ○ ○ ○

DUONG HOANG
LUC TAN THO

AUDIO STYLE TRANSFER

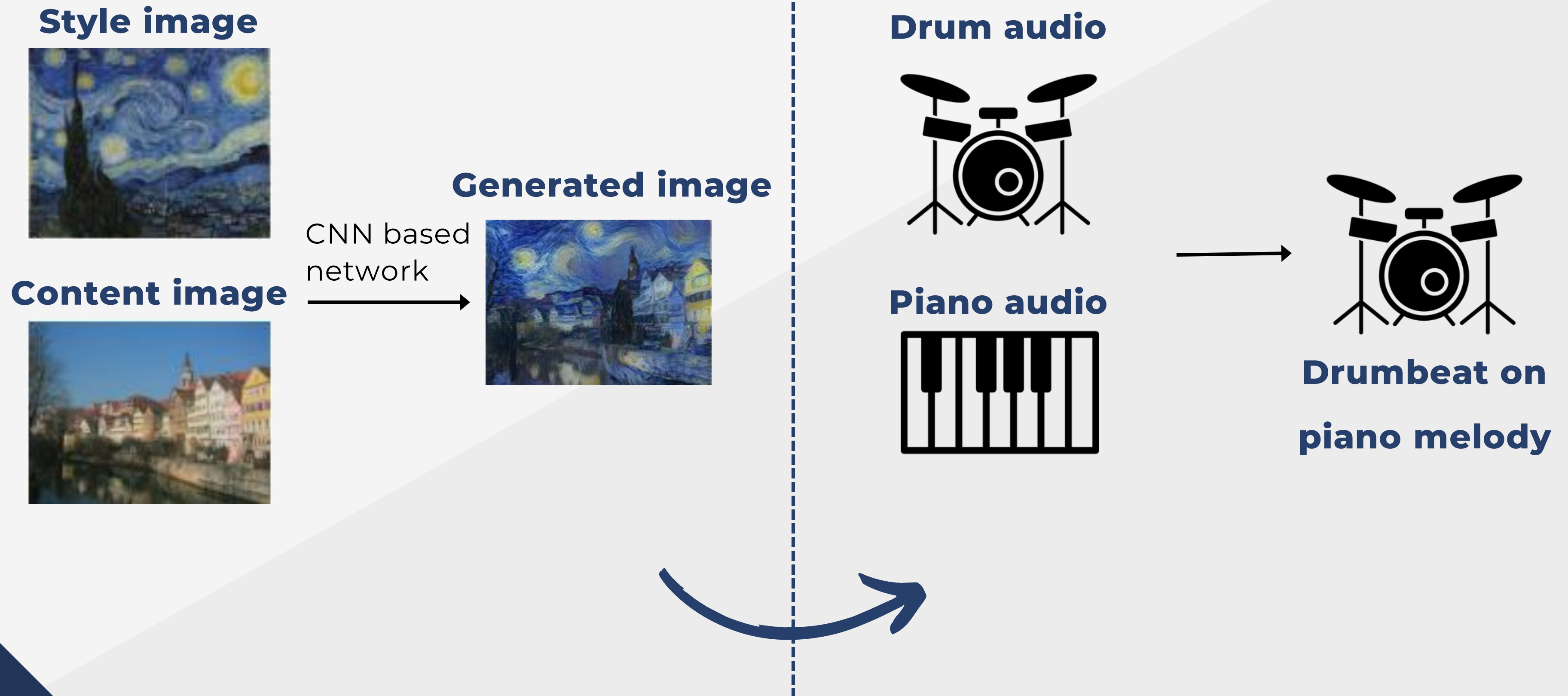
○ ○ ○ ○

TABLE OF CONTENTS

- From image to audio style
- Audio transfer using Vggish (Vgg for audio classification)
- Audio transfer using CNN14
- Shallow network for audio transfer

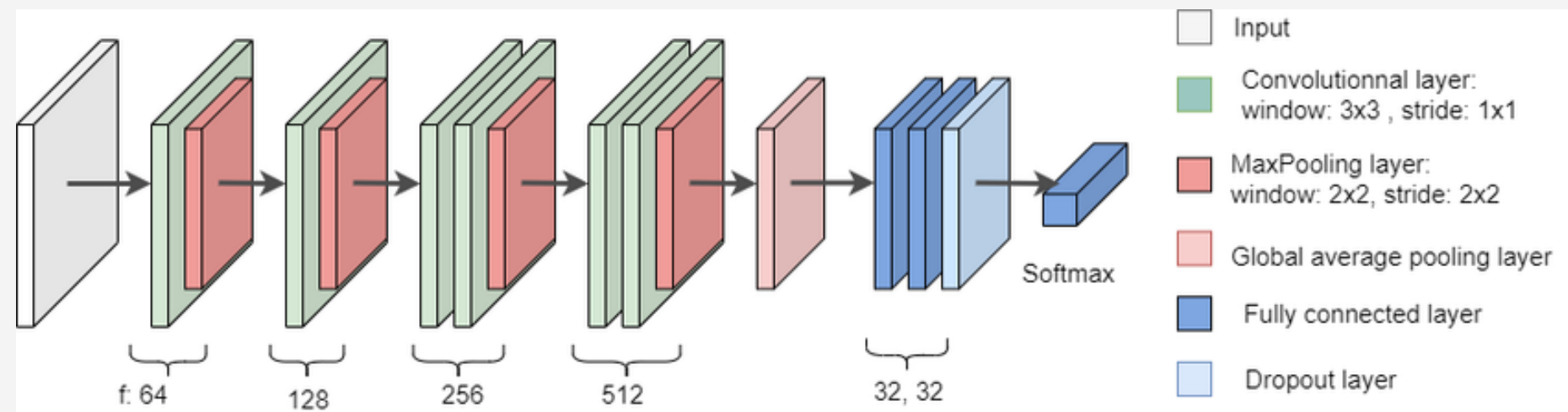


FROM IMAGE TO AUDIO STYLE



VGGISH FOR AUDIO TRANSFER

Vggish network



Same workflow as the original problem on image:

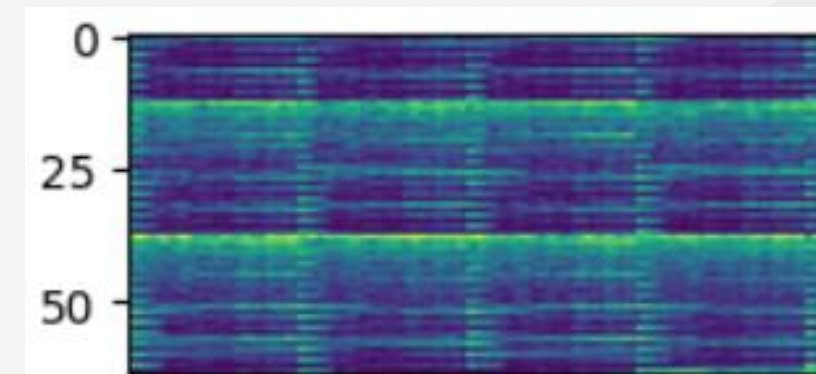
- Use a pre-trained CNN network for classification
- Choose layers for computing content and style loss
- Do backprop to generate stylized audio

Vggish:

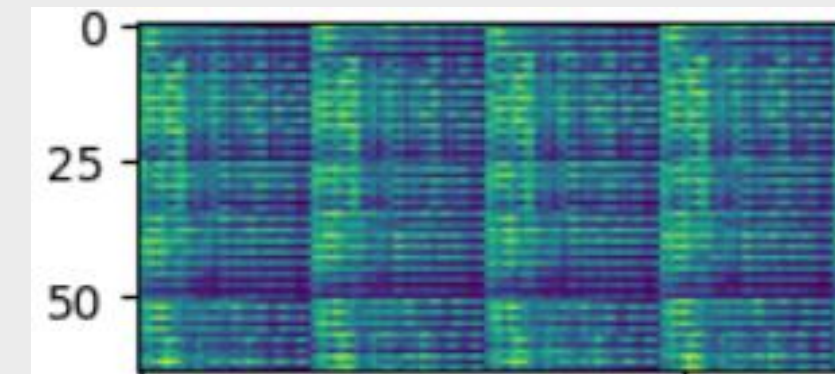
- Similar to Vgg
- Pre-trained model for audio classification
- Input: 2D mel spectrogram

Results

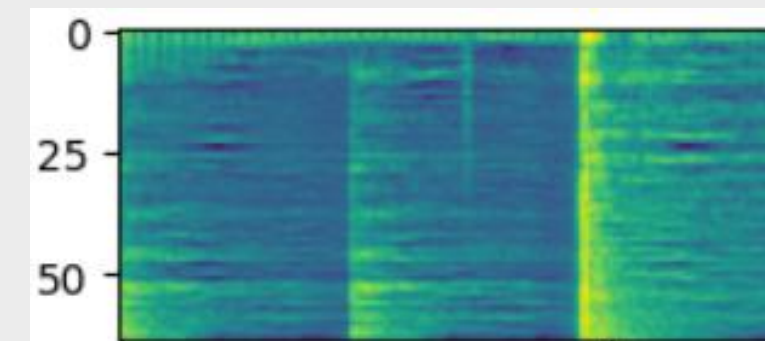
Content



Style



Generated audio



Not good, lost all content!

CNN14 FOR AUDIO TRANSFER

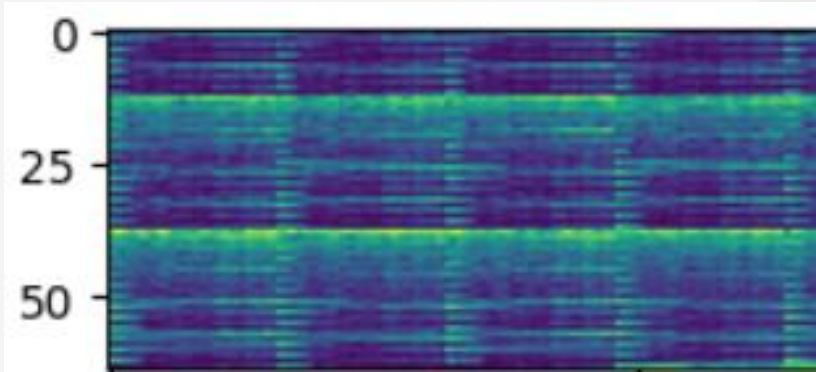
CNN14

| |
|--|
| $\left(\begin{matrix} 3 \times 3 @ 64 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| $\left(\begin{matrix} 3 \times 3 @ 128 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| $\left(\begin{matrix} 3 \times 3 @ 256 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| $\left(\begin{matrix} 3 \times 3 @ 512 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| Pooling 2×2 |
| $\left(\begin{matrix} 3 \times 3 @ 1024 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| Pooling 2×2 |
| $\left(\begin{matrix} 3 \times 3 @ 2048 \\ \text{BN, ReLU} \end{matrix} \right) \times 2$ |
| Global pooling |
| FC 2048, ReLU |
| FC 527, Sigmoid |

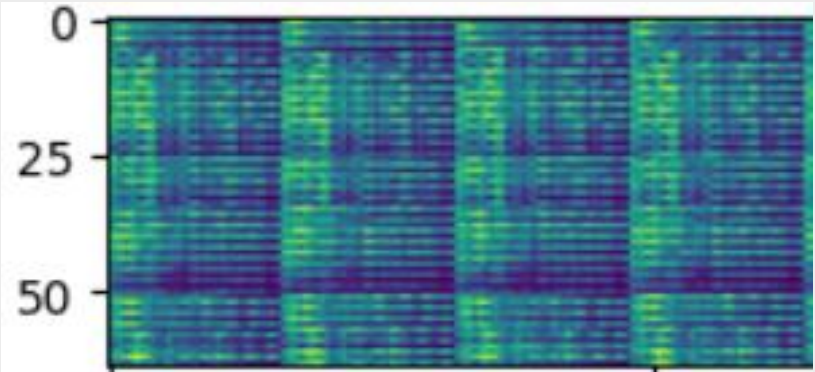
- CNN14:
- Pre-trained model for audio classification
 - Input: 2D log-mel spectrogram

Results

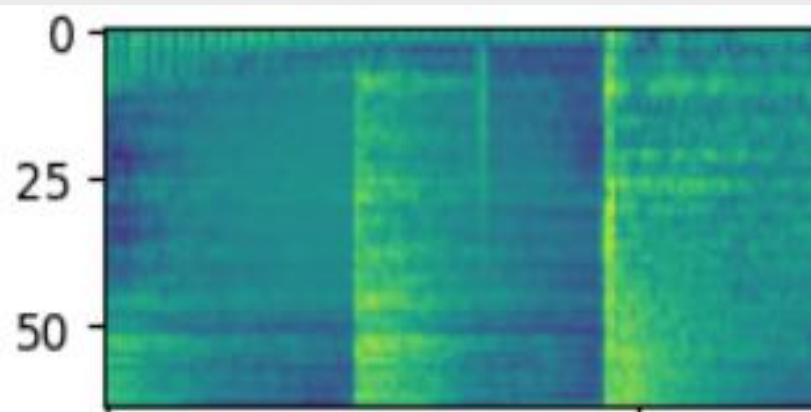
Content



Style



Generated audio



Still not good!

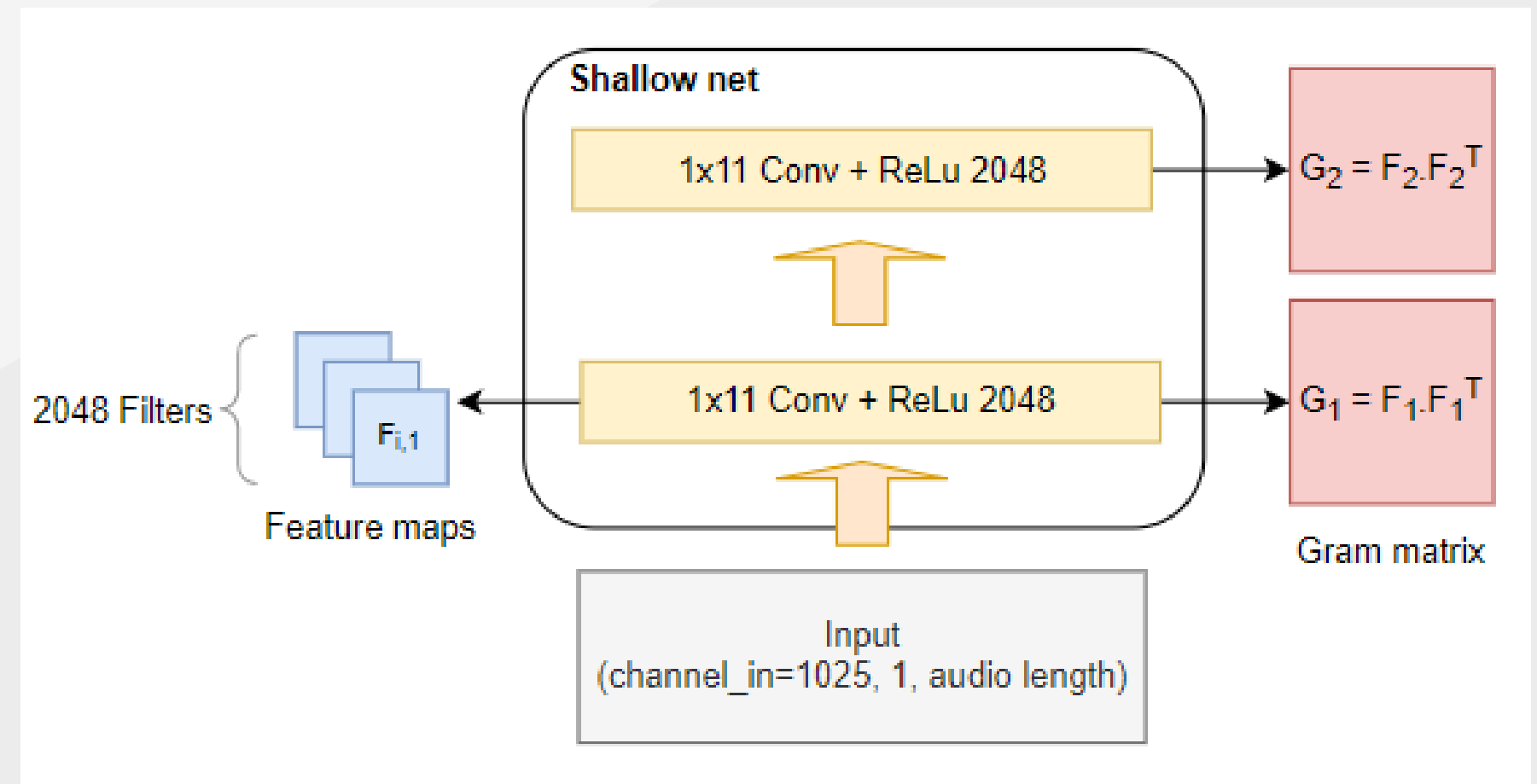
SHALLOW NET

Note:

- Visual perception and auditive perception may be different
- Pooling layer may not work for this audio problem
- Each frequency should not be interlinked with others

Shallow net

Based on Ulyanov's work, we define a shallow model:



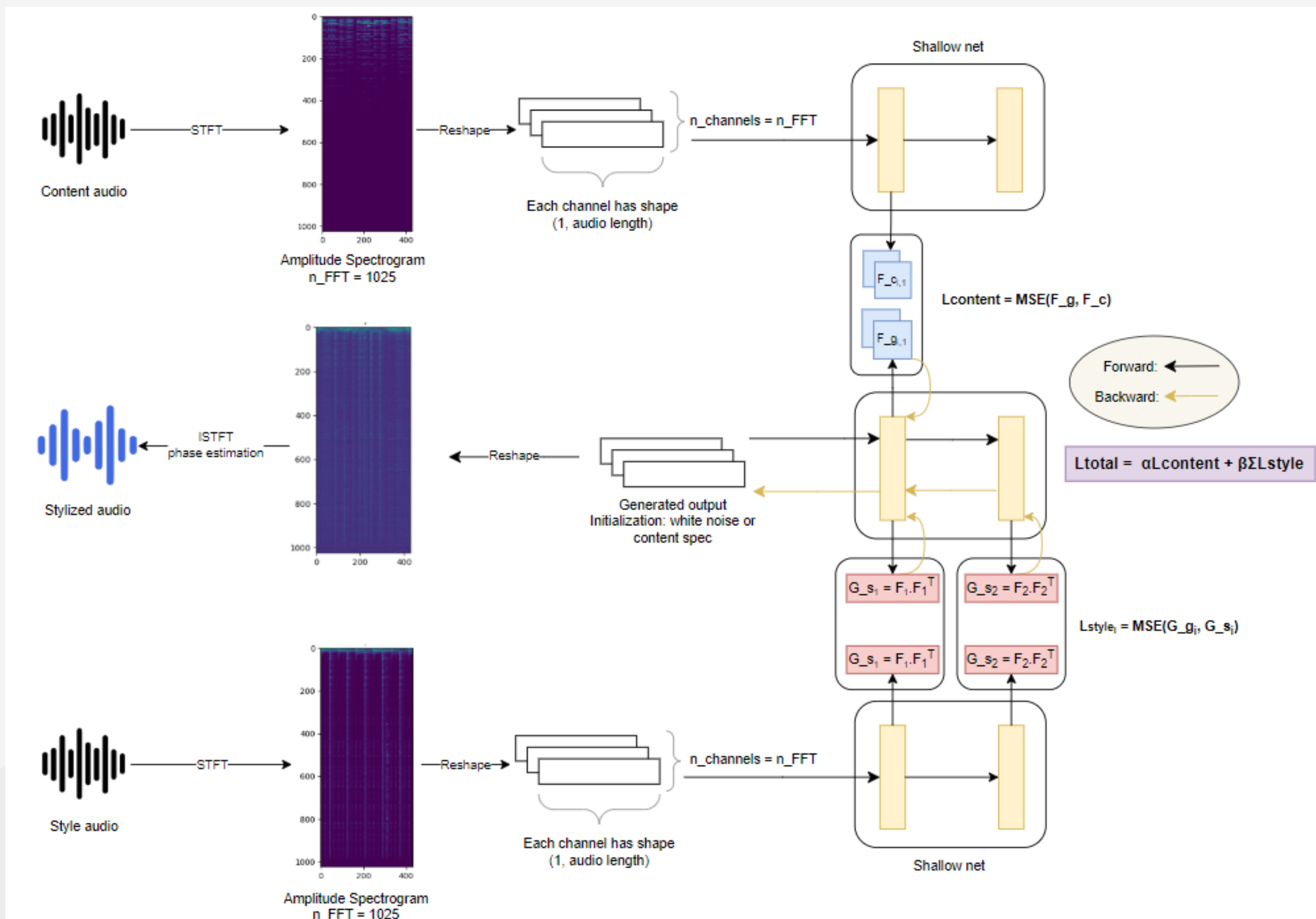
1D conv instead of 2D conv

2 conv layers, 2048 filters for each layer

Do not need to be trained, all weights are random

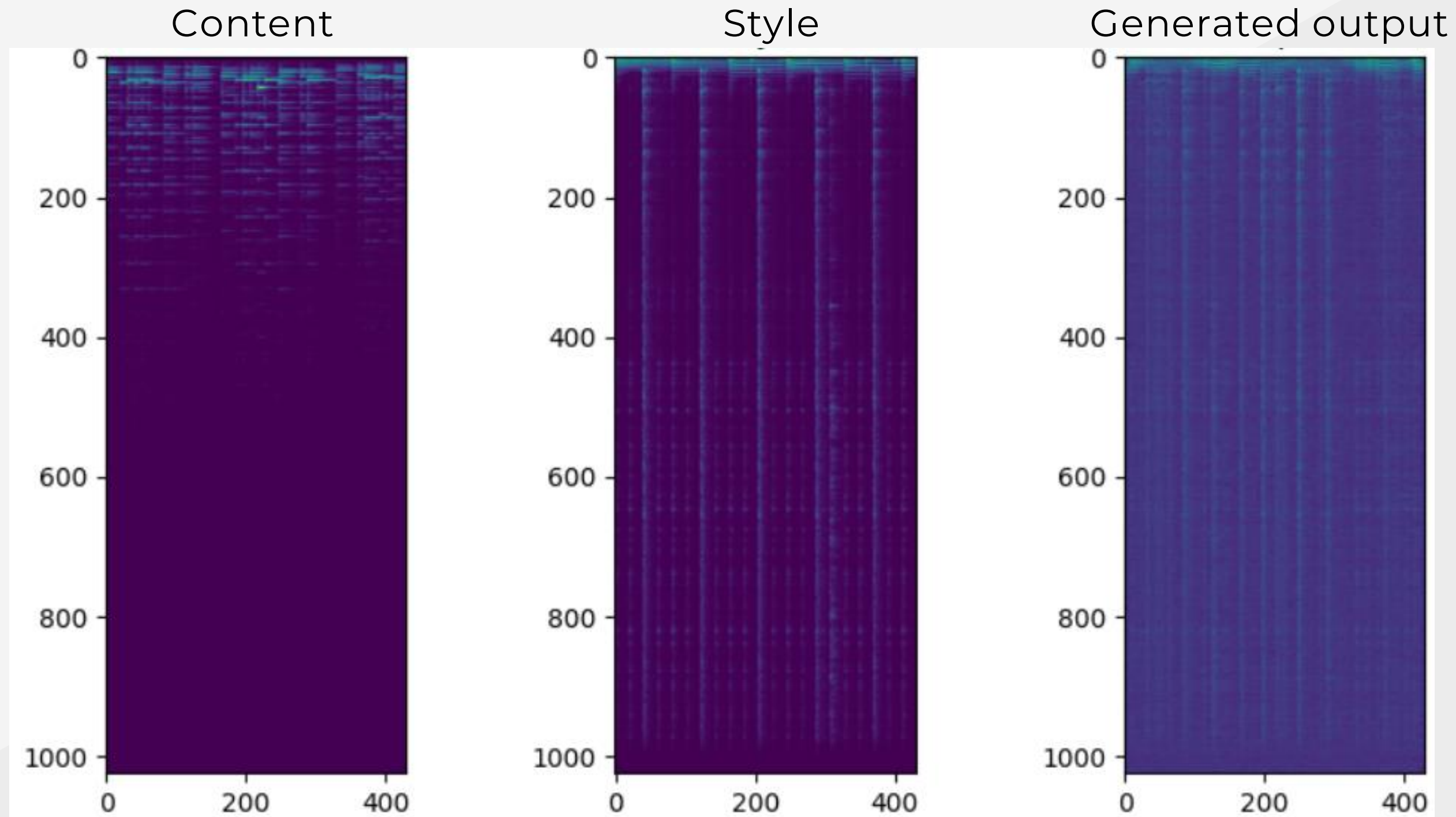
Input: reshaped amplitude spectrogram

SHALLOW NET















SHALLOW NET

Result



SHALLOW NET

More results:

| | Content | Style | Generated output |
|-----------------------------------|---|---|---|
| Piano to drum (same song) |  |  |  |
| Drum to piano (same song) |  |  |  |
| Drum to piano (different song) |  |  |  |
| Piano to String (same song) |  |  |  |

[Click here to access our github repository](#)

REFERENCES

- Image Style Transfer Using Convolutional Neural Networks
- CNN ARCHITECTURES FOR LARGE-SCALE AUDIO CLASSIFICATION
- <https://dmitryulyanov.github.io/>
- AUDIO STYLE TRANSFER, Eric Grinstein, Ngoc Q.K. Duong, Alexey Ozerov and Patrick Pérez
- PANNs: Large-Scale Pretrained Audio Neural Networks for Audio Pattern Recognition



**THANK YOU/
MERCI
BEAUCOUP**

