



StyleWav: Guiding Image Synthesis Using Audio

Phuc (Jerry) Ngo¹, Swami Sankaranarayanan², Phillip Isola³

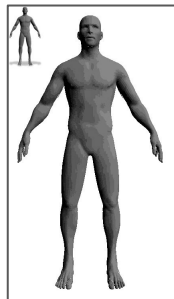
¹Department of Computer Science and Maths, Beloit College

^{2, 3}Department of Electrical Engineering & Computer Science, MIT

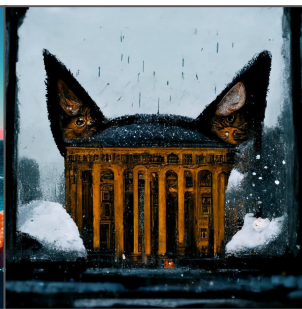
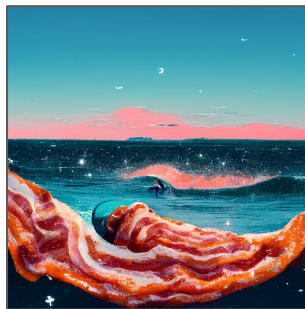


CLIP-Guiding Generative Algorithm

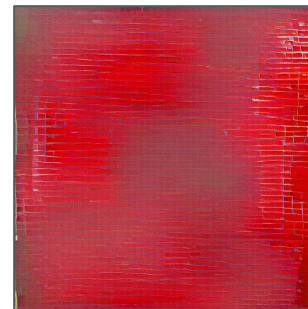
- Recent works in multimodality systems have enabled cross modality generation



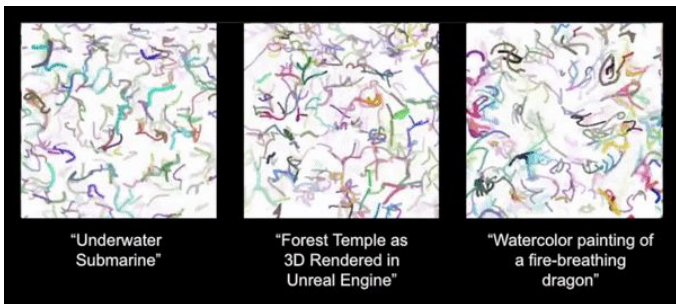
Text2Mesh, Michel et al, 2021



CLIP Guided Diffusion, MidJourney



CLIP + VQGAN, nerdyrodent



"Underwater Submarine"

"Forest Temple as 3D Rendered in Unreal Engine"

"Watercolor painting of a fire-breathing dragon"

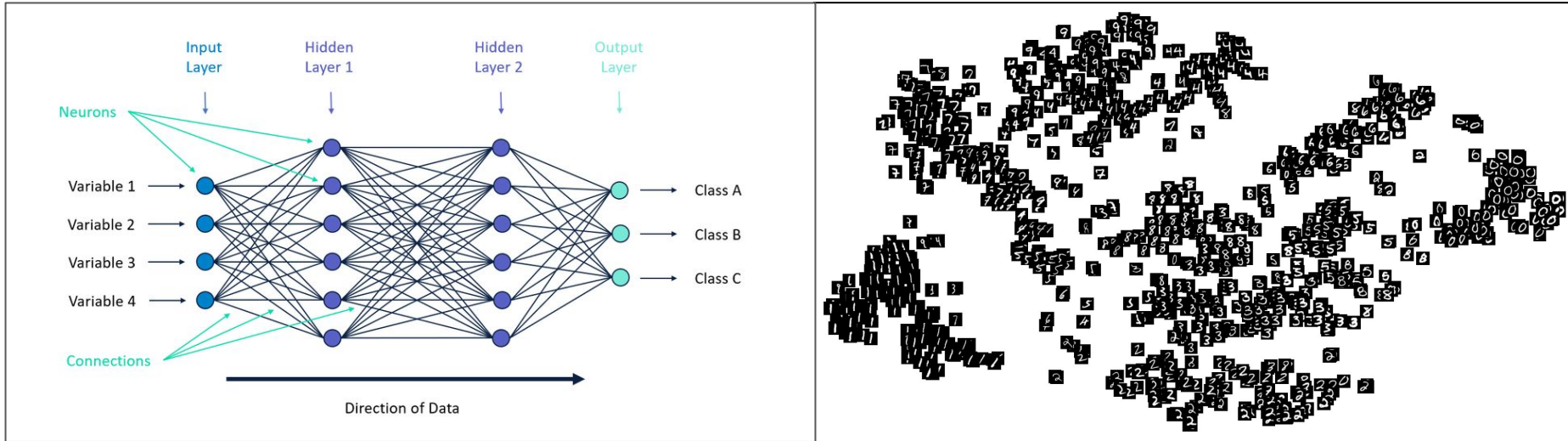
CLIPDraw, Frans et al



Pixray, dribnet

Neural Network and Representation

Representation : The second to last feature vector of neural network

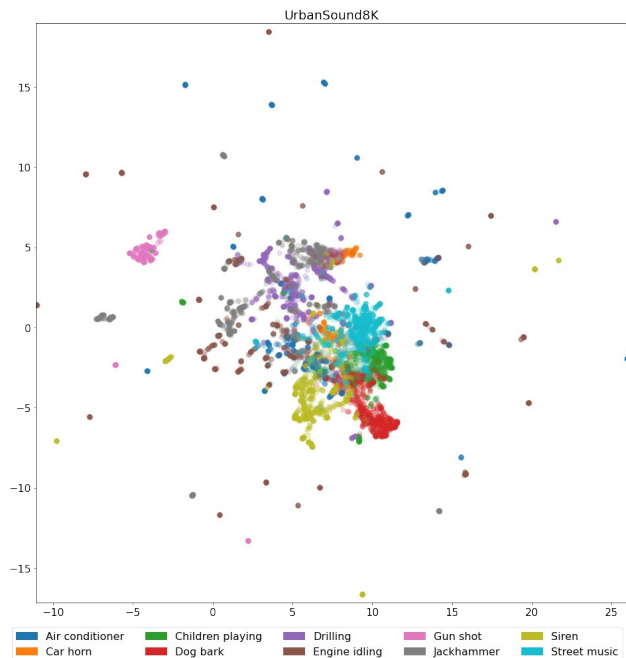


- Representation of the same object is clustered together
- The purpose of contrastive training is to further encourage representations of the same class to be near each other and push other classes' representations far away

Latent Space Visualization, [HackerNoon](#)
It's a No Brainer: An Introduction to Neural Networks, [Alteryx.com](#)

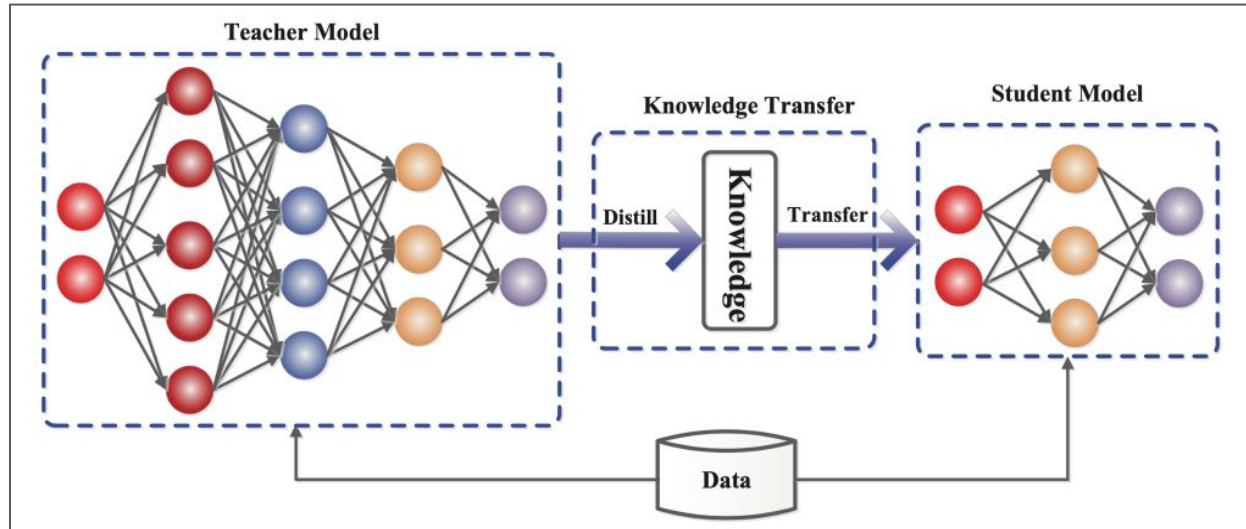
WAV2CLIP: Extending CLIP to Music

- This is a distilled model from CLIP
- It embed audio to the joint representation space of CLIP
- Could perform multiple downstream tasks

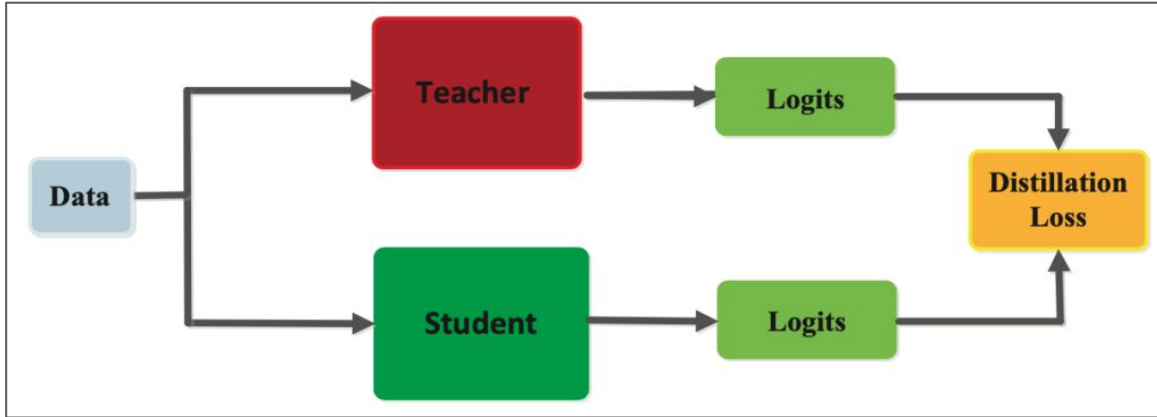


Knowledge Distillation

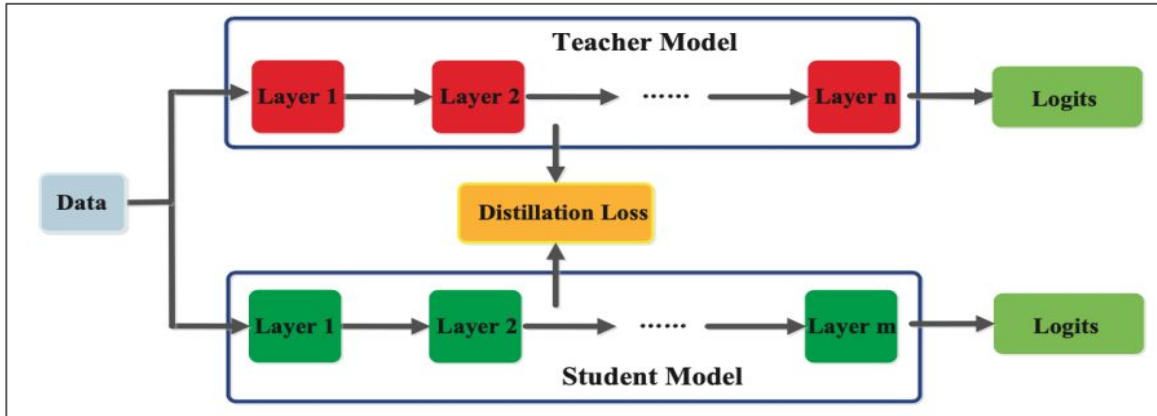
- Used to transfer knowledge from a bigger model to a smaller model
- Teach model to learn compressed representation without the loss of validity
- Tackle limited memory and computational capacity issue when deploying big deep learning model



Distillation Type



Response based



Feature based

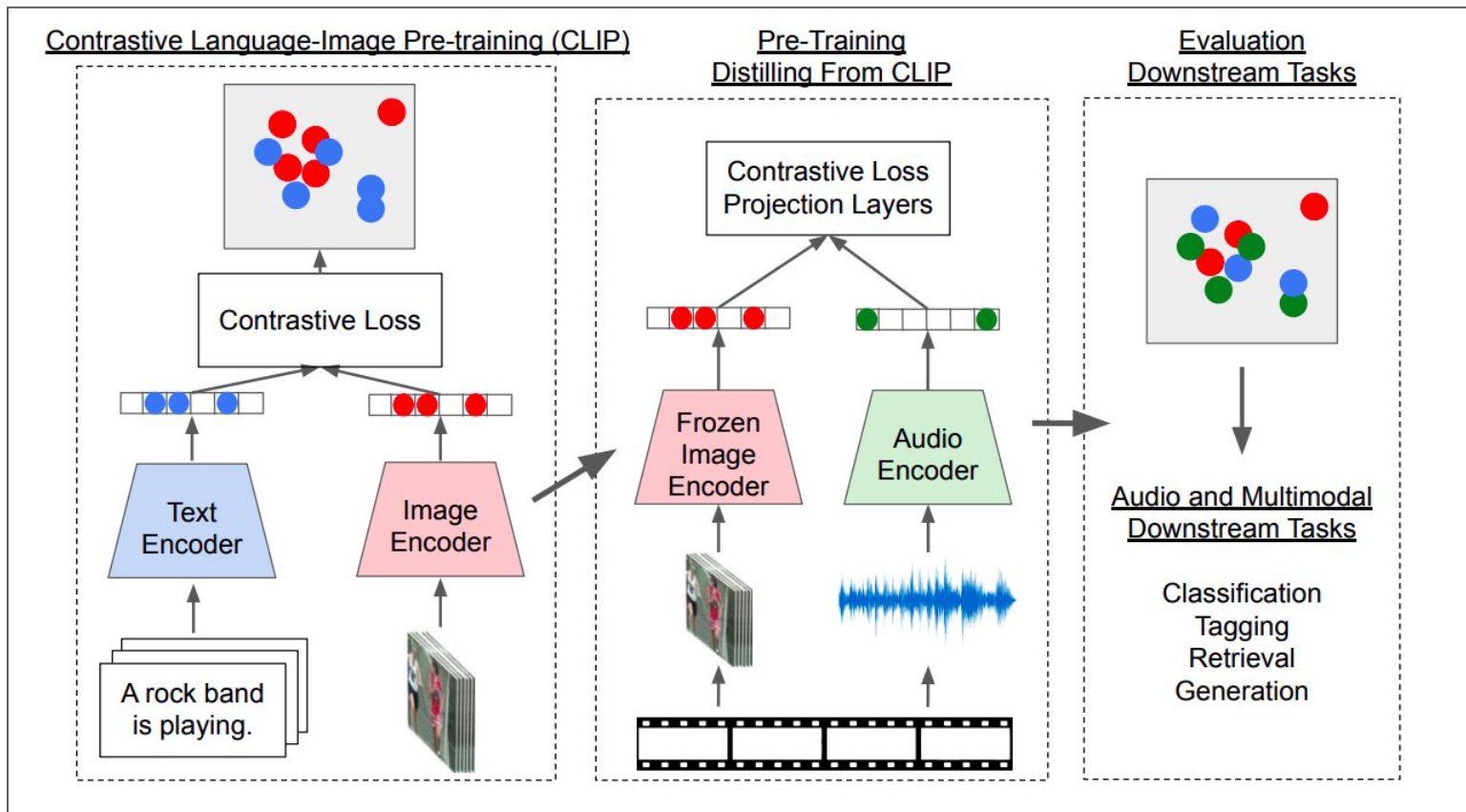
VGGSound Dataset

- For multimodal learning, we will need a dataset includes different modalities
- VGGSound is audio-visual Youtube video dataset. It includes:
 - ~200k 10-second clips
 - 309 class
- For distillation, they sample:
 - 5-second videos that has 150 frames each
 - Get CLIP embeddings for each frames

VGGSound, robots.ox.ac.uk



WAV2CLIP Distillation Process



$$CXLoss = L(f(\text{Image}), \text{Audio}) + L(\text{Image}, g(\text{Audio}))$$

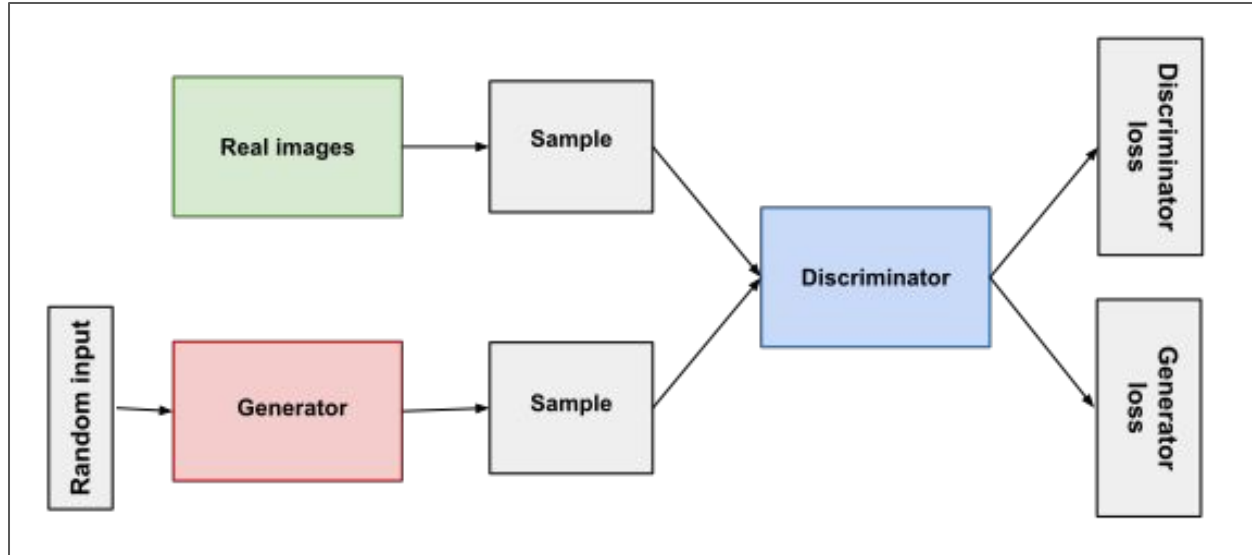
Sound to Image Generation

Text/Audio to Image Generation with VQGAN-CLIP



Question: Could we make use of the sound representation to generate more meaningful images?

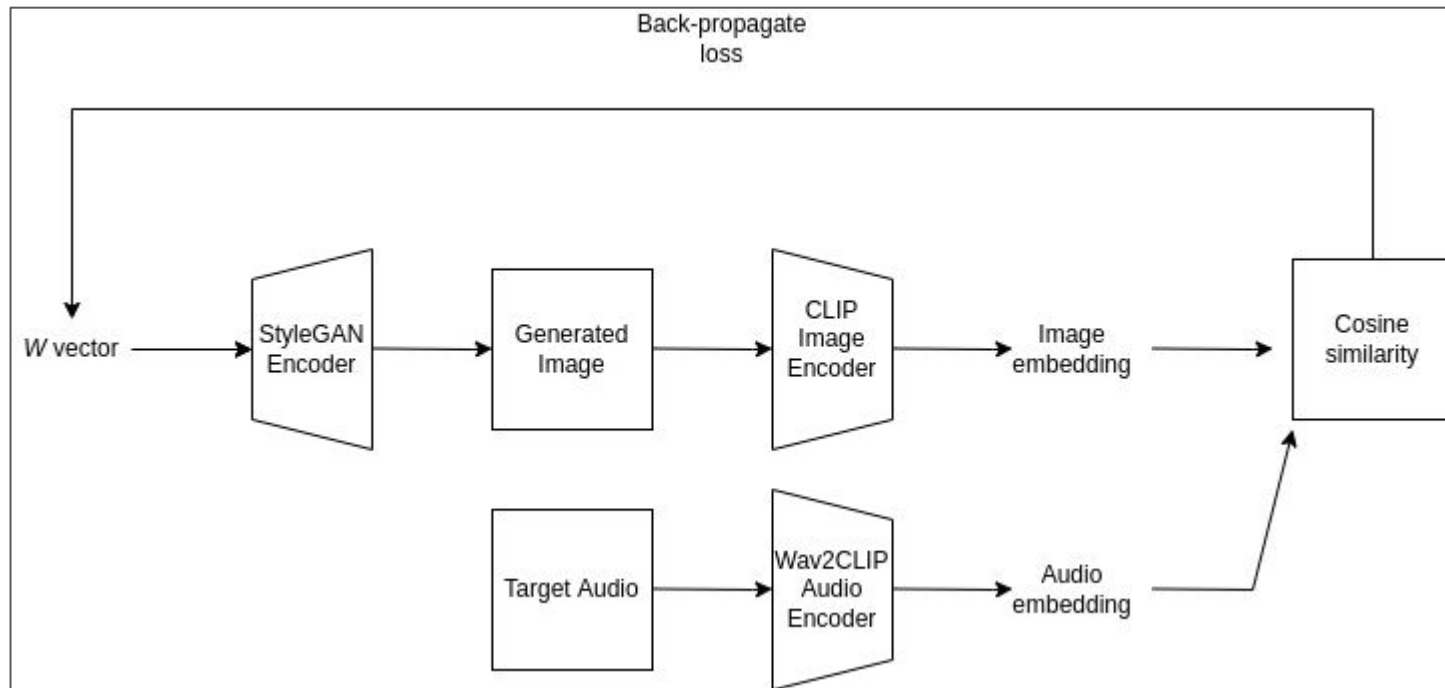
GAN Structure



StyleGAN Result



StyleWAV Structure



Loss Function:

$$\arg \min_{w \in W} D_C(G(w), t) + \lambda_{L2} \|w - w_s\| + \lambda_{ID} \mathcal{L}_{ID}(w)$$

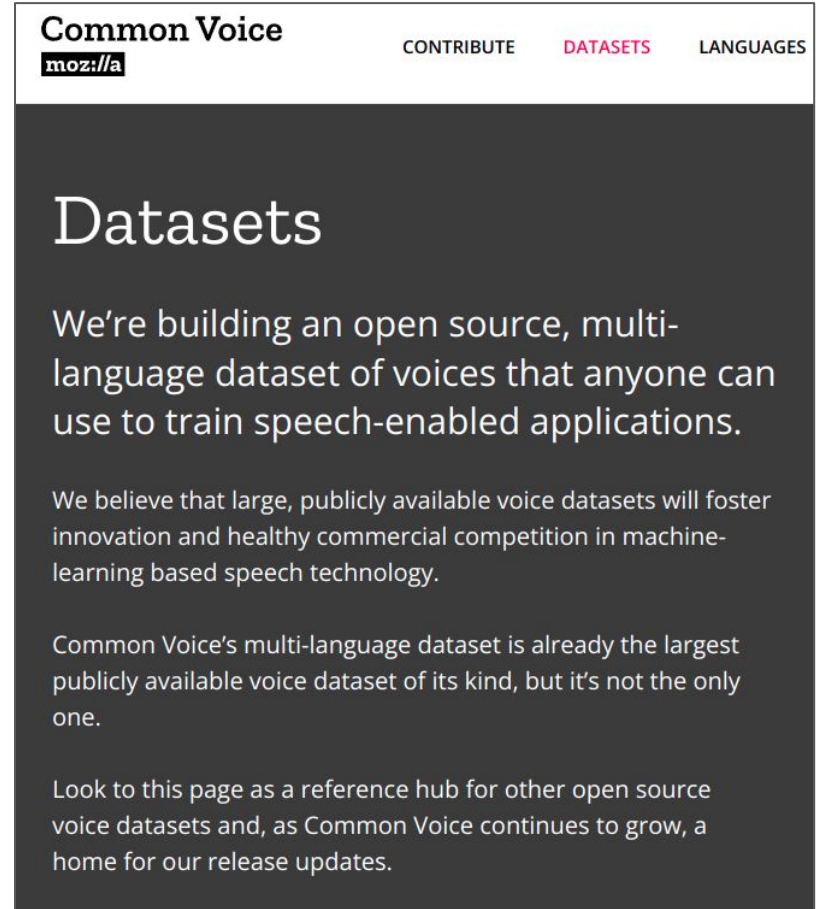
Where G is pretrained StyleGAN, D_C is the cosine distance. T is the audio representation. The two lambda is equal to zero when we do free generation.

$$\mathcal{L}_{ID} = D_C(R(G(w_s)), R(G(W)))$$

R is pretrained ArcFace model for face recognition

Image Generation

- 2 sex: male and female
- Each sex has 100 5-second audios from Mozilla Common Voice Dataset

A screenshot of the Mozilla Common Voice website's 'Datasets' page. The page has a dark grey background with white text. At the top, there is a navigation bar with the 'Common Voice' logo and the 'moz://a' icon on the left, and 'CONTRIBUTE', 'DATASETS', and 'LANGUAGES' links on the right. The main heading is 'Datasets' in a large, white, sans-serif font. Below the heading, there are three paragraphs of text in a smaller white font. The first paragraph states: 'We're building an open source, multi-language dataset of voices that anyone can use to train speech-enabled applications.' The second paragraph states: 'We believe that large, publicly available voice datasets will foster innovation and healthy commercial competition in machine-learning based speech technology.' The third paragraph states: 'Common Voice's multi-language dataset is already the largest publicly available voice dataset of its kind, but it's not the only one.' The final paragraph states: 'Look to this page as a reference hub for other open source voice datasets and, as Common Voice continues to grow, a home for our release updates.'

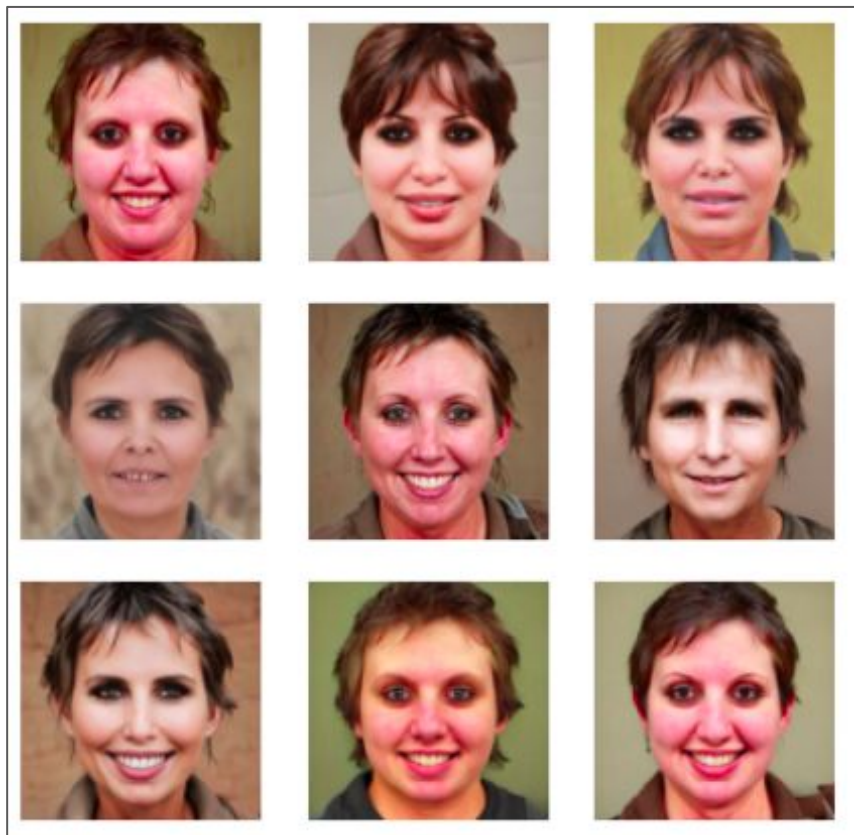
Male Audio



Female Audio



StyleWav Snapshot



T-test

- Perform Welch's t-test on two population of similarity score between two prompts and the same material
- Two mean are significantly different if p-value $\leq 5\%$

	Man Photos	Female Photos
this is a photo of a boy	0.26	0.241
this is a photo of a girl	0.252	0.256
have significantly two different means	yes	yes

Thank you for listening
ngoph@beloit.edu