

**Explicitly disentangling
speaker from speech**

**Implicit disentanglement all
aspects of speech**

This requires that we explicitly design
how to discard speaker information...

What if we can *learn* separation between speaker information
and the rest of speech \Rightarrow implicit disentanglement

Improving implicit speech disentanglement with GANs

Matthew Baas

Supervisor: Herman Kamper

August 2022



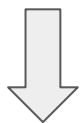
Disentanglement

Explicit disentanglement

*Speech synthesis
model designer after
much studying...*



Speaker identity is
different from the rest of
speech by **X**



Design cool new model using idea **X** to
disentangle and **modify speaker identity**

Limitation: human expert designed demarcations of
different characteristics of speech

Can we try learn them instead?

Implicit disentanglement



Good disentanglement
that is still easy to
control?



Design cool new model that
should discover a controllable representation
that disentangles **speaker identity**.

Task

1. Minimal assumptions about speech

2. Tractable training

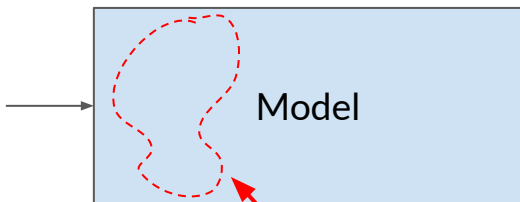
3. Practical for inference



**Unconditional speech
synthesis task**

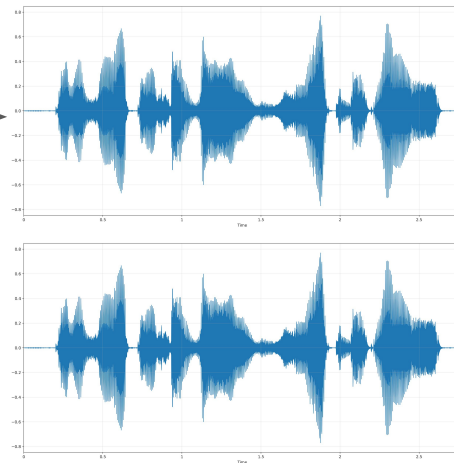
Common
setup

Known distribution



Latent space

We also
want

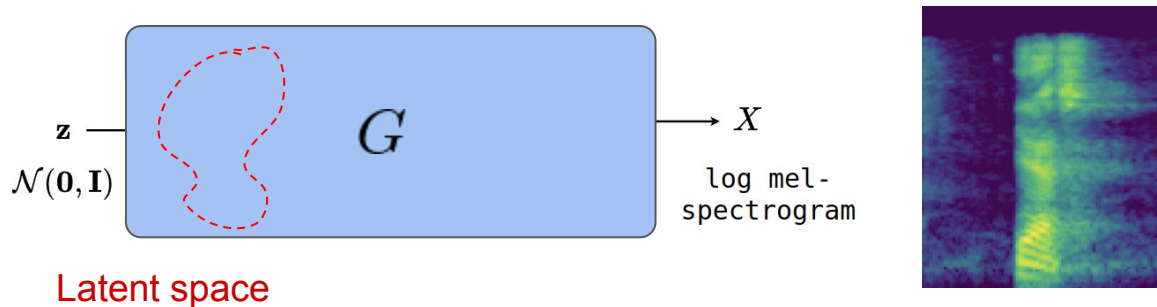


Related efforts

Similar effort in **unconditional image synthesis**:

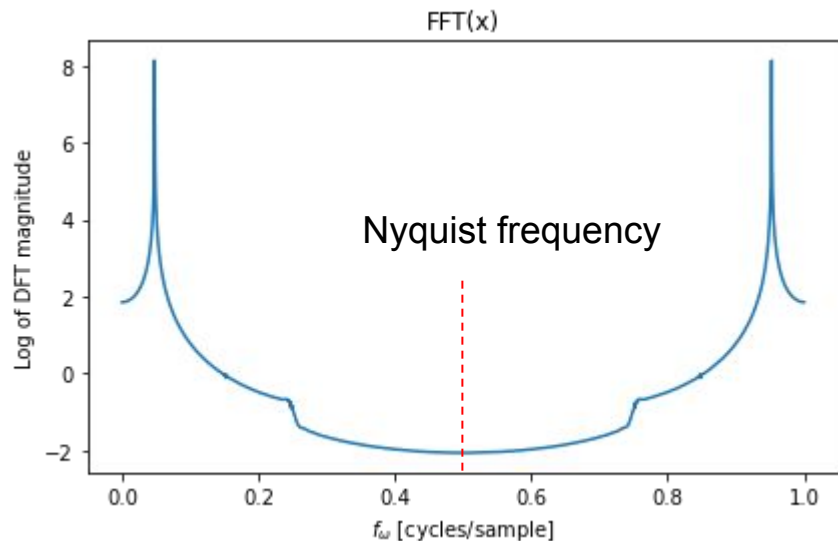
- Modelling a latent space from noise
 - StyleGAN
 - BigGAN
 - PPOGAN
- Encourage disentanglement
 - StyleGAN2
 - StyleGAN3
- Latent space metrics
 - FID
 - Path length
 - Linear separability
- Inverse mapping methods
 - Noisy projection
 - Pivot tuning

Overview of our approach



- Google Speech Commands dataset
- Great synthesis quality and diversity
- **Many tricks!**
- Apply to unseen tasks zero-shot
 - Voice conversion
 - Speech editing
 - Speech enhancement
 - Speaker verification

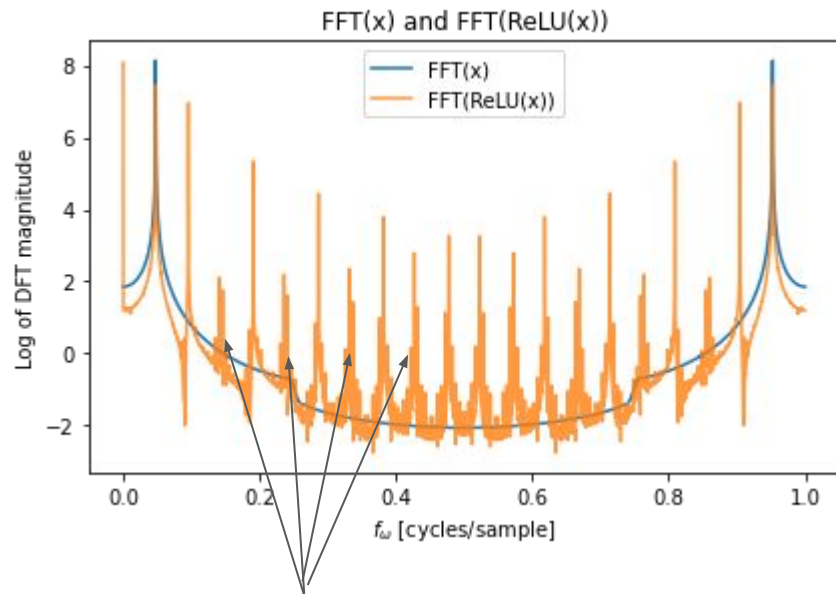
Anti-aliasing in GANs



Aliasing causes:

- GAN to hide artifacts in aliasing
- Poorer reconstruction quality when finally vocoding to audio

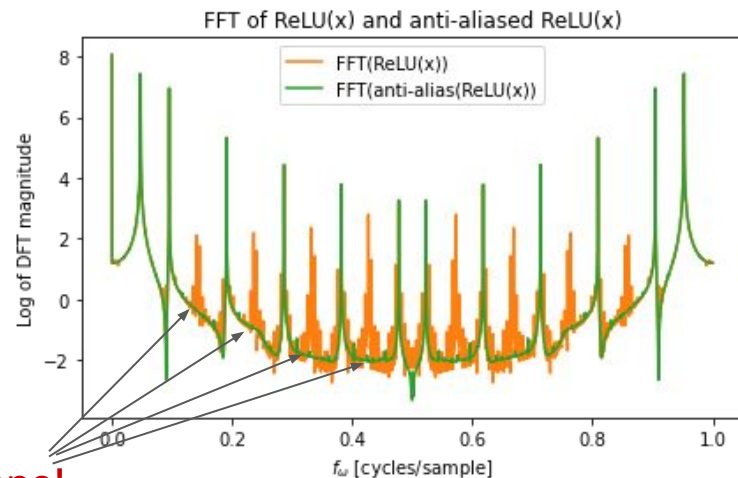
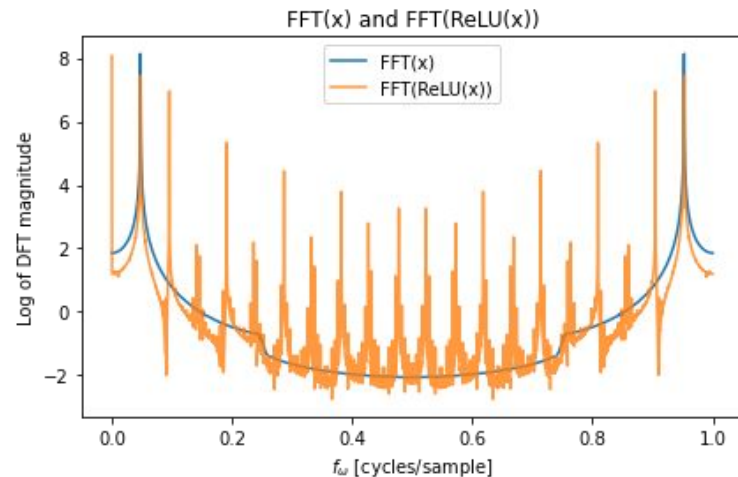
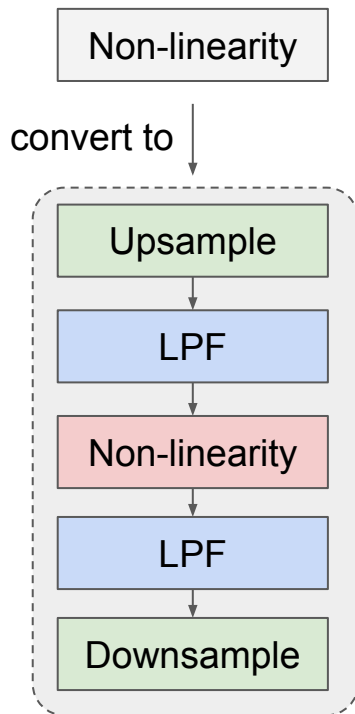
x is single sinusoid activation.



Aliasing! \Rightarrow need anti-aliasing filters!

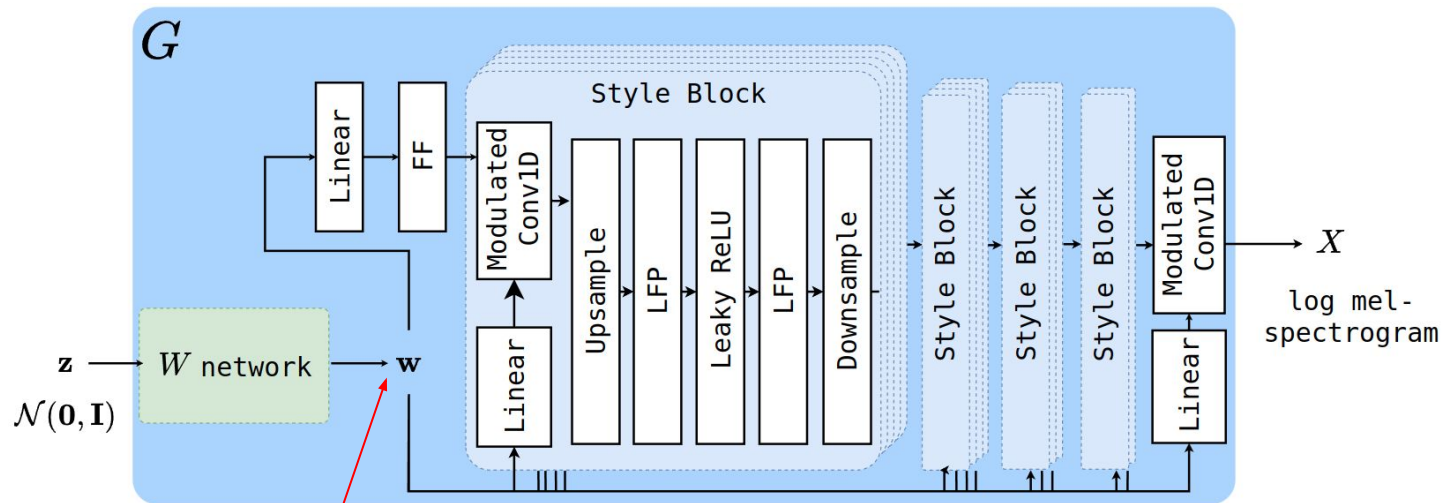
Anti-aliasing in GANs

StyleGAN3 proposed:



Aliasing gone!

Model & GAN tricks

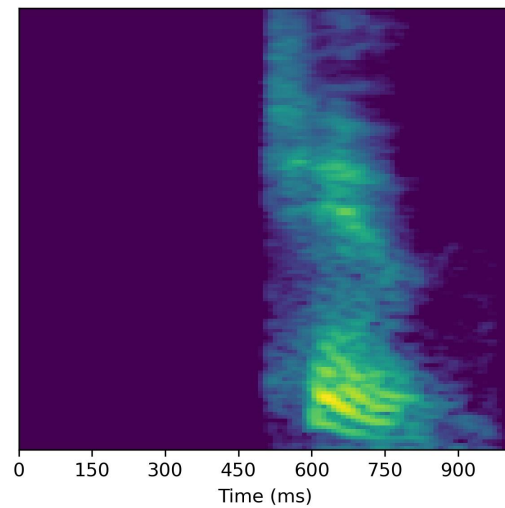
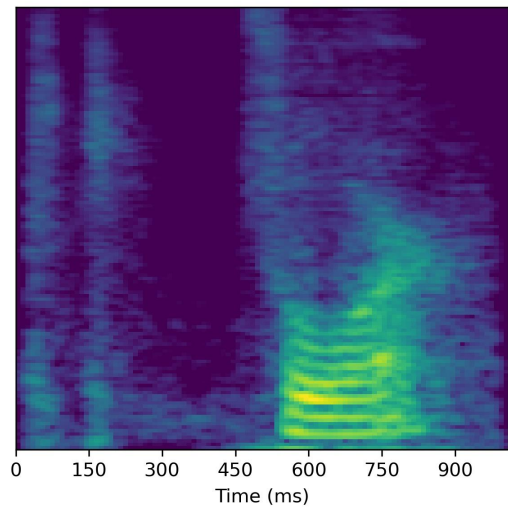
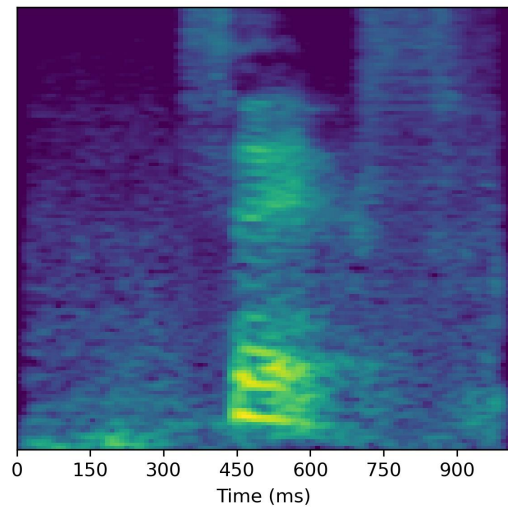


Disentangled
latent space

Tricks:

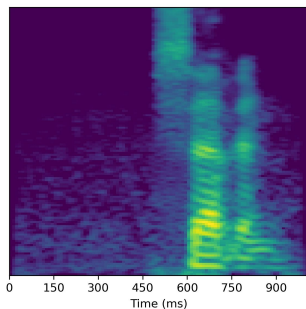
- R1 regularization
- EMA generator weights
- Adaptive discriminator augmentation
- Equalized learning rates
- Lower W network learning rate
- Adaptive discriminator updates

Unconditional samples



Voice conversion

Source utterance

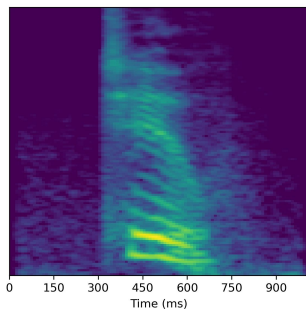


projection

w_1

Course styles

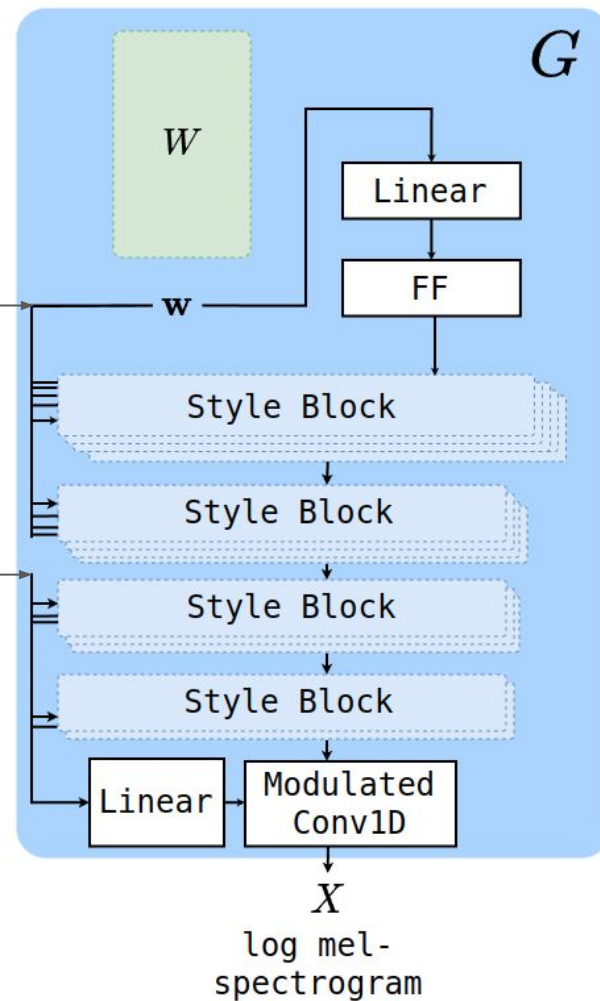
Utterance with target speaker



projection

w_2

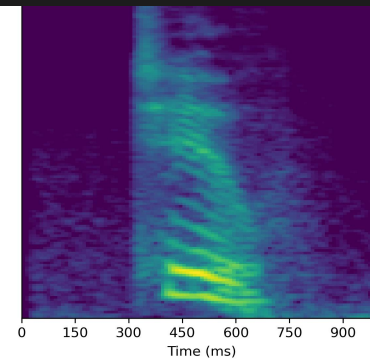
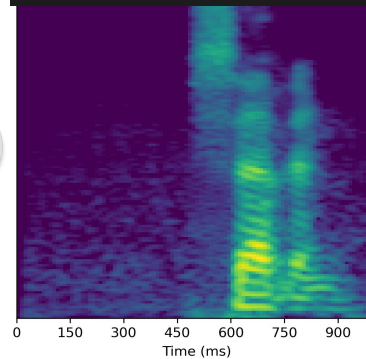
Fine styles



Voice conversion

Projected X_1

Projected X_2



Course styles: w_1

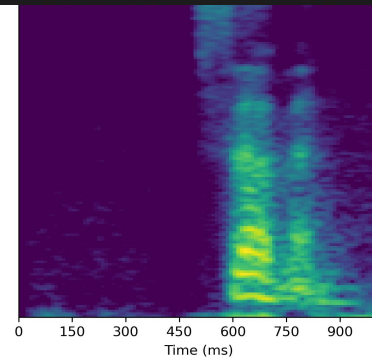
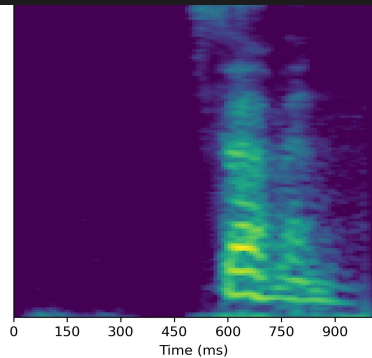
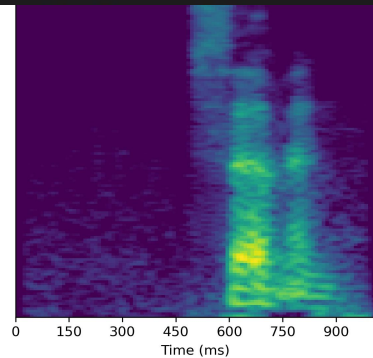
Course styles: w_1

Course styles: w_1

Fine styles: $w_1 + 1.0(w_2 - w_1)$

Fine styles: $w_1 + 1.5(w_2 - w_1)$

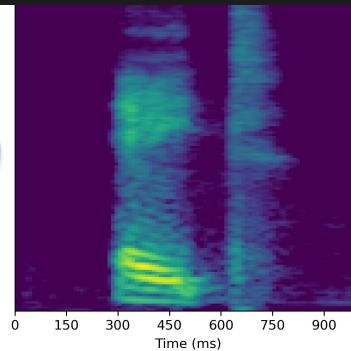
Fine styles: $w_1 + 2.0(w_2 - w_1)$



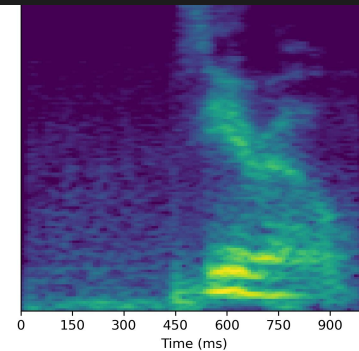
Digit conversion

Just do the opposite!

Projected X_1

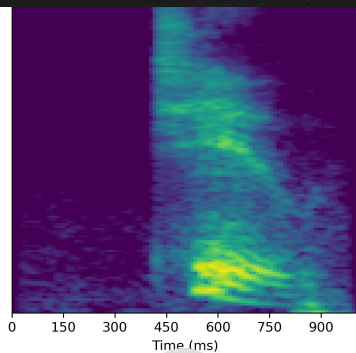


Projected X_2



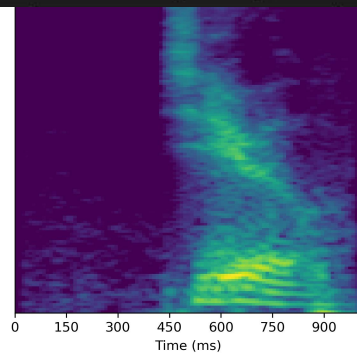
Course styles: $\mathbf{w}_1 + 0.8(\mathbf{w}_2 - \mathbf{w}_1)$

Fine styles: \mathbf{w}_1



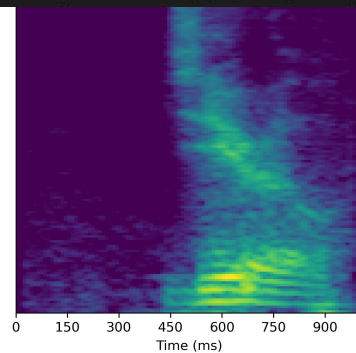
Course styles: $\mathbf{w}_1 + 0.95(\mathbf{w}_2 - \mathbf{w}_1)$

Fine styles: \mathbf{w}_1

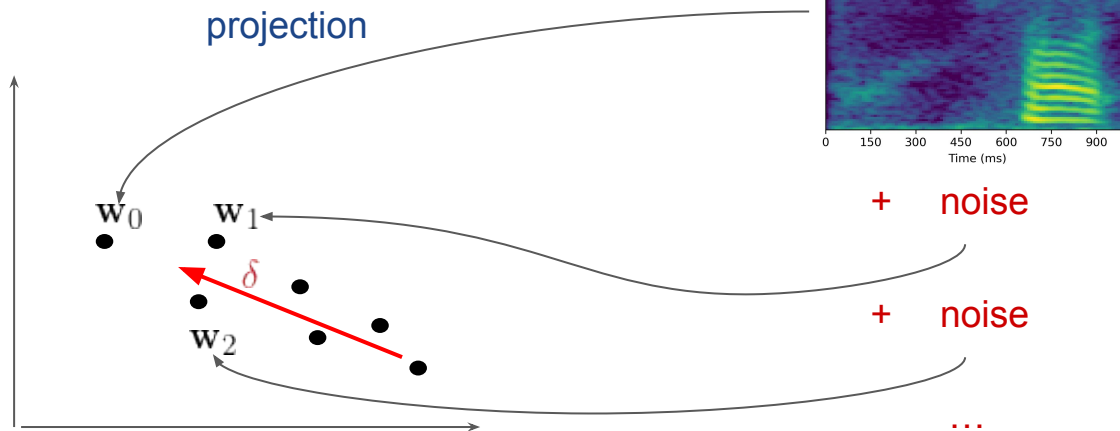


Course styles: $\mathbf{w}_1 + 1.0(\mathbf{w}_2 - \mathbf{w}_1)$

Fine styles: \mathbf{w}_1



Speech enhancement



We now know the
direction of decreasing
noise in W-space

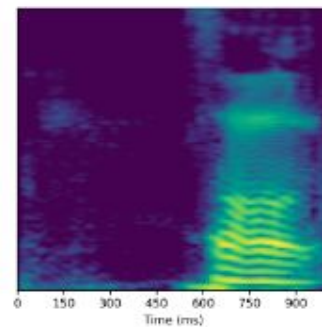
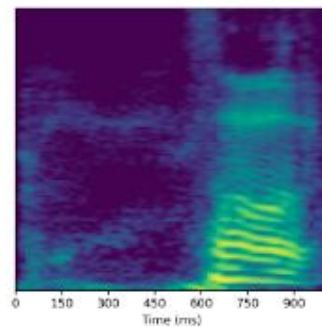
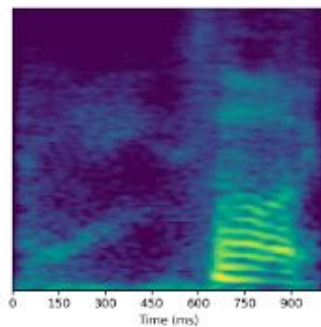
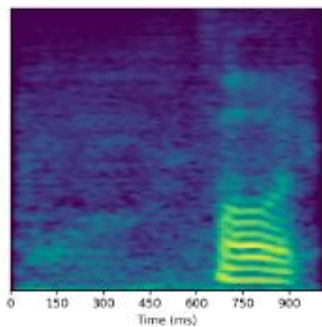
δ

All styles: $w_0 - 3\delta$

All styles: $w_0 + 3\delta$

All styles: $w_0 + 6\delta$

All styles: $w_0 + 9.0\delta$



Limitations & future directions

1. Other synthesis tasks phrased as latent manipulations / perceptual tests?
2. Does not work well on out-of-domain data
3. Struggles to scale
4. Projection method not sophisticated

Thank you