

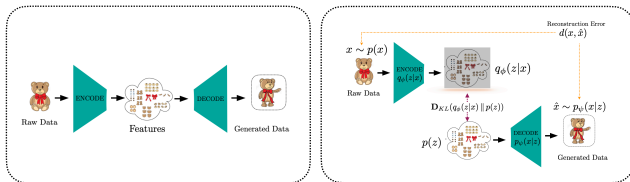
11-695: AI Engineering

Conditional Generation

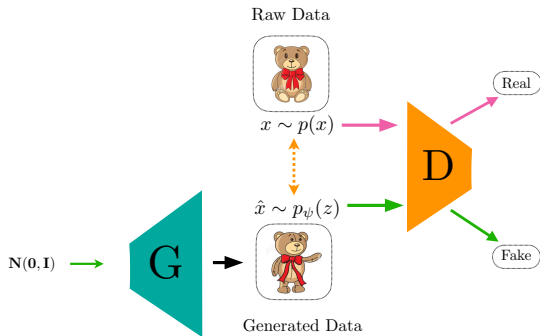
LTI/SCS

Spring 2020

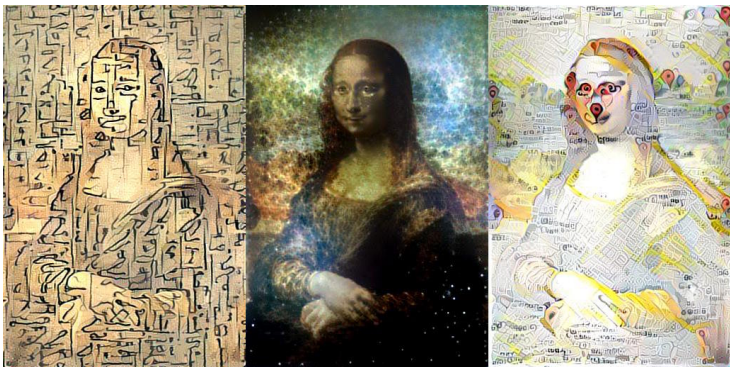
- 1 Generative Models: Review
- 2 Conditional Generative Models
- 3 Conditional Generation as Image/Video Style Transfer



- Assume the model is well-trained. Take MNIST as an example
- If we want to generate digit 1
 - AE can do with feeding 1 into E, but it's not generative
 - VAE: we need to sample from random noise \rightarrow no control!

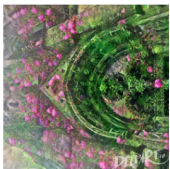
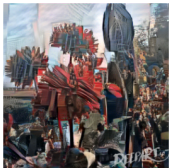


- Same as VAE, we generate from a noise \rightarrow no control!
- Upside: they both can generate completely new, cool images



- Can we generate images of digits from a single signal, *e.g.* from a scalar input?
- Can we adopt, say Van Gogh's style into images w/o formal drawing training?

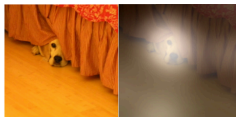
Image credit: Gene Kogan



- Can we change the style of input based on an random template?



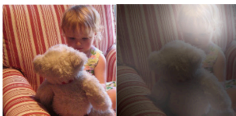
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

- Can we generate description based on images?
- Or vice versa?

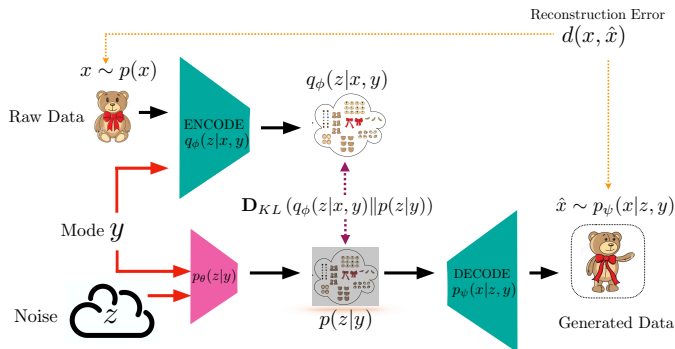
- 1 Generative Models: Review
- 2 Conditional Generative Models
- 3 Conditional Generation as Image/Video Style Transfer

- So far: we generate from a latent representation

$$z \sim p(z), \quad \hat{x} = p(x|z) \approx G_\psi(z) \quad (1)$$

- Now, we also have an attribute y as part of input:

$$z \sim p(z|y), \quad \hat{x} = p(x|z, y) \approx G_\psi(z, y) \quad (2)$$

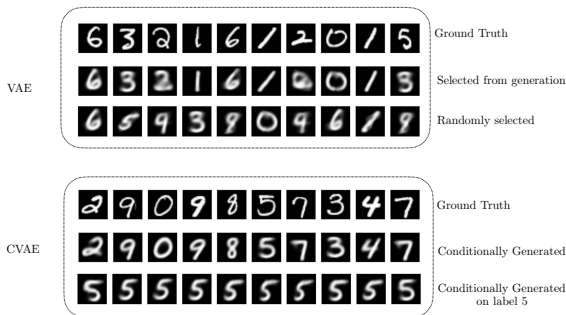


- VAE objective: $\text{ELBO} = \mathbf{E}_{q_\phi(z|x)}[\log p(x|z)] - \mathbf{D}_{KL}(q_\phi(z|x) || p(z))$
- Now CVAE:

$$\text{ELBO} = \mathbf{E}_{q_\phi(z|x, y)}[\log p(x|z, y)] - \mathbf{D}_{KL}(q_\phi(z|x, y) || p(z|y)) \quad (3)$$

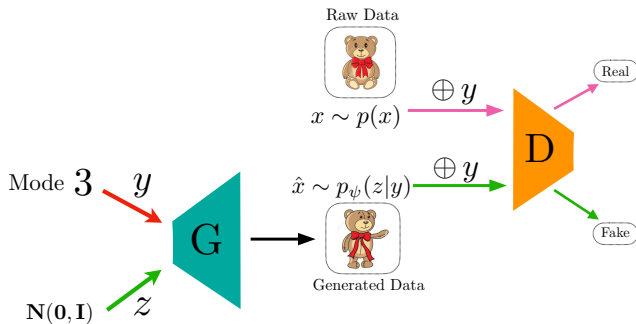
¹ <https://arxiv.org/abs/1406.5298>

² <https://pdfs.semanticscholar.org/3f25/e17eb717e5894e0404ea634451332f85d287.pdf>



- How to edit VAE to make CVAE:
 - Fuse y and input to Encoder, *e.g.* as simple as concatenation
 - Fuse y and noise with learnable weights before Decoder
- Results seem sharper (we implicitly make use of labels)
- And we can now control the the modes of generation!

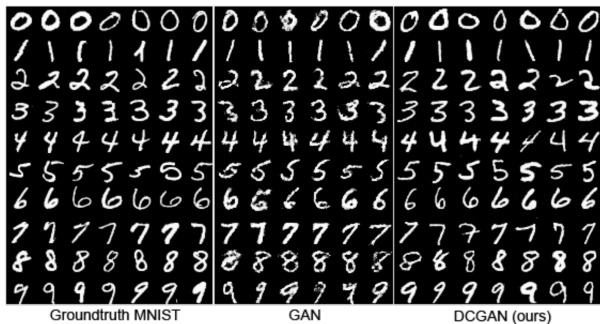
Image credit: Agustinus Kristiadi



- GAN: $\min_\psi \max_\theta \mathbb{E}_{x \sim p(x)} [\log D_\theta(x)] + \mathbb{E}_{z \sim p(z)} [1 - \log D_\theta(G_\psi(z))]$
- Now CGAN:

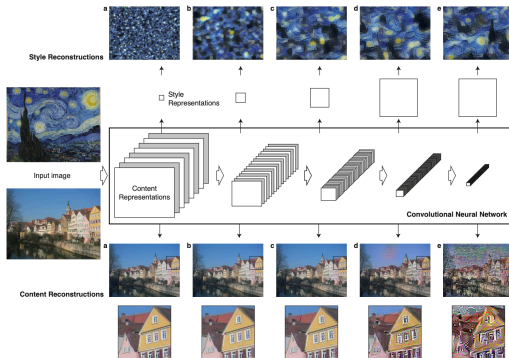
$$\min_\psi \max_\theta \mathbb{E}_{x \sim p(x)} [\log D_\theta(x|y)] + \mathbb{E}_{z \sim p(z|y)} [1 - \log D_\theta(G_\psi(z|y))] \quad (4)$$

³<https://arxiv.org/pdf/1411.1784.pdf>



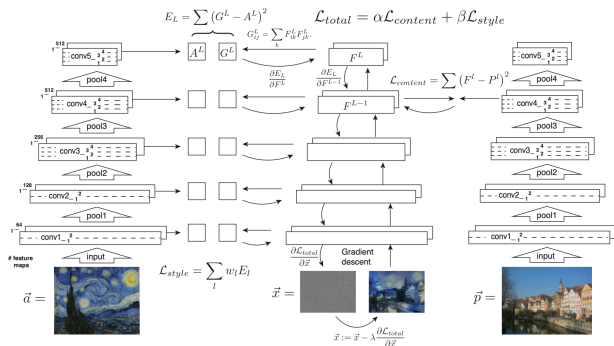
- How to edit GAN to make CGAN: similar to CVAE
- And also, we can now control the modes of generation

- ① Generative Models: Review
- ② Conditional Generative Models
- ③ Conditional Generation as Image/Video Style Transfer

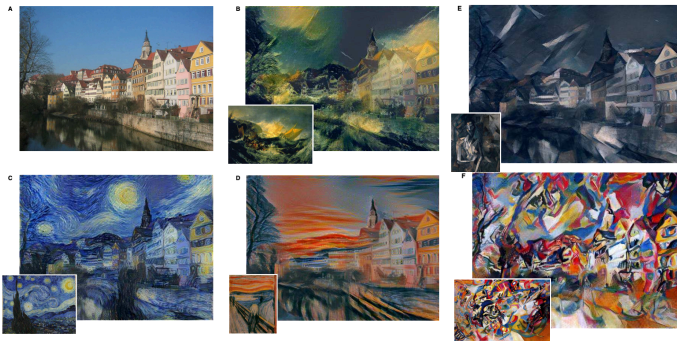


- Style Image and Content Image features are extracted by VGG

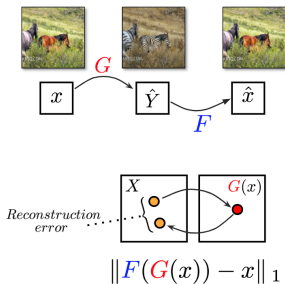
⁴[https:](https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/Gatys_Image_Style_Transfer_CVPR_2016_paper.pdf)



- White noise x is modeled into 2 components: content and style
- Each will have the respective loss
- Both losses are optimized together



- Training converges when x matches content and style

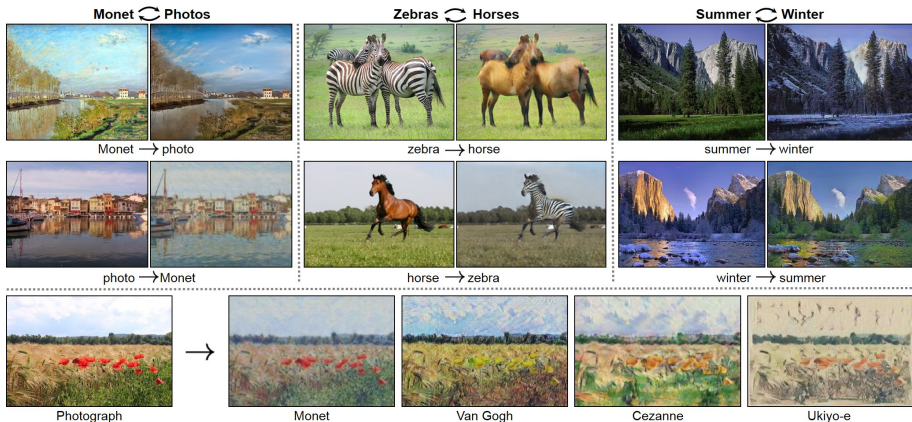


- Cycle loss is based on 2 phases: forward and back transfers
- Provides symmetric alignment intuition

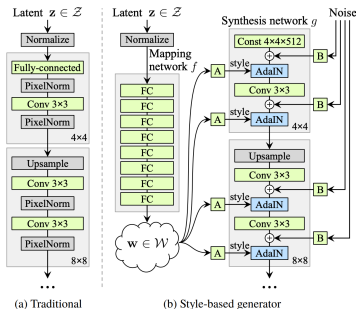


- 1 Generator,
- 2 Discriminators (1 per each style).
- Augmented by cycle loss

⁵<https://arxiv.org/pdf/1703.10593.pdf>

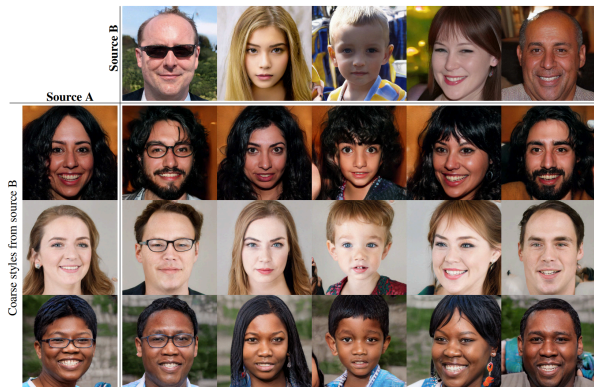


- Is also known as unpaired image-to-image translation



- 8-layer FNN mapping network f transforms z into styles w
- Each layer control each details: *e.g.* hair, pose, face, ...
- Adaptive instance norm (AdaIN) + noise are added to CNN layers
- G takes a constant instead of z directly

⁶<https://arxiv.org/pdf/1812.04948.pdf>



- Use two latent codes z_1 and z_2 during training
- Regulate contribution between them

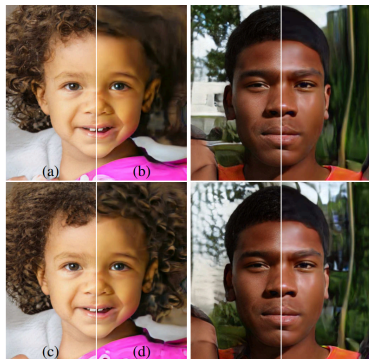
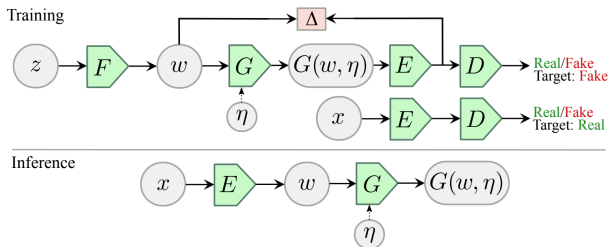
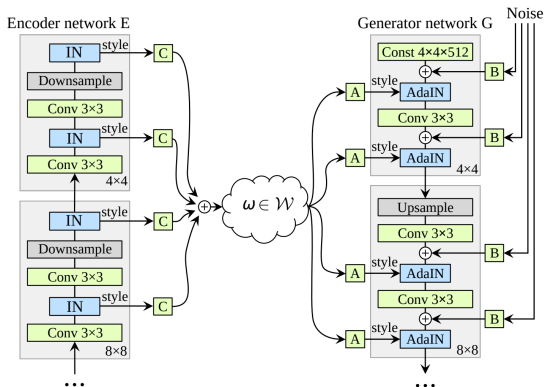


Figure 5. Effect of noise inputs at different layers of our generator. (a) Noise is applied to all layers. (b) No noise. (c) Noise in fine layers only ($64^2 - 1024^2$). (d) Noise in coarse layers only ($4^2 - 32^2$). We can see that the artificial omission of noise leads to featureless “painterly” look. Coarse noise causes large-scale curling of hair and appearance of larger background features, while the fine noise brings out the finer curls of hair, finer background detail, and skin pores.



- Idea: break GAN into smaller components
- Decompose: $G = G \circ F$ and $D = D \circ E$
- G and E act like Decoder and Encoder in auto-encoder

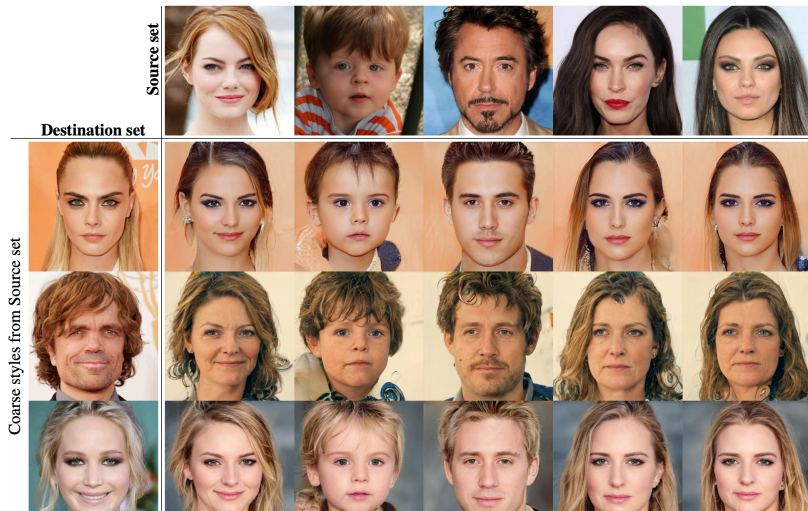


- Now E and G act the same roles as f and g networks in StyleGAN

⁷ <https://arxiv.org/pdf/2004.04467.pdf>



Figure 6: **FFHQ generations.** Generations with StyleALAE trained on FFHQ [24] at 1024×1024 .



- Multi-Content GAN for Few-Shot Font Style Transfer⁸
- Deep Photo Style Transfer⁹
- FSGAN: Subject Agnostic Face Swapping and Reenactment¹⁰
- Unsupervised Multimodal Video-to-Video Translation via Self-Supervised Learning¹¹
- Unpaired Photo-to-manga Translation Based on The Methodology of Manga Drawing¹²
- [▶ Code: Neural Style Transfer](#)
- [▶ Code: CycleGAN](#)

⁸ <https://arxiv.org/pdf/1712.00516v1.pdf>

⁹ <https://arxiv.org/pdf/1703.07511v3.pdf>

¹⁰ <https://arxiv.org/pdf/1908.05932.pdf>

¹¹ <https://arxiv.org/pdf/2004.06502.pdf>

¹² <https://arxiv.org/pdf/2004.10634.pdf>