

# Style transfer, GANs

Vlad Shakhuro



4 December 2025

# Outline

1. Metrics in image generation

2. Style transfer

3. Unconditional generation with GANs

4. Conditional generation with GANs



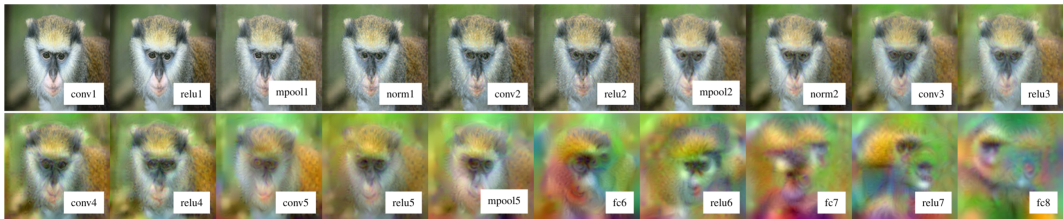
# Reconstructing images from neural features

Initialize  $x$  with white noise

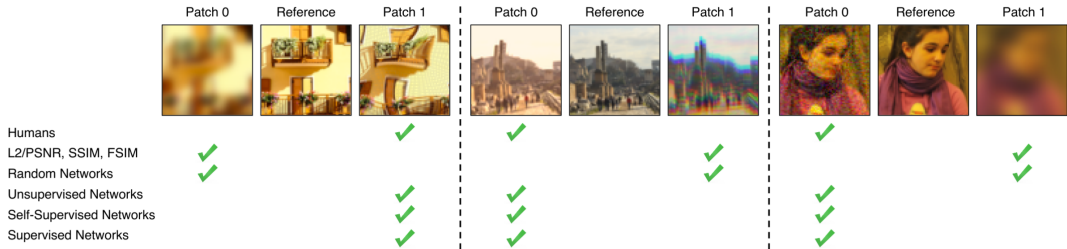
Optimize  $x$  using following loss:

$$x^* = \arg \min_{x \in \mathbb{R}^{H \times W \times C}} \|\Phi(x) - \Phi_0\|^2 + \lambda R(x)$$

$$R(x) = \sum_{i,j} (x_{i,j+1} - x_{i,j})^2 + (x_{i+1,j} - x_{i,j})^2$$



# LPIPS: comparing two images



Zhang et al. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. CVPR 2018

# Human comparison

Q: Which one is a real artwork?

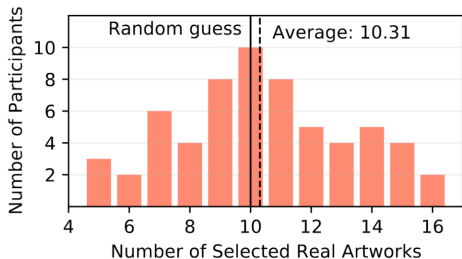


Figure 9. **(Left)** An example of our Artistic Style Transfer Confusion Test. Only 40.6% participants successfully distinguished the real artwork in this example. The answer can be found in our supplementary material. **(Right)** The statistical results with a total of 61 participants, where each participant is asked 20 questions.

# Inception score (IS)

For generated image  $x$ :

- $p(y|x)$  should have low entropy  
(generated image is somehow confidently classified)
- $\int p(y|x = G(z))dz$  should have high entropy  
(generated images are varied)

Inception score:

$$\text{IS}(X_s) = \mathbb{E}_{x \in X_s} KL(p(y|x) \parallel p(y))$$

Label prediction  $p(y|x)$  is computed using Inception model

**Drawback:** real images aren't used for computing metric

# Fréchet Inception Distance (FID)

Assume that image features computed with Inception model have normal distribution. Compute Fréchet (also called Wasserstein-2) distance between two gaussians:

$$\text{FID}(X_r, X_s) = \left\| \mu_{X_r} - \mu_{X_s} \right\|^2 - \text{Tr} \left( \Sigma_{X_r} + \Sigma_{X_s} - 2\sqrt{\Sigma_{X_r} \Sigma_{X_s}} \right)$$

## Drawbacks:

1. Inception embeddings aren't normally distributed
2. Estimating  $(2048 \times 2048)$ -dimensional covariance matrices from a small sample can lead to large errors
3. Has a bias that depends on the  $X_s$  model

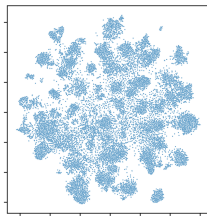


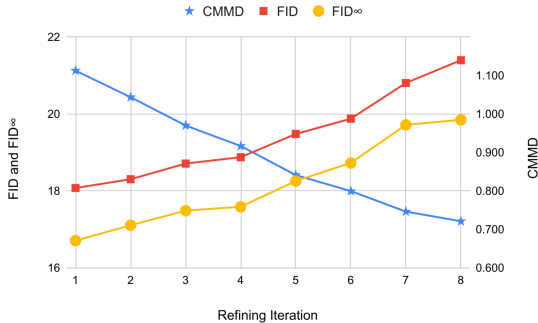
Figure 2. t-SNE visualization of Inception embeddings of the COCO 30K dataset. Note that even in the reduced-dimensional 2D representation, it is easy to identify that embeddings have multiple modes and do not follow a multivariate normal distribution.

Heusel et al. GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. NIPS 2017  
Chong, Forsyth. Effectively Unbiased FID and Inception Score and where to find them. CVPR 2020

# CLIP and Maximum Mean Discrepancy (CMMD)

Use CLIP embeddings (more diverse than Inception embs) and MMD distance:

$$\text{CMMD}(X_r, X_s) = \mathbb{E}_{x_r, x'_r} k(x_r, x'_r) + \mathbb{E}_{x_s, x'_s} k(x_s, x'_s) - 2 \mathbb{E}_{x_r, x_s} k(x_r, x_s)$$



# Outline

1. Metrics in image generation

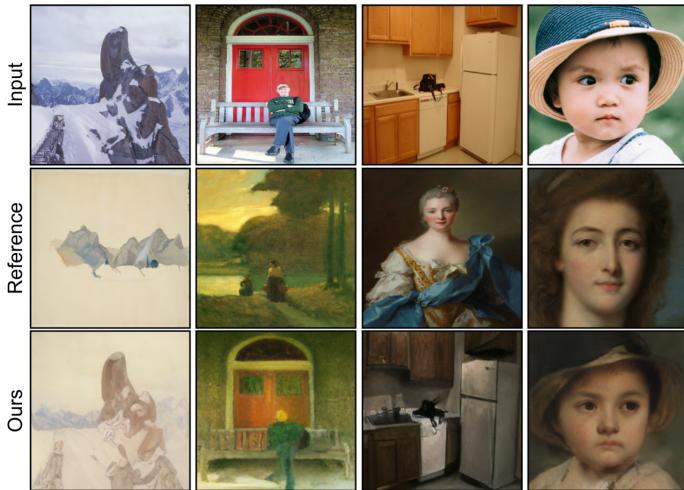
2. Style transfer

3. Unconditional generation with GANs

4. Conditional generation with GANs

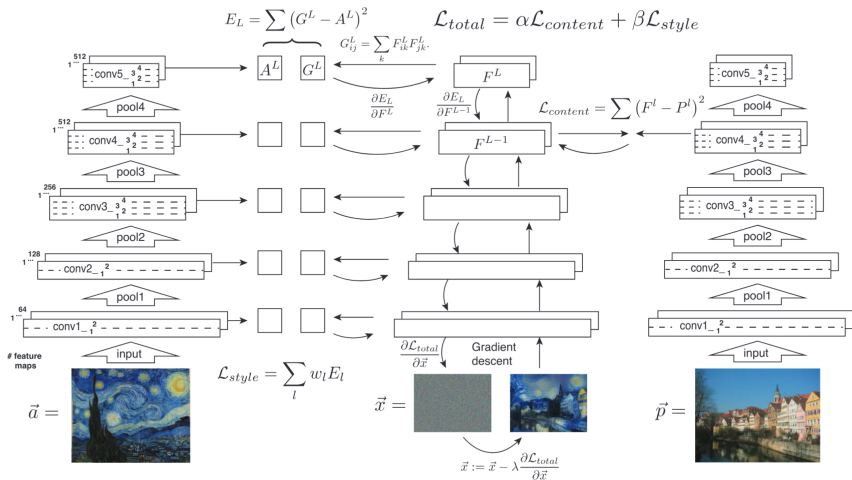
# Style transfer

Photo  $\rightarrow$  Art

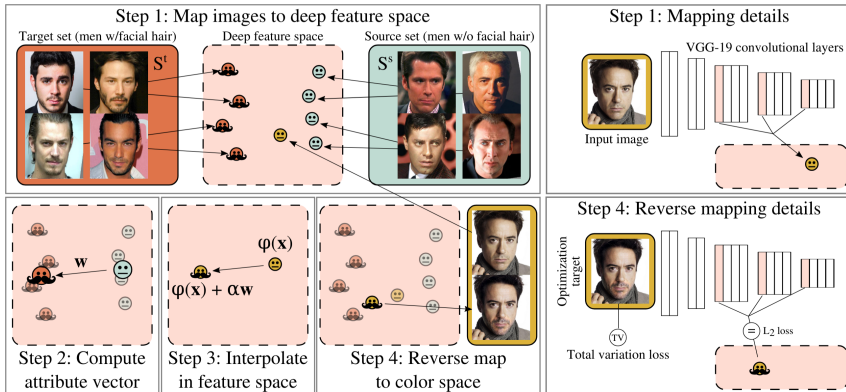




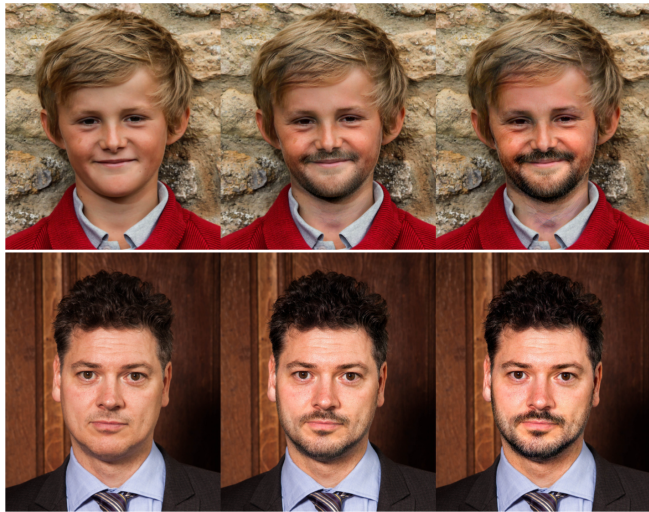
# Generating stylized images from noise



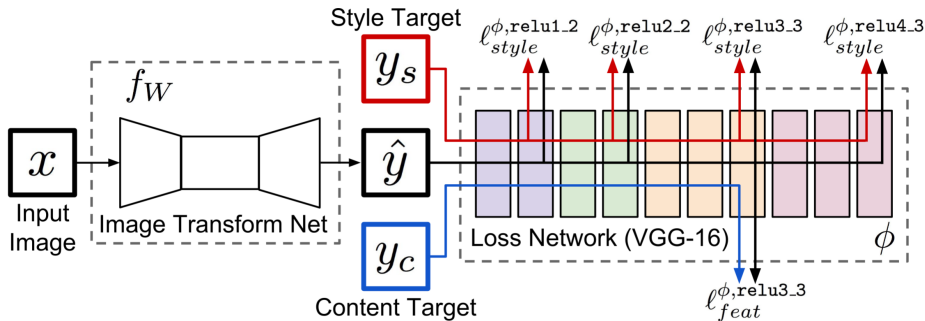
# Deep Feature Interpolation



# Deep Feature Interpolation results



# Training a neural network for a single style



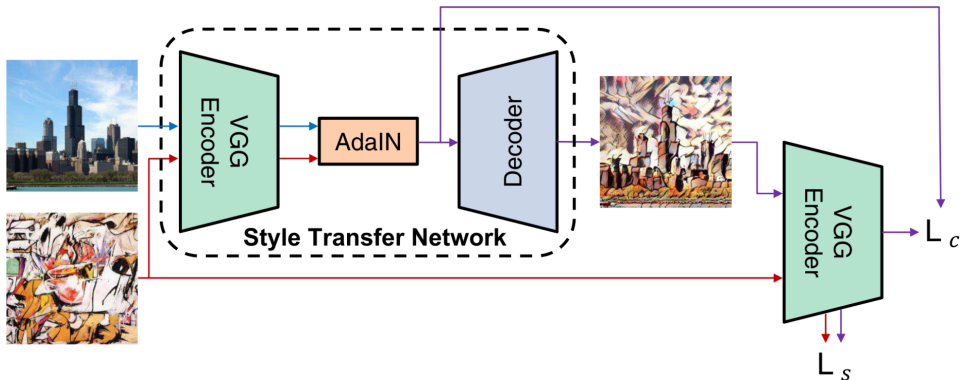
**Fig. 2.** System overview. We train an *image transformation network* to transform input images into output images. We use a *loss network* pretrained for image classification to define *perceptual loss functions* that measure perceptual differences in content and style between images. The loss network remains fixed during the training process.

# Instance Normalization



Ulyanov et al. Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis. ICCV 2017

# Adaptive Instance Normalization



$$AdaIN(x, y) = \sigma(y) \left( \frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$

# Outline

1. Metrics in image generation
2. Style transfer
3. Unconditional generation with GANs
4. Conditional generation with GANs

# Datasets: FFHQ



70k images of people with permissive license  
 $1024 \times 1024$  resolution



# Datasets: LandscapeHQ

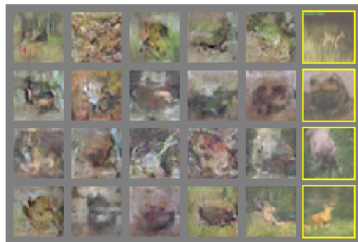
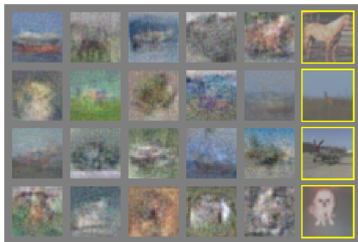
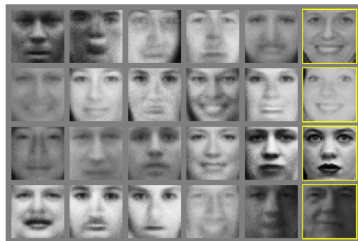
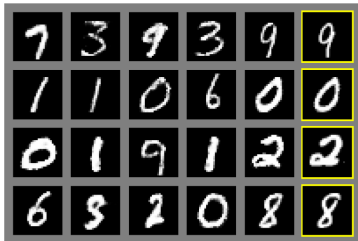


90k images of landscapes from Unsplash and Flickr with high resolution

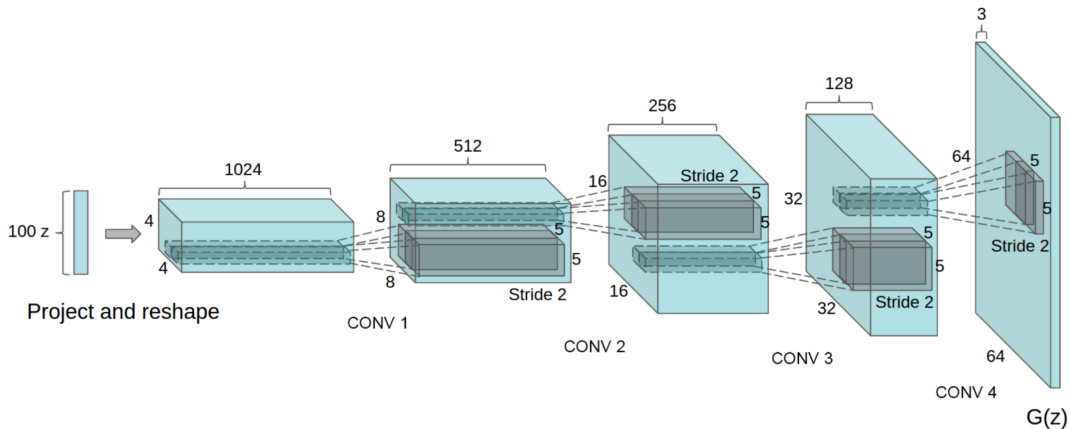
Skorokhodov et al. Aligning Latent and Image Spaces to Connect the Unconnectable. ICCV 2021

# GAN

# GAN results



# DCGAN



Radford, Metz. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. ICLR 2016

# Wasserstein loss function

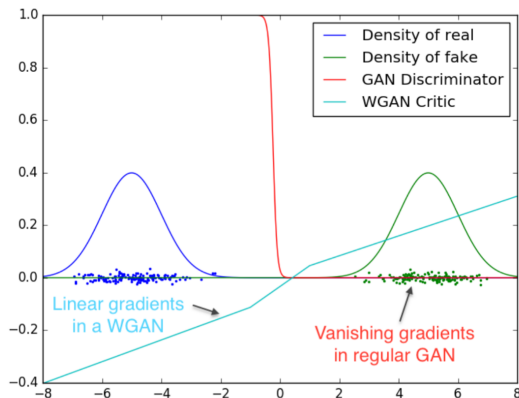


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.

# Wasserstein GAN training procedure

---

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

---

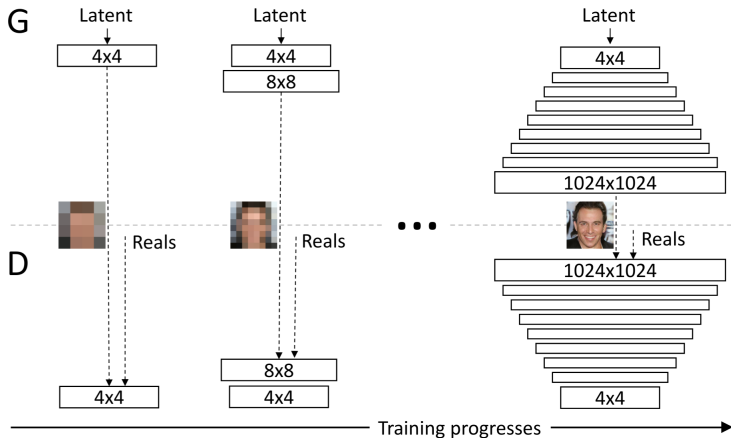
**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, g_\theta)$ 
12: end while
```

---

# Progressive GAN



Karras et al. Progressive Growing of GANs for Improved Quality, Stability, and Variation. ICLR 2018

# Progressive GAN

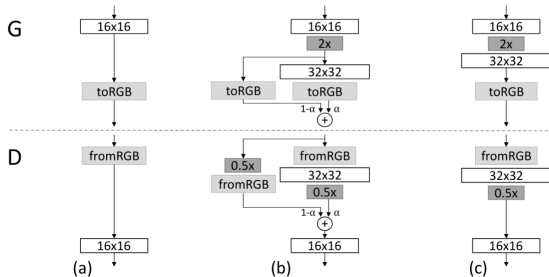
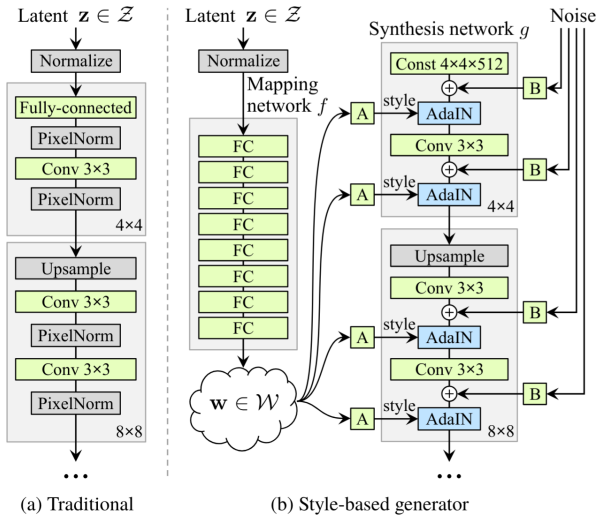


Figure 2: When doubling the resolution of the generator (G) and discriminator (D) we fade in the new layers smoothly. This example illustrates the transition from  $16 \times 16$  images (a) to  $32 \times 32$  images (c). During the transition (b) we treat the layers that operate on the higher resolution like a residual block, whose weight  $\alpha$  increases linearly from 0 to 1. Here  $2\times$  and  $0.5\times$  refer to doubling and halving the image resolution using nearest neighbor filtering and average pooling, respectively. The `toRGB` represents a layer that projects feature vectors to RGB colors and `fromRGB` does the reverse; both use  $1 \times 1$  convolutions. When training the discriminator, we feed in real images that are downsampled to match the current resolution of the network. During a resolution transition, we interpolate between two resolutions of the real images, similarly to how the generator output combines two resolutions.



# StyleGAN



# Style mixing in StyleGAN



# Outline

1. Metrics in image generation
2. Style transfer
3. Unconditional generation with GANs
4. Conditional generation with GANs

# Superresolution



Ground Truth



Bicubic



Ours ( $\ell_{pixel}$ )



SRCNN [11]

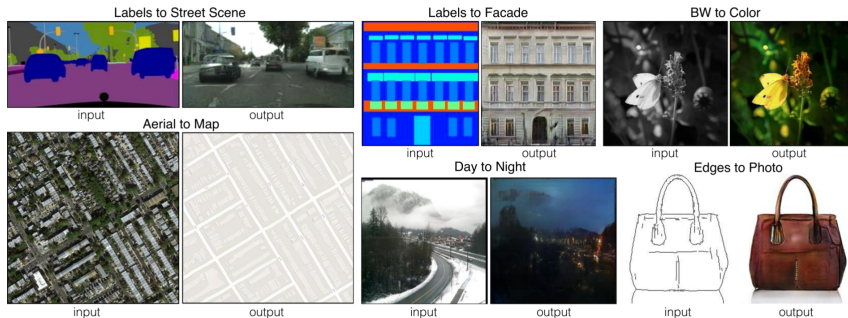
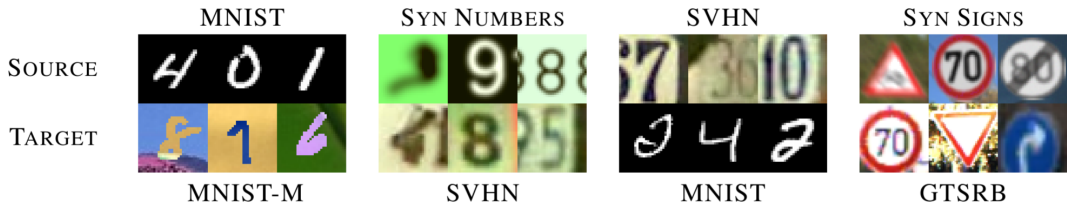


Ours ( $\ell_{feat}$ )

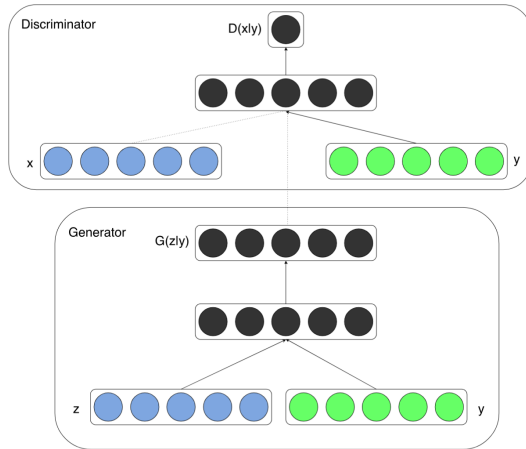
# Inpainting



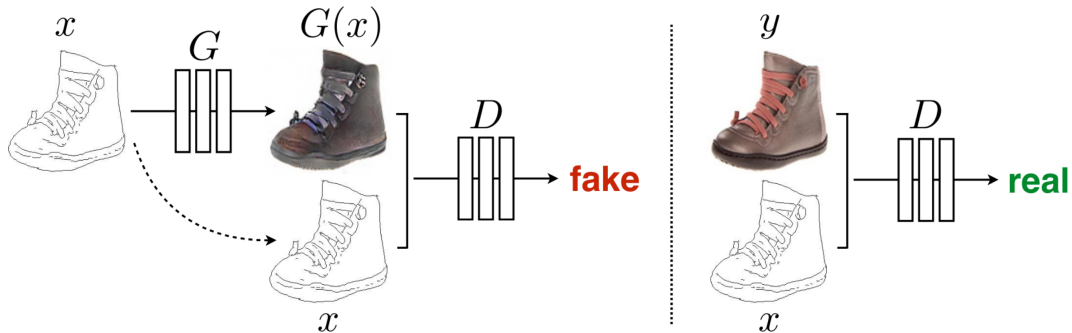
# Domain adaptation



# cGAN



# pix2pix



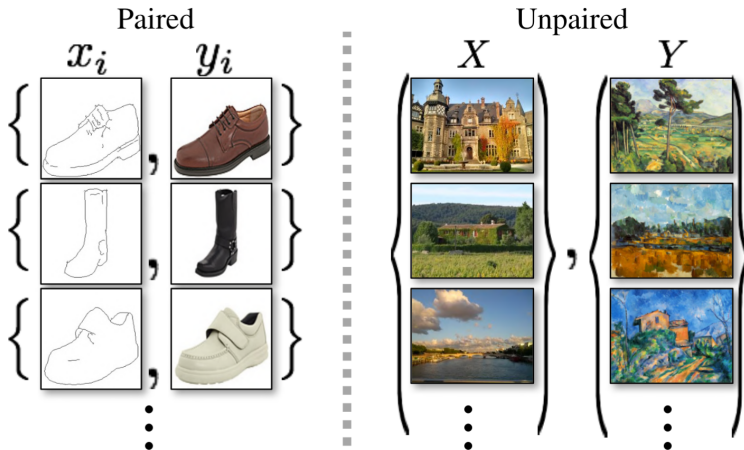
Isola et al. Image-to-Image Translation with Conditional Adversarial Networks. CVPR 2017



# pix2pix results



# Unpaired data



# CycleGAN

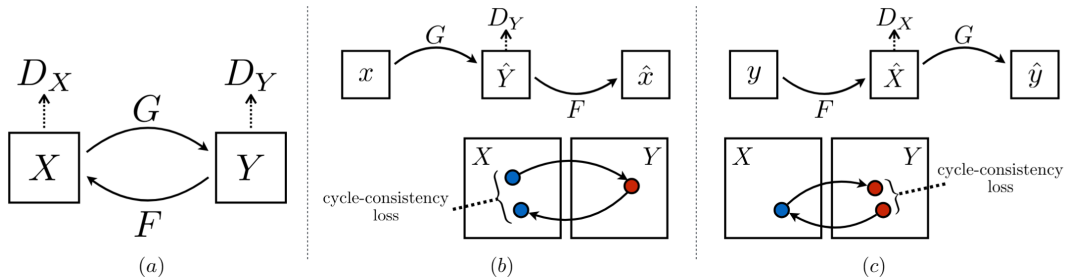
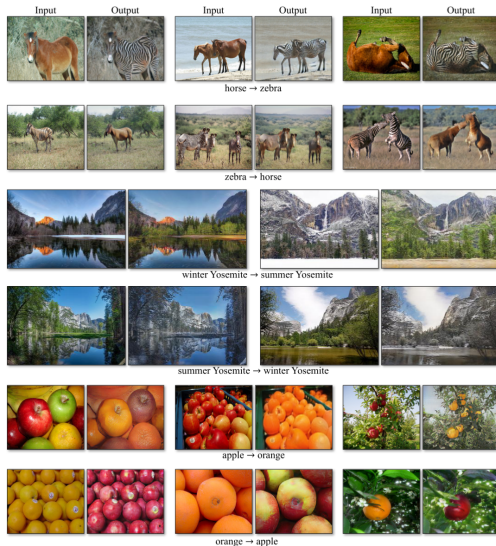
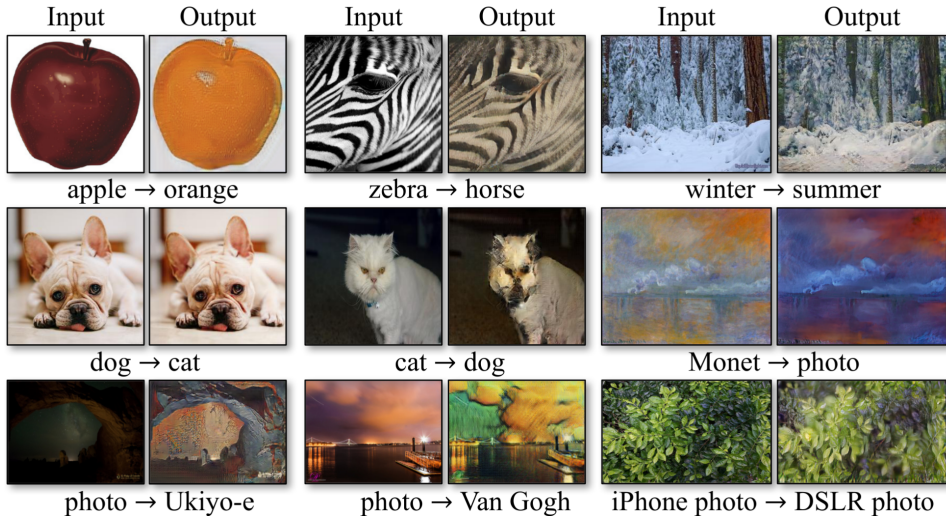


Figure 3: (a) Our model contains two mapping functions  $G : X \rightarrow Y$  and  $F : Y \rightarrow X$ , and associated adversarial discriminators  $D_Y$  and  $D_X$ .  $D_Y$  encourages  $G$  to translate  $X$  into outputs indistinguishable from domain  $Y$ , and vice versa for  $D_X$  and  $F$ . To further regularize the mappings, we introduce two *cycle consistency losses* that capture the intuition that if we translate from one domain to the other and back again we should arrive at where we started: (b) forward cycle-consistency loss:  $x \rightarrow G(x) \rightarrow F(G(x)) \approx x$ , and (c) backward cycle-consistency loss:  $y \rightarrow F(y) \rightarrow G(F(y)) \approx y$

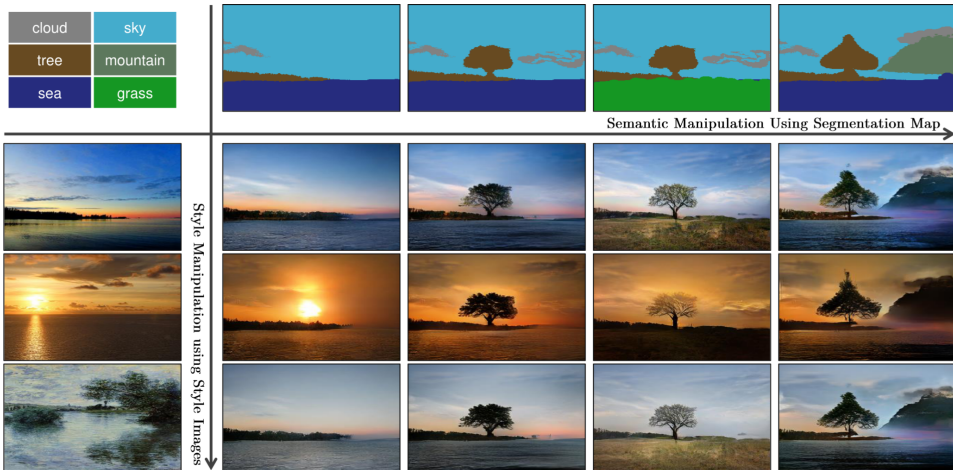
# CycleGAN results



# CycleGAN failures



# SPADE



# AdaIN in SPADE

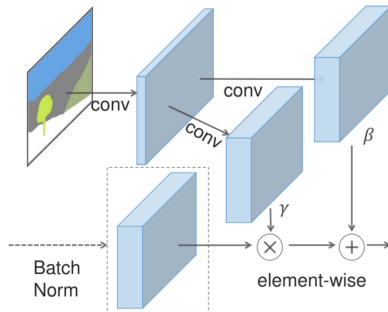


Figure 2: In the SPADE, the mask is first projected onto an embedding space and then convolved to produce the modulation parameters  $\gamma$  and  $\beta$ . Unlike prior conditional normalization methods,  $\gamma$  and  $\beta$  are not vectors, but tensors with spatial dimensions. The produced  $\gamma$  and  $\beta$  are multiplied and added to the normalized activation element-wise.

# SPADE

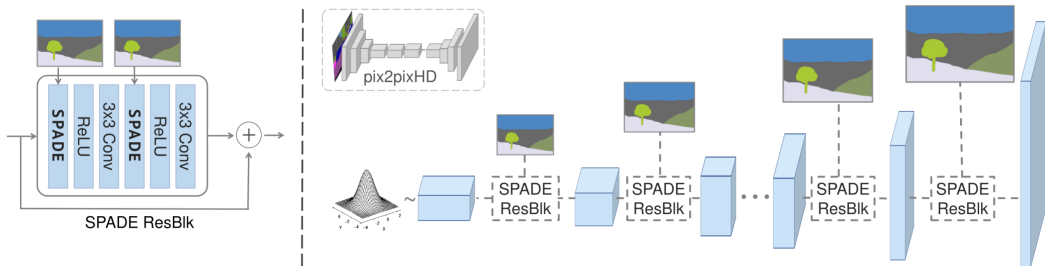
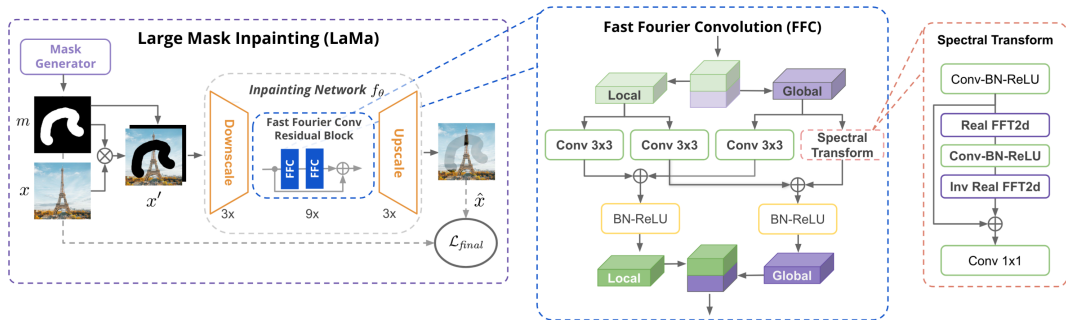


Figure 4: In the SPADE generator, each normalization layer uses the segmentation mask to modulate the layer activations. (left) Structure of one residual block with the SPADE. (right) The generator contains a series of the SPADE residual blocks with upsampling layers. Our architecture achieves better performance with a smaller number of parameters by removing the downsampling layers of leading image-to-image translation networks such as the pix2pixHD model [48].

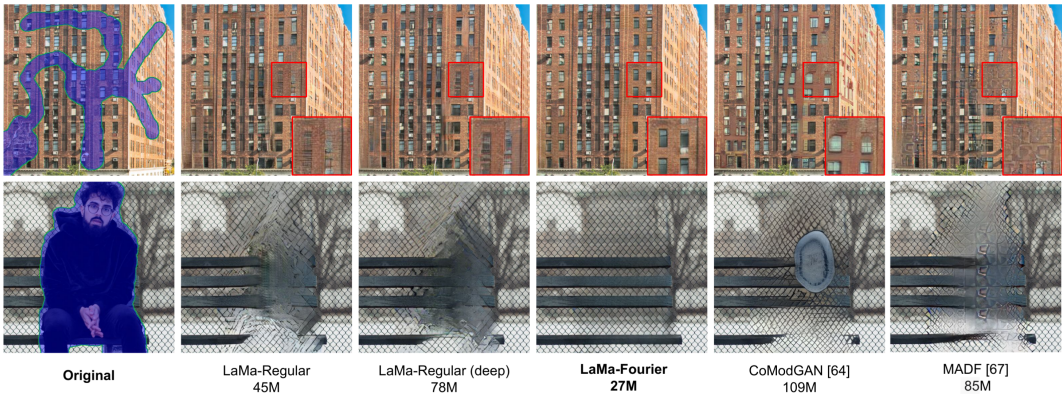


# LAMA



Suvorov et al. Resolution-robust Large Mask Inpainting with Fourier Convolutions. WACV 2022

# LAMA results



# LAMA high-res results

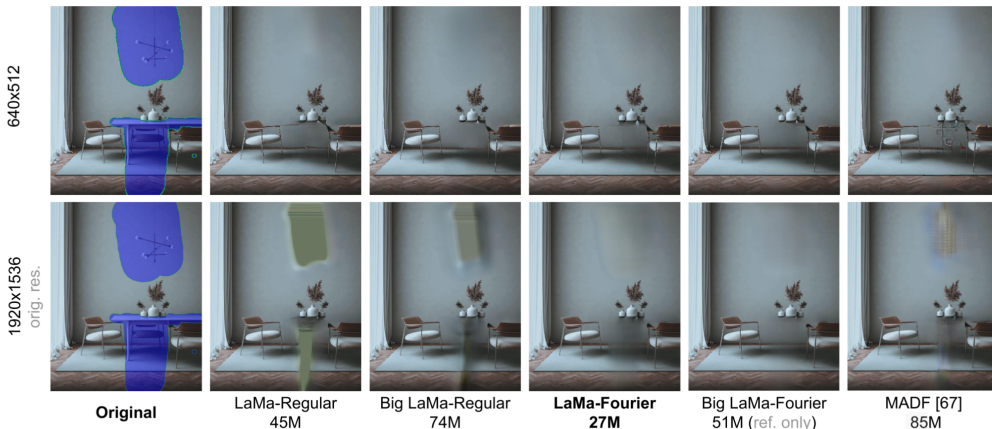
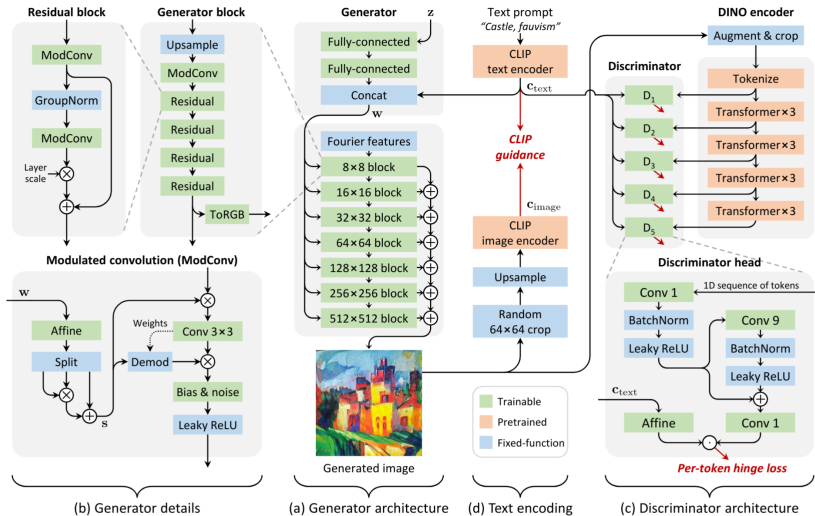


Figure 5: Transfer of inpainting models to a higher resolution. All LaMa models were trained using  $256 \times 256$  crops from  $512 \times 512$ , and MADF [67] was trained on  $512 \times 512$  directly. As the resolution increases, the models with regular convolutions swiftly start to produce critical artifacts, while FFC-based models continue to generate semantically consistent image with fine details. More negative and positive examples of our 51M model can be found at [bit.ly/3k0gaIK](https://bit.ly/3k0gaIK).

# StyleGAN-T



# StyleGAN-T



A 4k DSLR photo of a cute lion cub floating in a bowl of honey.

The Tower of Babel by J.M.W. Turner



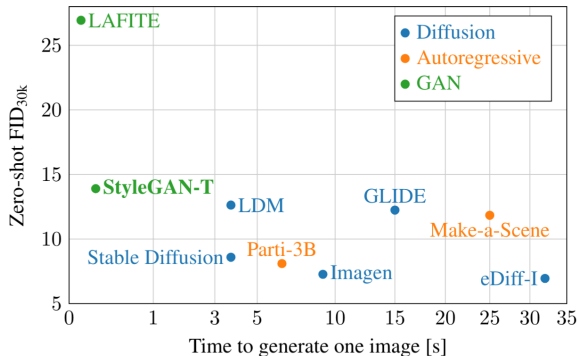
A fog rolling into new york



A forest rendered in the Unreal Engine.

A painting of a fox in the style of starry night.

# Comparison with diffusion models



*Figure 1. **Quality vs. speed** in large-scale text-to-image synthesis. StyleGAN-T greatly narrows the quality gap between GANs and other model families while generating samples at a rate of 10 FPS on an NVIDIA A100. The  $y$ -axis corresponds to zero-shot FID on MS COCO at  $256 \times 256$  resolution; lower is better.*

# Conclusion

We reviewed following topics:

- reconstructing images from neural features
- various metrics on top of neural features used for image comparison
- style transfer: optimization-based and training networks for single and several styles
- unconditional image generation with GANs
- conditional image generation with GANs