

10707

Deep Learning: Spring 2021

Andrej Risteski

Machine Learning Department

Lecture 16:

Applications of GANs, normalizing flows

Conditional GANs (Mirza-Osindero '14)

What if we want to have a notion of “class” and have class-labeled data?



Figure 2: Generated MNIST digits, each row conditioned on one label

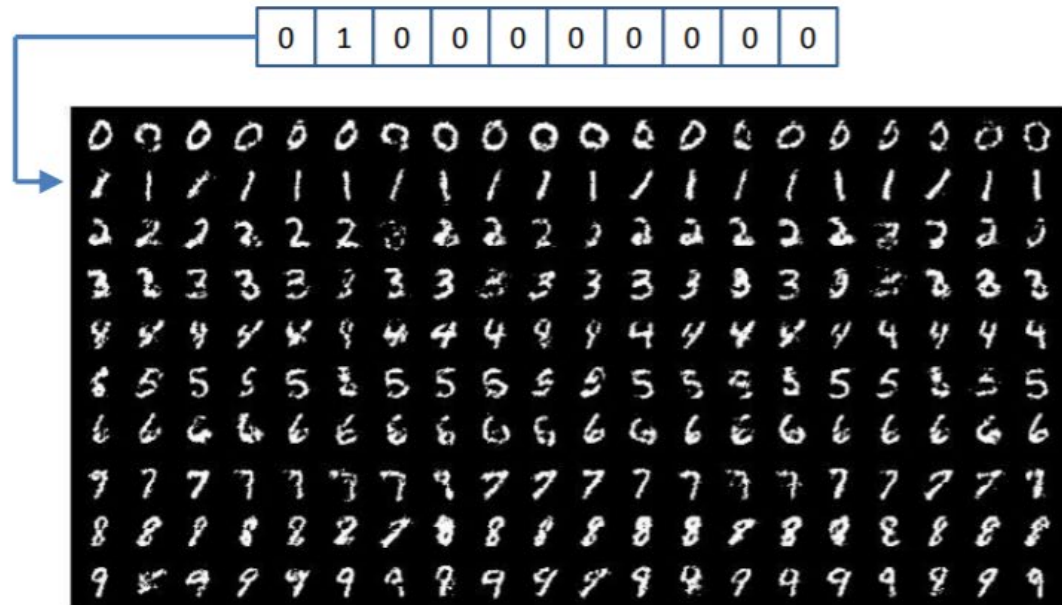
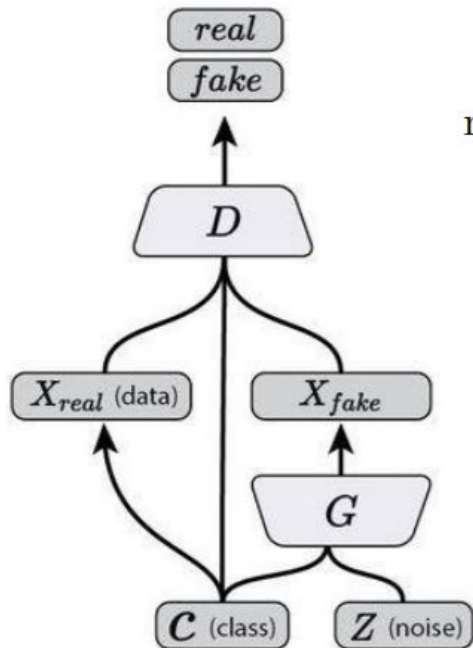
Figure from Mirza-Osindero '14



Conditional GANs (Mirza-Osindero '14)

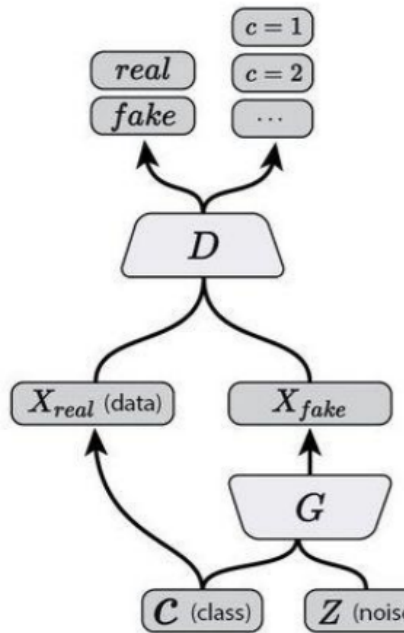
Discriminator and generator receive an image as well as a class label. The **loss** is then:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z}|\mathbf{y})))]$$



Auxillary classifier GANs (Odena et al '16)

Idea: give discriminator **only** image, and ask it to both predict label, and tell real/fake.



Let us define the “**distinguishing**” loss and the **class prediction** loss:

$$L_S = E[\log P(S = \text{real} \mid X_{\text{real}})] + E[\log P(S = \text{fake} \mid X_{\text{fake}})]$$

$$L_C = E[\log P(C = c \mid X_{\text{real}})] + E[\log P(C = c \mid X_{\text{fake}})]$$

Discriminator is trained to maximize $L_S + L_C$

Generator is trained to maximize $L_C - L_S$

(i.e. discriminator tries to both distinguish and predict label; generator tries to fool discriminator and generate “classifiable” images)

Main benefit: class-independent latent embedding; stabler training

Auxillary classifier GANs (Odena et al '16)



Figure 1. 128×128 resolution samples from 5 classes taken from an AC-GAN trained on the ImageNet dataset. Note that the classes shown have been selected to highlight the success of the model and are not representative. Samples from all ImageNet classes are linked later in the text.

Figure from Odena et al '16

Superresolution from a single image (Ledig et al '17)

Task: estimate a high-resolution version of an image, given a low-resolution version of it.

Training data: high-resolution images, along with manually created low-res versions (e.g. by Gaussian filtering).

Superresolution from a single image (Ledig et al '17)

Outperforms approaches based on “pre-designed” features/losses.



Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

Superresolution from a single image (Ledig et al '17)

Combines l2 loss on VGG-19 features with adversarial loss:

$$l_{VGG/i,j}^{SR} = \frac{1}{W_{i,j}H_{i,j}} \sum_{x=1}^{W_{i,j}} \sum_{y=1}^{H_{i,j}} (\phi_{i,j}(I^{HR})_{x,y} - \phi_{i,j}(G_{\theta_G}(I^{LR}))_{x,y})^2$$

$\phi_{i,j}$ The feature map for the j-th convolution (after activation) before the i-th maxpooling layer.

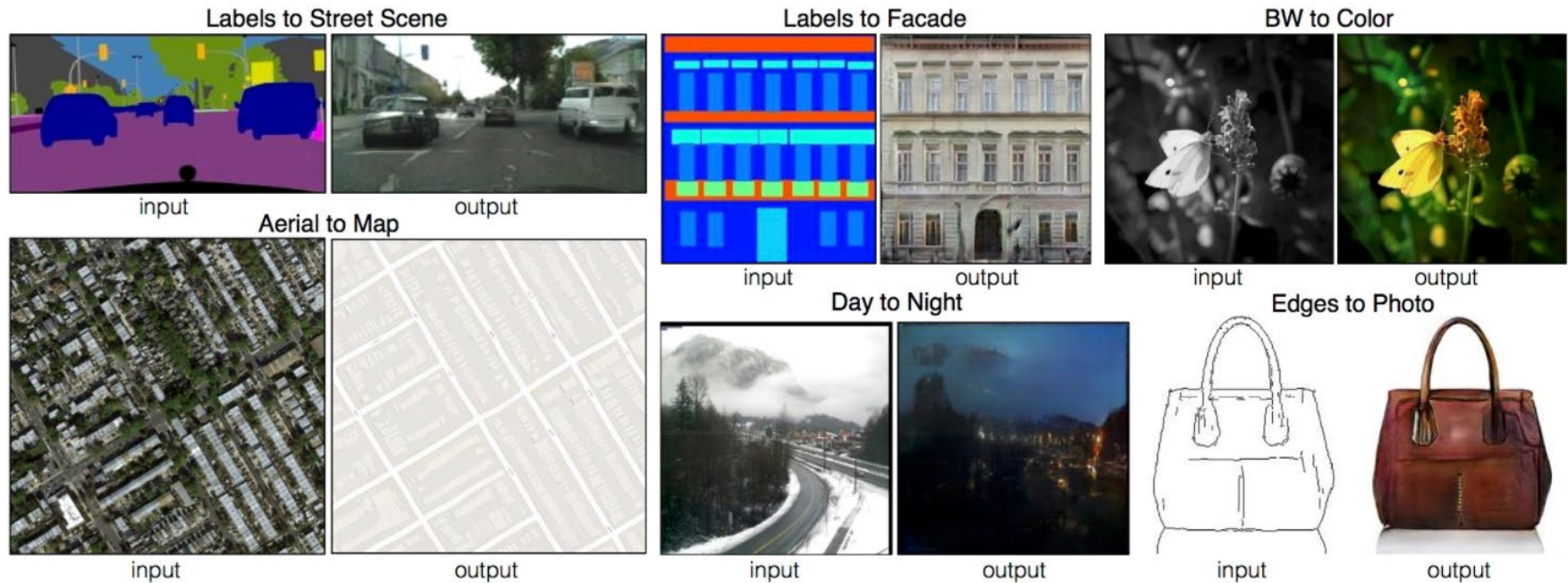
$$\min_{\theta_G} \max_{\theta_D} \mathbb{E}_{I^{HR} \sim p_{\text{train}}(I^{HR})} [\log D_{\theta_D}(I^{HR})] + \mathbb{E}_{I^{LR} \sim p_G(I^{LR})} [\log(1 - D_{\theta_D}(G_{\theta_G}(I^{LR})))]$$



Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

Image-to-image translation (Pix2Pix, Isola et al '17)

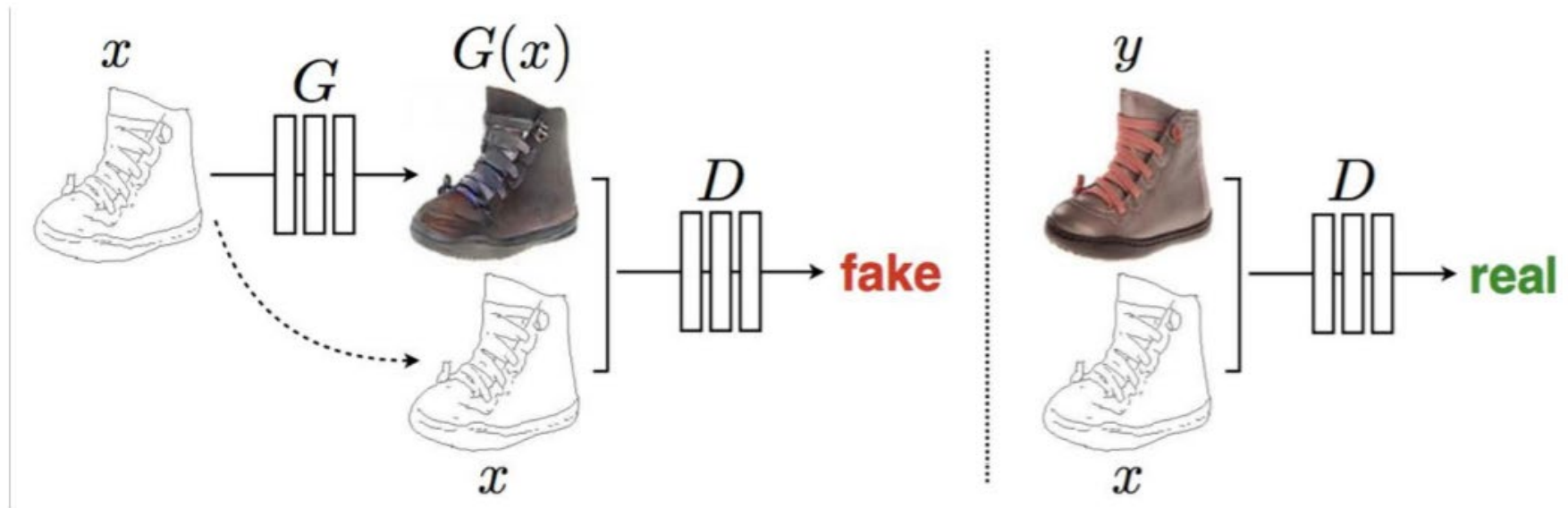
Goal: “translate” one “style” of image to another, given as training paired images of the styles.



Example results on several image-to-image translation problems. In each case we use the same architecture and objective, simply training on different data.

Image-to-image translation (Pix2Pix, Isola et al '17)

Loss: generator tries to “translate”; discriminator tries to distinguish whether it’s a “real” translation or “fake” translation.



$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))]$$

Image-to-image translation (Pix2Pix, Isola et al '17)



Figure 16: Example results of our method on automatically detected edges→handbags, compared to ground truth.

Image-to-image translation (Pix2Pix, Isola et al '17)

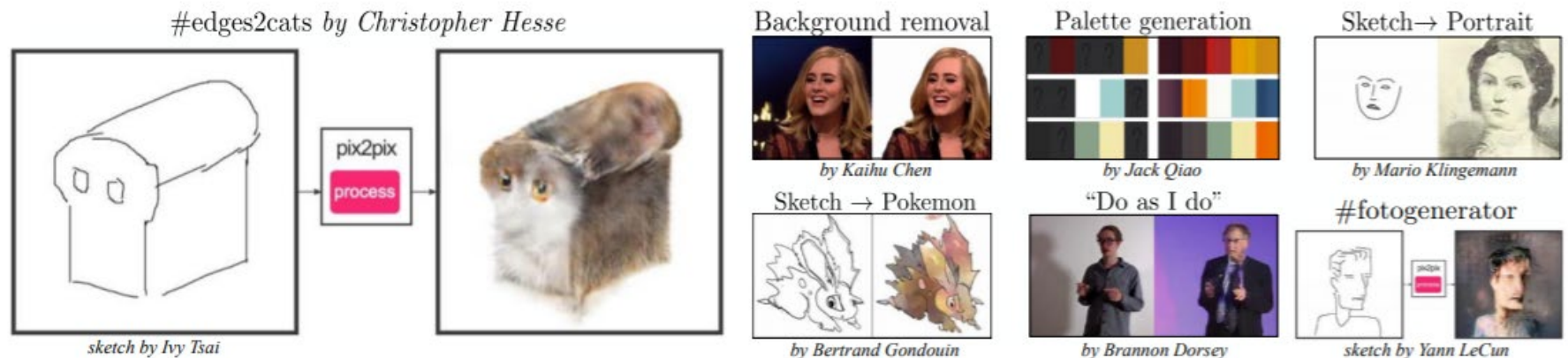


Figure 11: Example applications developed by online community based on our `pix2pix` codebase: `#edges2cats` [3] by Christopher Hesse, *Background removal* [6] by Kaihu Chen, *Palette generation* [5] by Jack Qiao, *Sketch → Portrait* [7] by Mario Klingemann, *Sketch → Pokemon* [1] by Bertrand Gondouin, *"Do As I Do" pose transfer* [2] by Brannon Dorsey, and `#fotogenerator` by Bosman et al. [4].

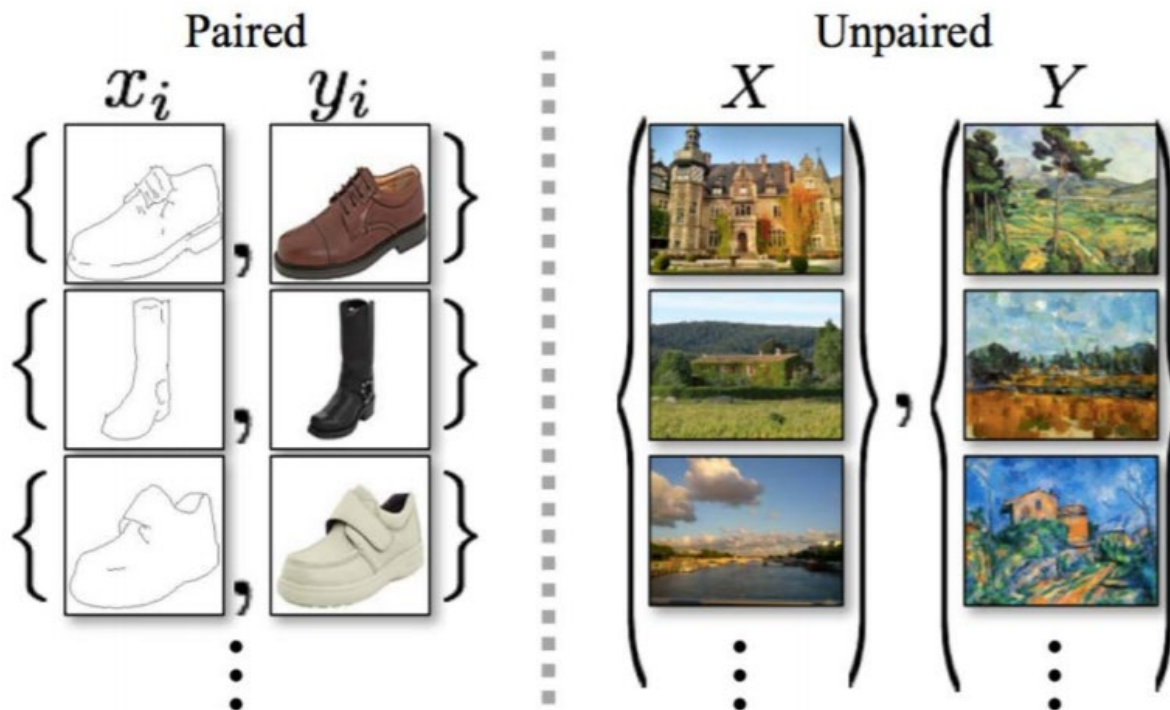
Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

Major problem with prior approach: we need **paired** samples.

(How do you “translate” a photograph to the style of Van Gogh?

No “paired” up photos...)

Can we solve “unpaired” version of problem?

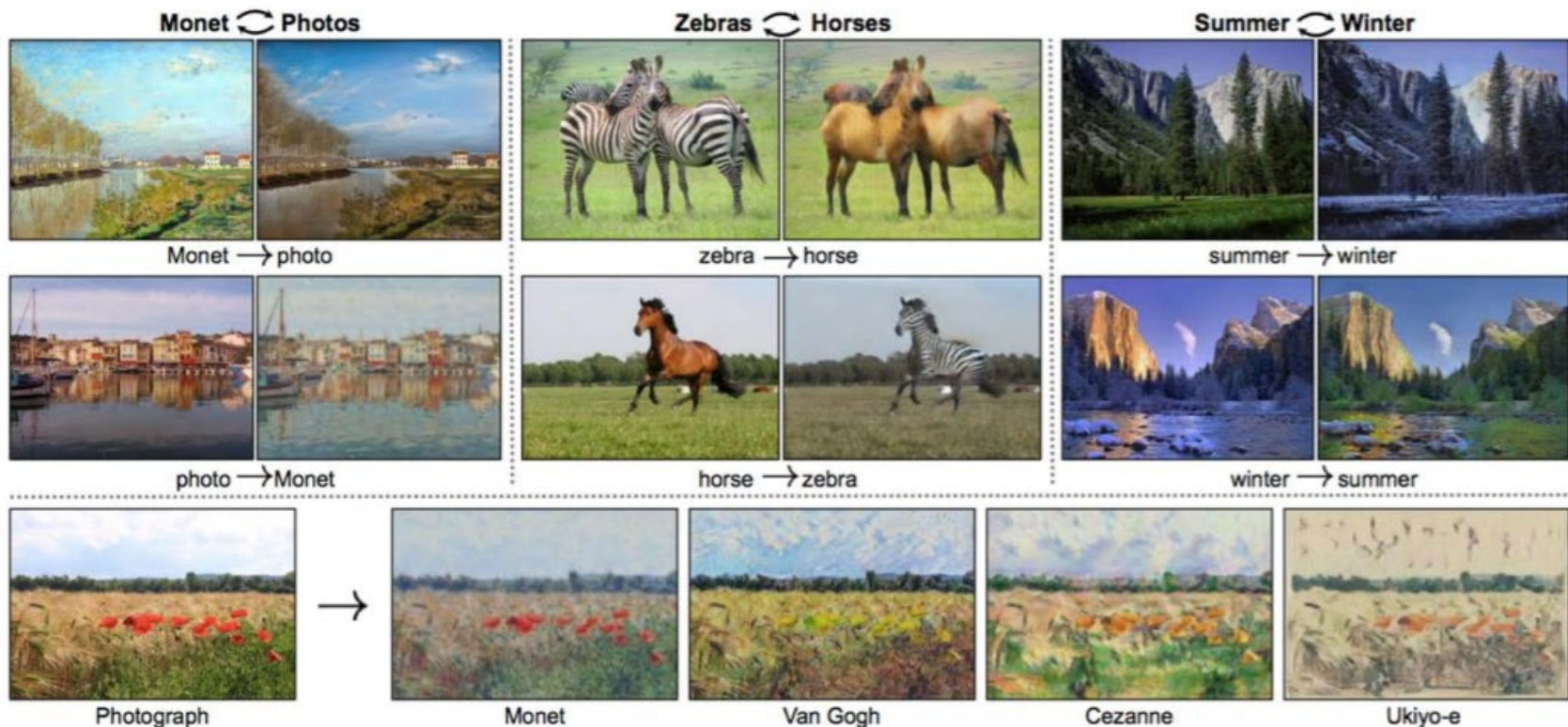


Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

Major problem with prior approach: we need **paired** samples.

(How do you “translate” a photograph to the style of Van Gogh?

No “paired” up photos...)



Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

Main idea:

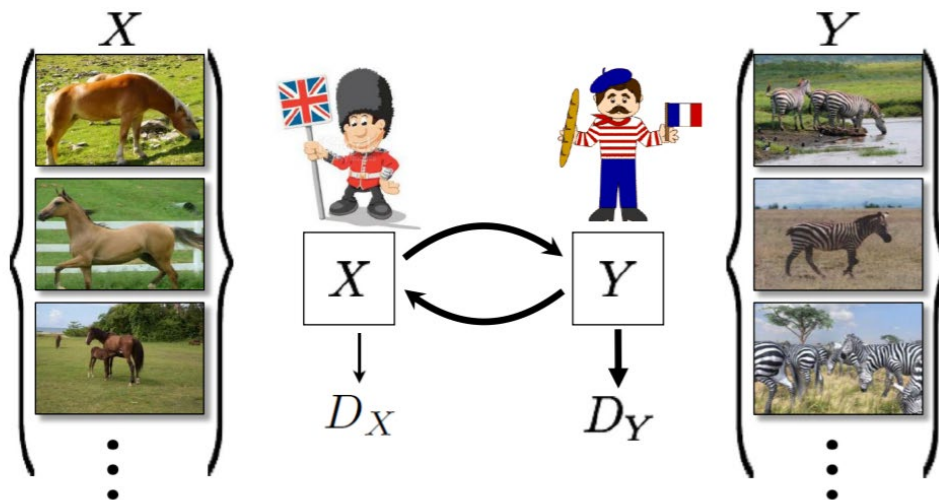
Train *two* generators F, G , s.t.

F translates from domain X to domain Y ,

G translates from domain Y to domain X .

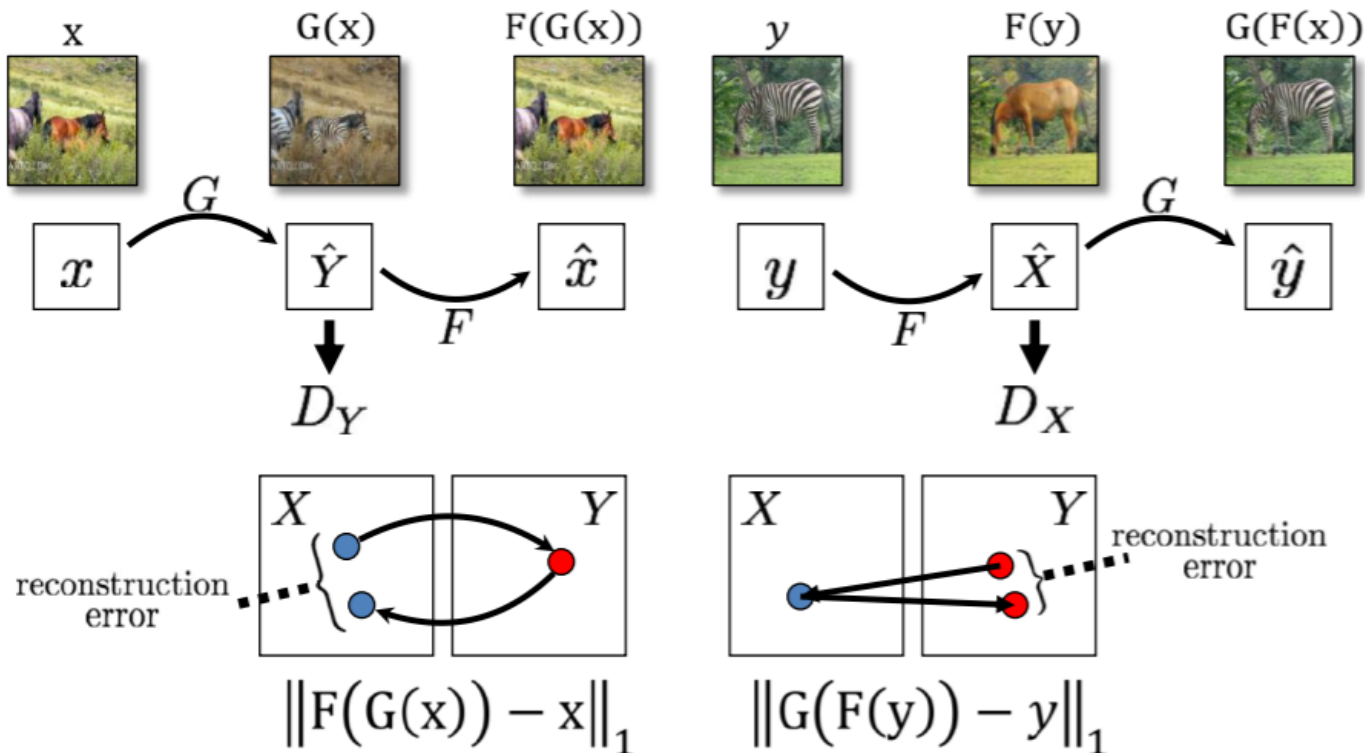
Discriminators D_X, D_Y trying to recognize domains X, Y .

Requirement: $F(G(x)) \approx x, G(F(X)) \approx x$



Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

The CycleGAN loss:



Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

The CycleGAN loss:

$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) &= \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] \\ &\quad + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))] \end{aligned} \quad \left| \quad \begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ &\quad + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \end{aligned}$$

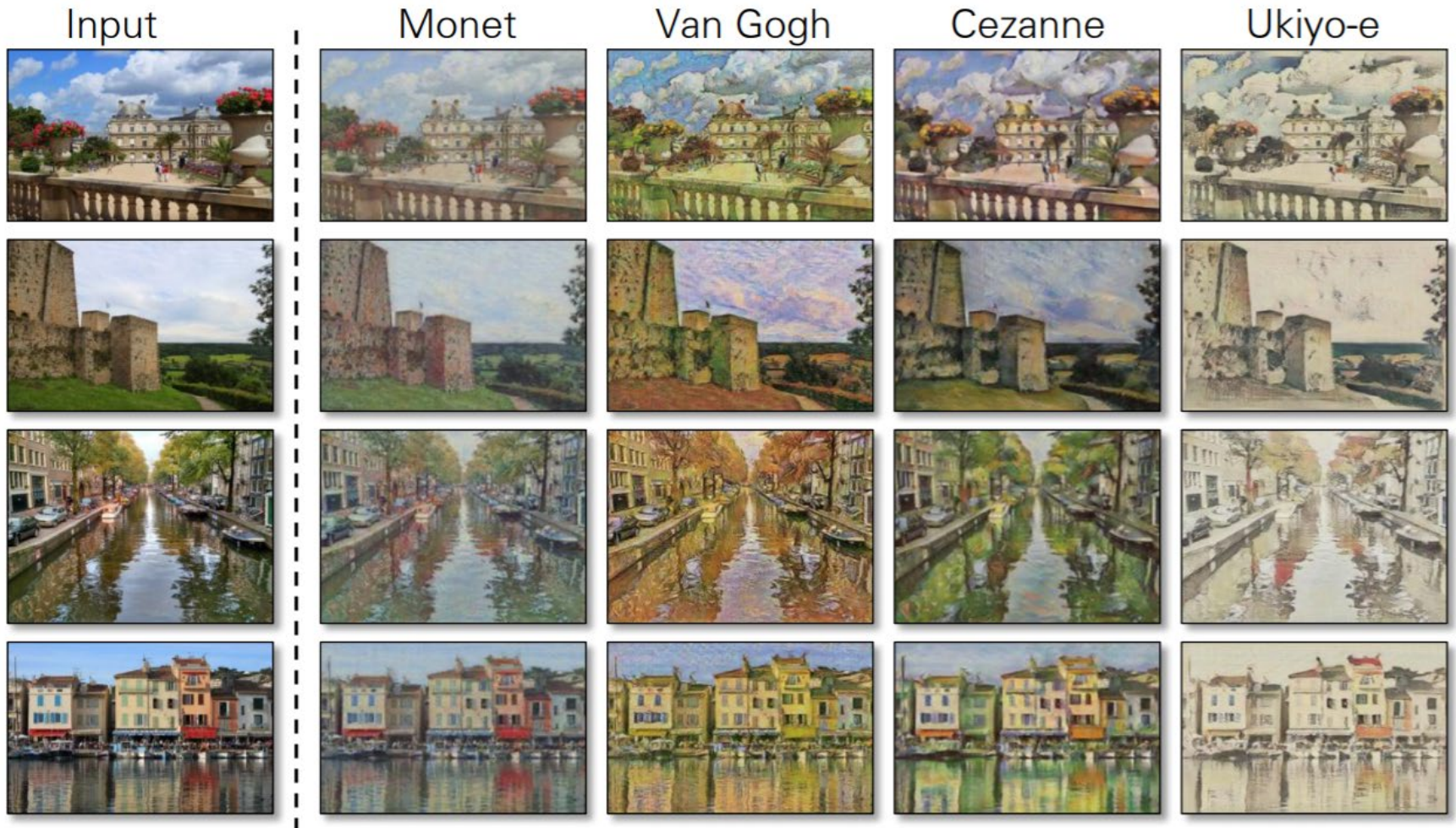
Putting them together, we get:

$$\begin{aligned} \mathcal{L}(G, F, D_X, D_Y) &= \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) \\ &\quad + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) \\ &\quad + \lambda \mathcal{L}_{\text{cyc}}(G, F), \end{aligned}$$

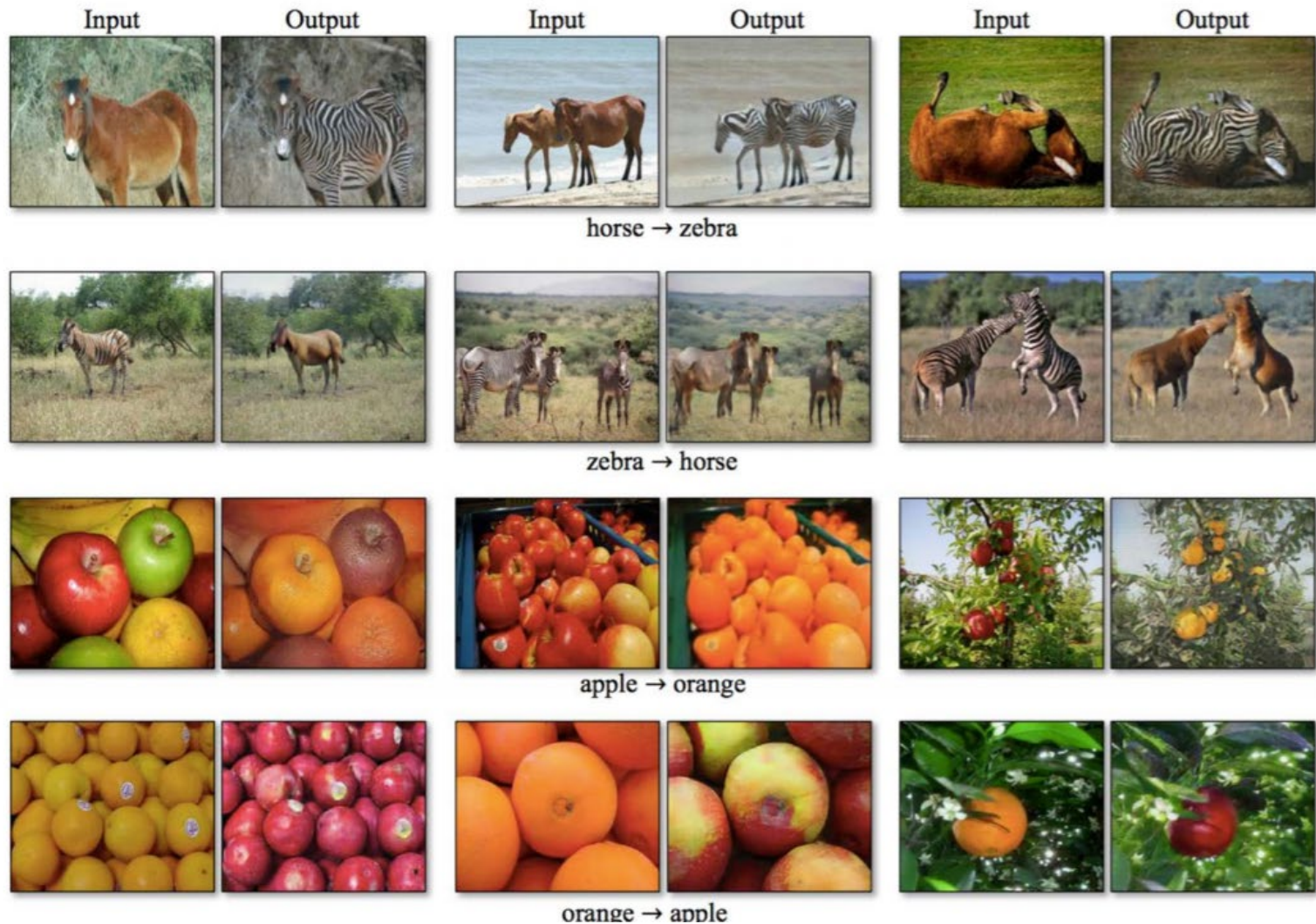
We are trying to optimize:

$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}(G, F, D_X, D_Y).$$

Unpaired image-to-image translation (CycleGAN, Zhu et al '17)



Unpaired image-to-image translation (CycleGAN, Zhu et al '17)



Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

Can definitely fail...

Input



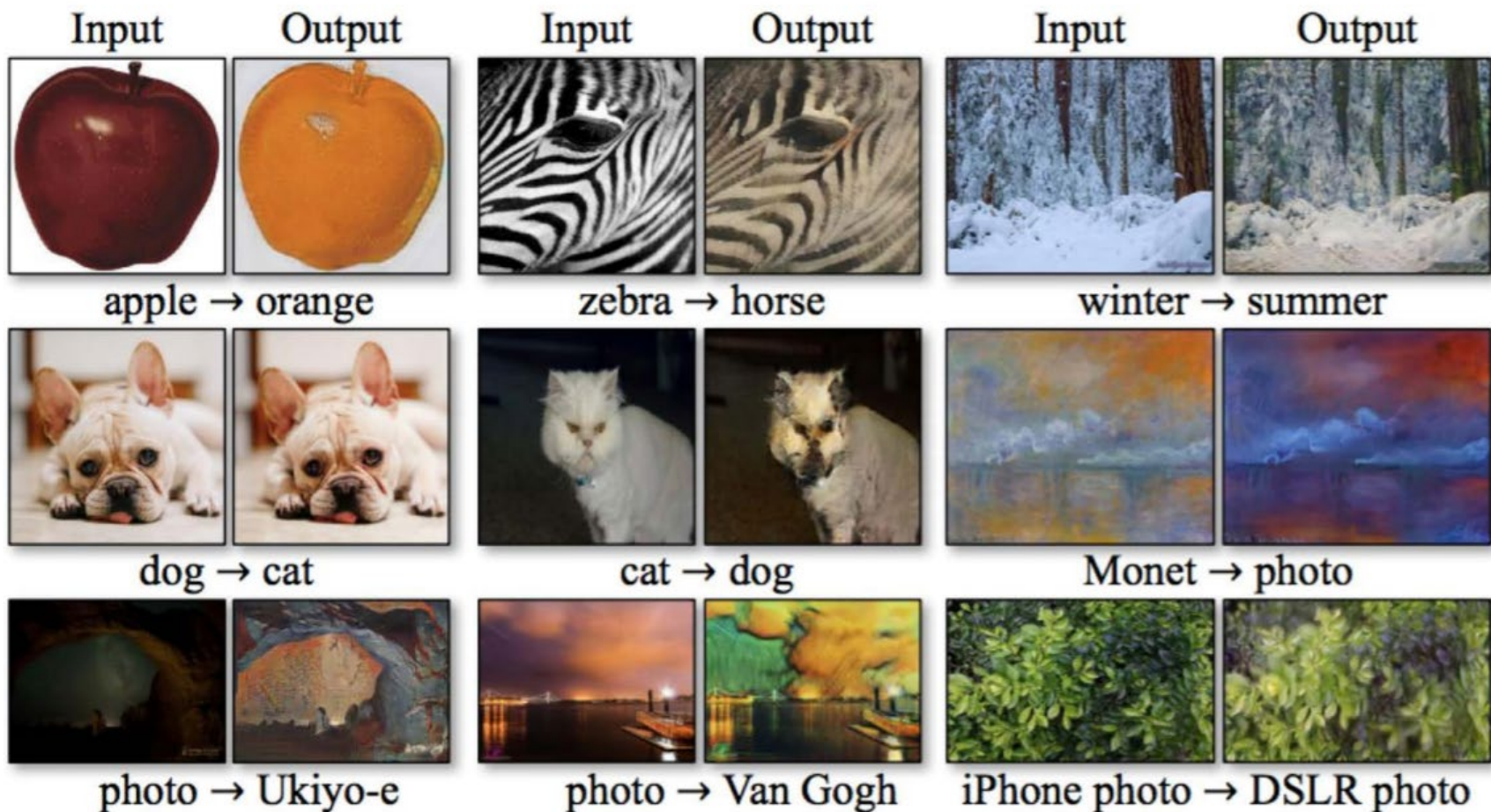
Output



horse → zebra

Unpaired image-to-image translation (CycleGAN, Zhu et al '17)

Can definitely fail...



Applications to other domains too

Brain lesion segmentation (Bowles et al 2018)

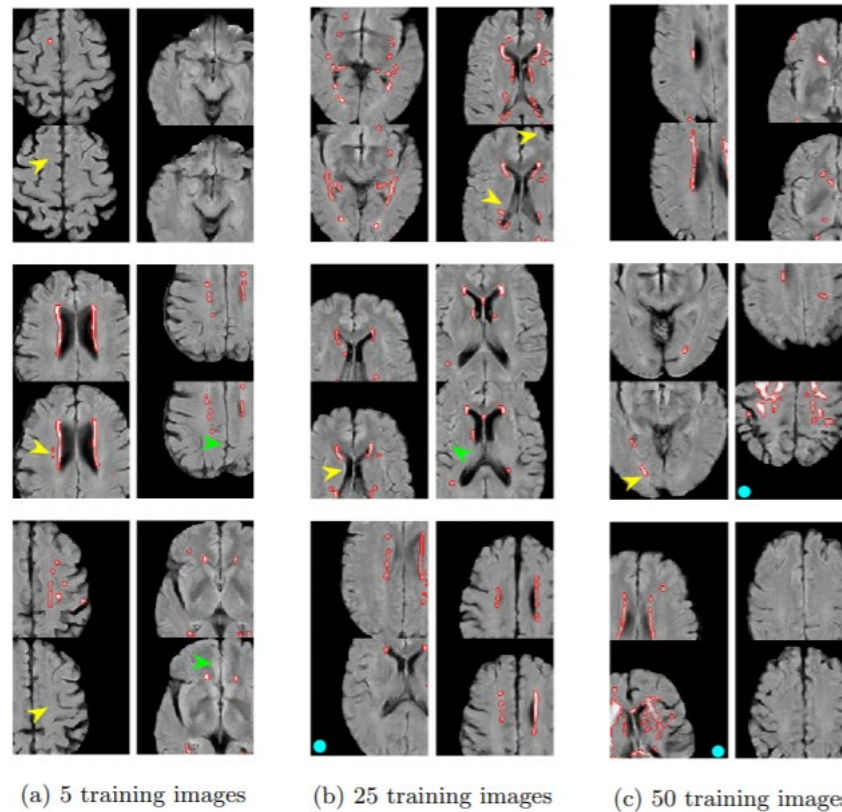
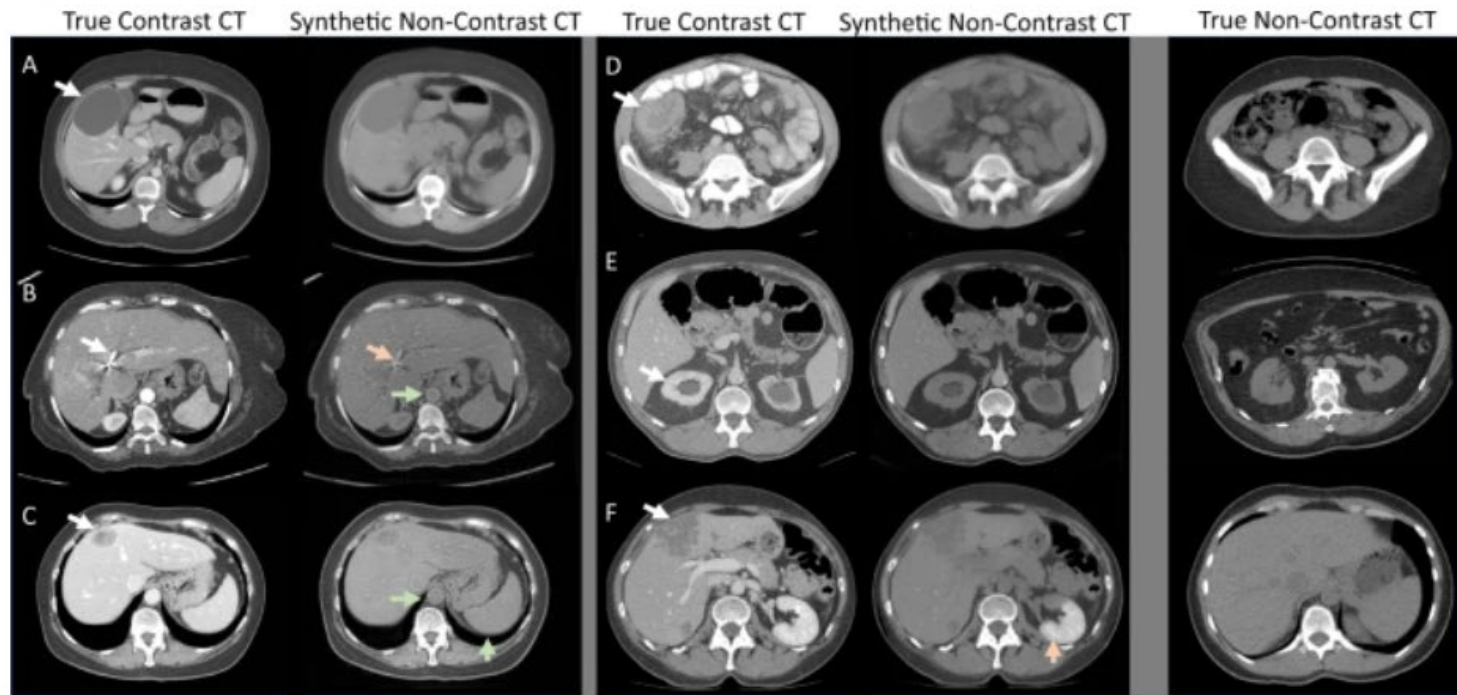


Fig. 3: Synthetic images (top of pair) with their nearest neighbours in the training set (bottom of pair) from GANs trained on patches from 5, 25 and 50 real MR images. Some local signs of successful augmentation are indicated using green (same lesions, different anatomy) and yellow (same anatomy, different lesions) arrows, and novel images (new anatomy and lesions) are shown with blue dots.

Applications to other domains too

CT scan segmentation (Sanfort et al, Nature 2019)

Figure 1



Examples of true IV contrast CT scans (left column) and synthetic non-contrast CT scans generated by a CycleGAN. The rightmost column shows unrelated example non-contrast images. Overall the synthetic non-contrast images appear convincing - even when significant abnormalities are present in the contrast CT scans.