

Image-to-image Translation

Given two domain the goal is to translate image from one possible representation to another.

$$\mathbf{x} \sim p(\mathbf{x}|\mathbf{y})$$

$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})$$

Paired image-to-image translation

$$\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})$$

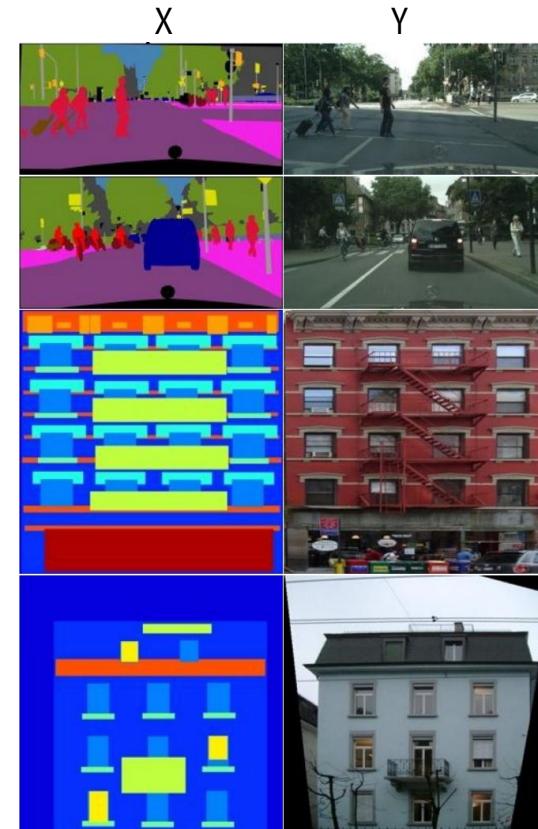


Image-to-image Translation

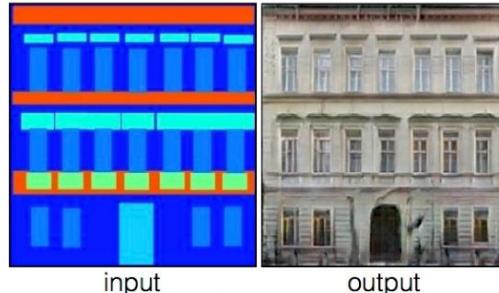
Labels to Street Scene



input

output

Labels to Facade



input

output

BW to Color



input

output

Aerial to Map



input

output

Day to Night



input

output

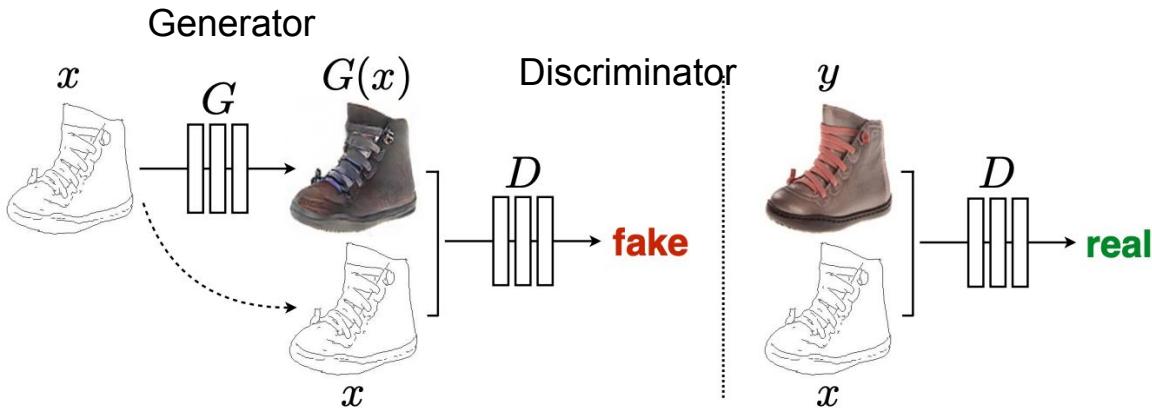
Edges to Photo



input

output

Image-to-image Translation: Pix2Pix



Combined GAN-loss and reconstruction loss:

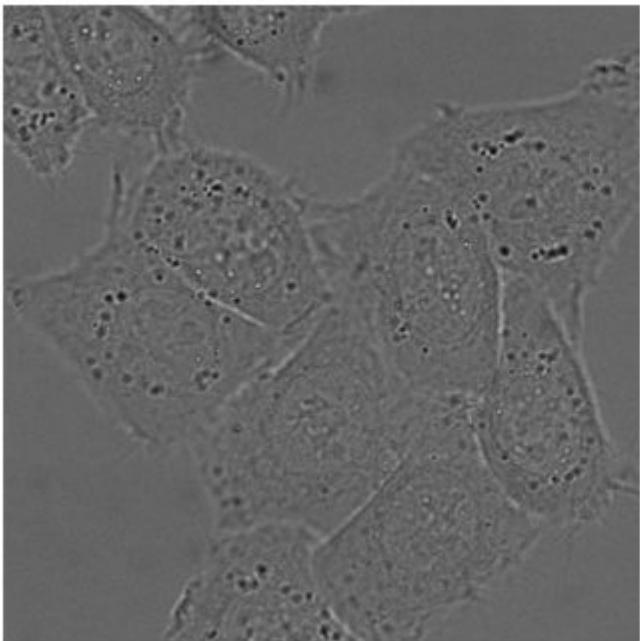
$$\begin{aligned}\mathcal{L}_{cGAN}(G, D) = & \mathbb{E}_{x,y}[\log D(x,y)] + \\ & \mathbb{E}_{x,z}[\log(1 - D(x, G(x, z)))]\end{aligned}$$

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1]$$

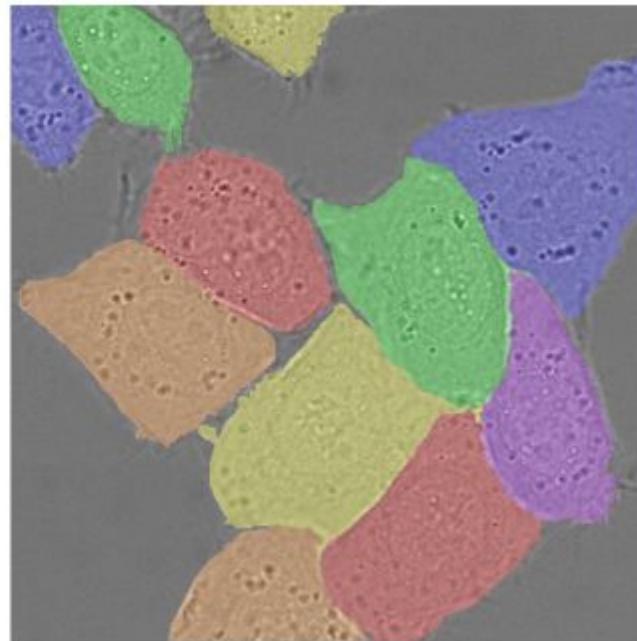
$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$

Image-to-image Translation: Pix2Pix

a



b

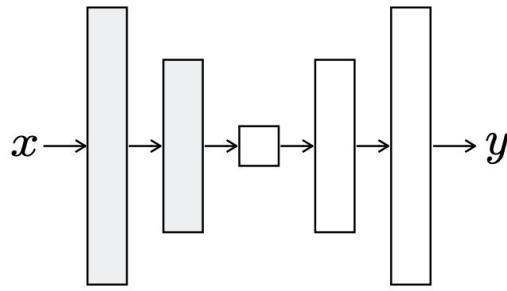


U-Net: Convolutional Networks for Biomedical Image Segmentation Olaf Ronneberger, Philipp Fischer, Thomas Brox MICCAI 2015

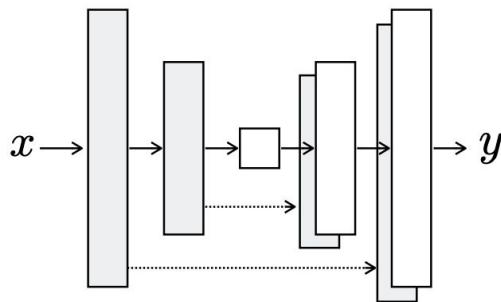
Image-to-image Translation: Pix2Pix

Skip connections in generator

Encoder-decoder



U-Net

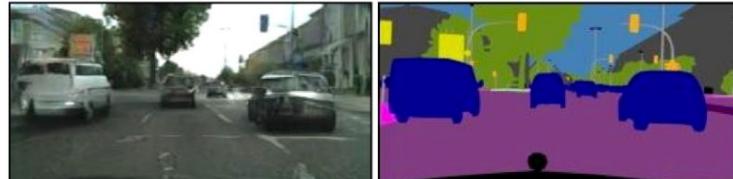


U-Net: Convolutional Networks for Biomedical Image Segmentation Olaf Ronneberger, Philipp Fischer, Thomas Brox MICCAI 2015

Pix2Pix: Ablations

Loss	Per-pixel acc.	Per-class acc.	Class IOU
Encoder-decoder (L1)	0.35	0.12	0.08
Encoder-decoder (L1+cGAN)	0.29	0.09	0.05
U-net (L1)	0.48	0.18	0.13
U-net (L1+cGAN)	0.55	0.20	0.14

Table 2: FCN-scores for different generator architectures (and objectives), evaluated on Cityscapes labels↔photos. (U-net (L1-cGAN) scores differ from those reported in other tables since batch size was 10 for this experiment and 1 for other tables, and random variation between training runs.)



Pix2Pix: Ablations

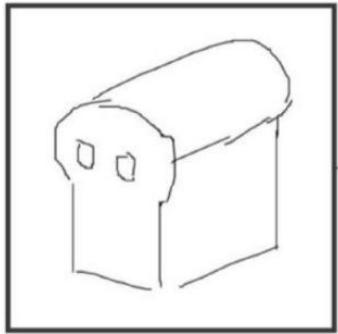
Discriminator receptive field	Per-pixel acc.	Per-class acc.	Class IOU
1×1	0.39	0.15	0.10
16×16	0.65	0.21	0.17
70×70	0.66	0.23	0.17
286×286	0.42	0.16	0.11

Table 3: FCN-scores for different receptive field sizes of the discriminator, evaluated on Cityscapes labels→photos. Note that input images are 256×256 pixels and larger receptive fields are padded with zeros.



Pix2Pix: Applications

#edges2cats by Christopher Hesse



pix2pix
process



sketch by Ivy Tsai

Background removal



by Kaihu Chen

Palette generation



by Jack Qiao

Sketch → Pokemon



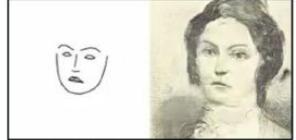
by Bertrand Gondouin

“Do as I do”



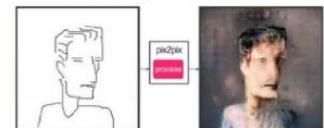
by Brannon Dorsey

Sketch → Portrait



by Mario Klingemann

#fotogenerator



sketch by Yann LeCun

Unpaired Image-to-image Translation

Given two domain the goal is to translate image from one possible representation to another.

$$\mathbf{x} \sim p(\mathbf{x}|\mathbf{y})$$

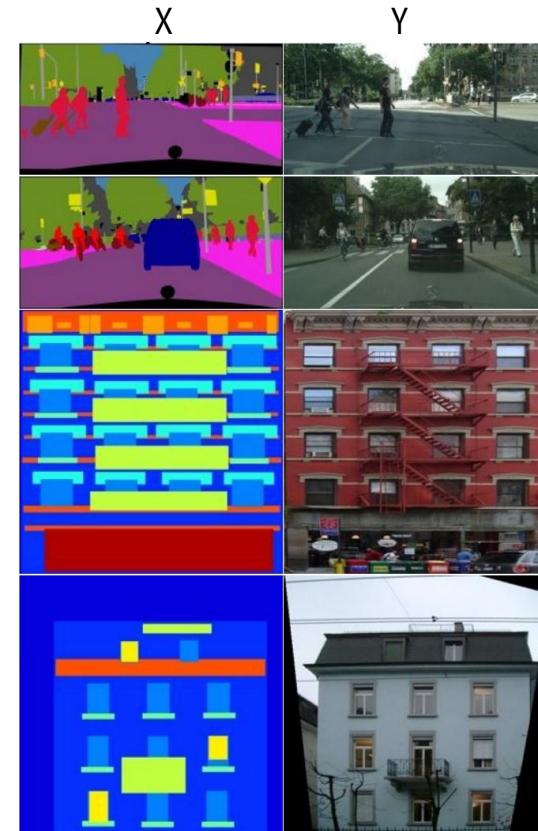
$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})$$

Paired image-to-image translation

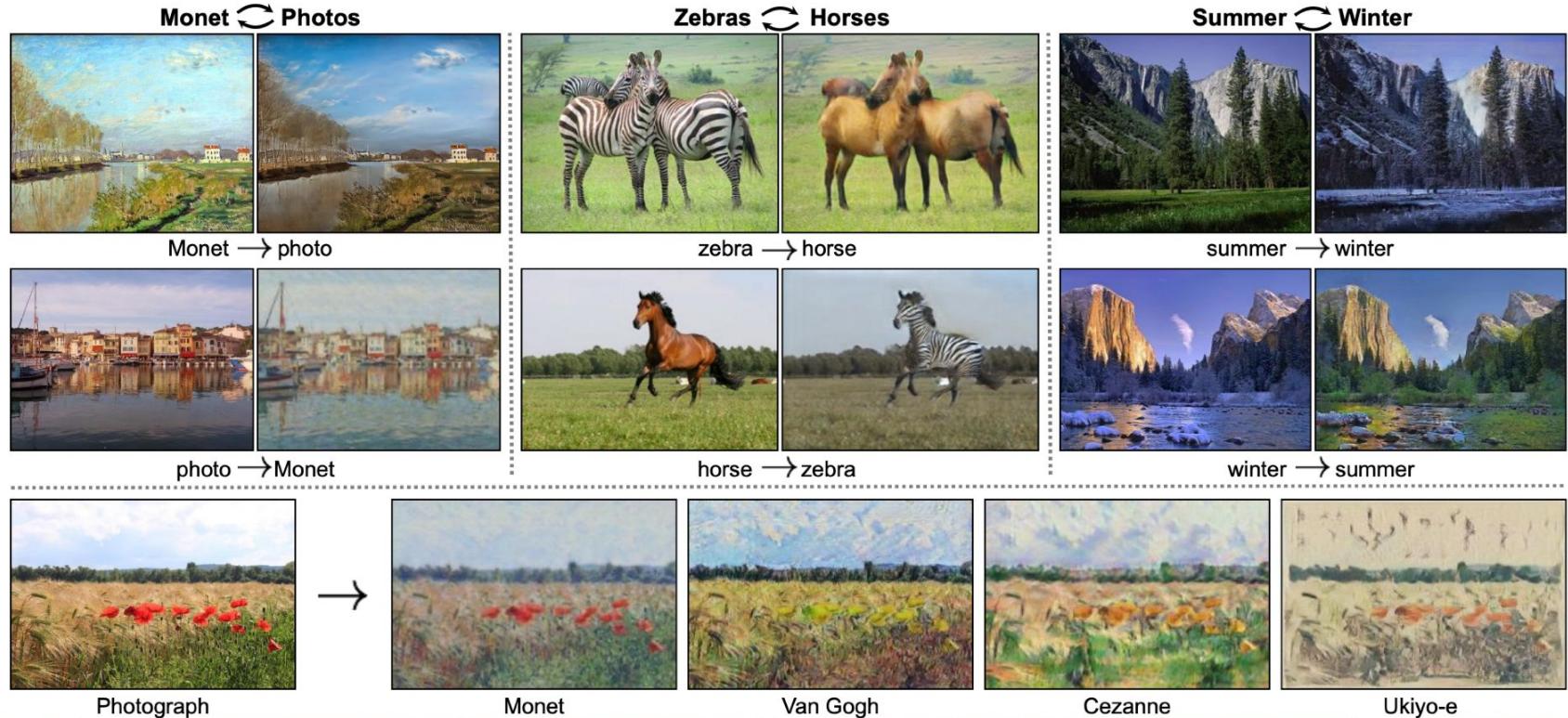
$$\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})$$

Unpaired

$$\mathbf{x} \sim p(\mathbf{x}), \mathbf{y} \sim p(\mathbf{y})$$



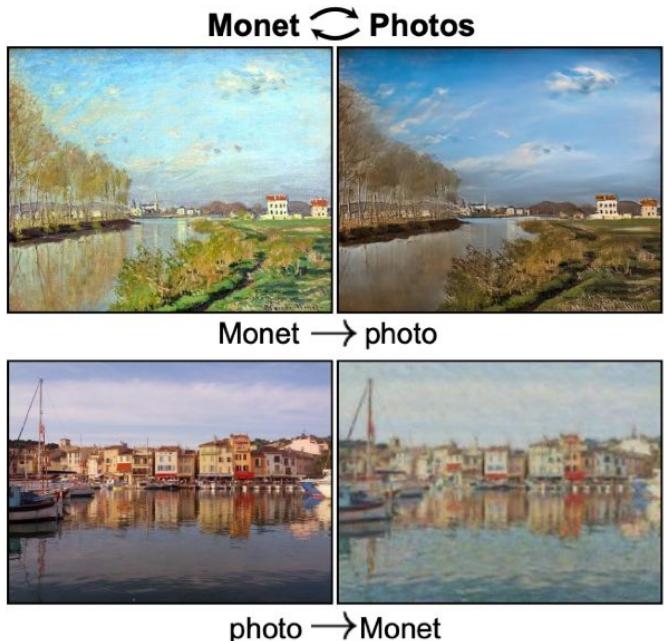
Unpaired Image-to-image Translation



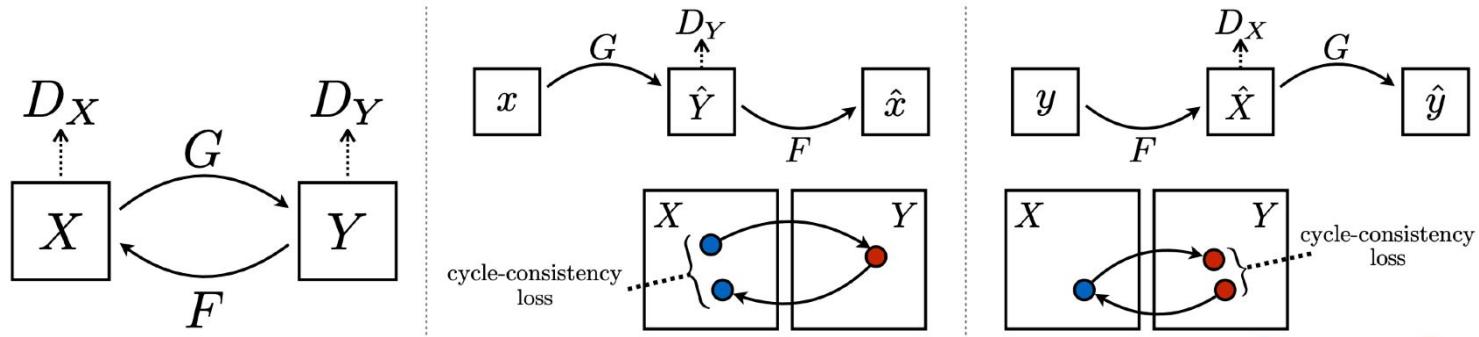
Unpaired Image-to-image Translation

To solve it, constraints are necessary:

- Cycle-consistency constraint
- Weight-sharing constraint
- Geometry-consistency constraint



CycleGAN



Adversarial loss:

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

Cycle-consistency loss:

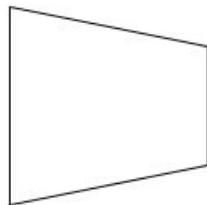
$$\begin{aligned}\mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]\end{aligned}$$

Full objective:

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

CycleGAN

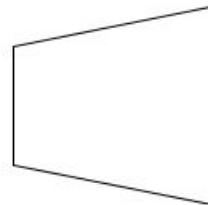
Downsampling



Strided conv

Batch-norm

Upsampling



ResBlocks

Conv

Batch-norm

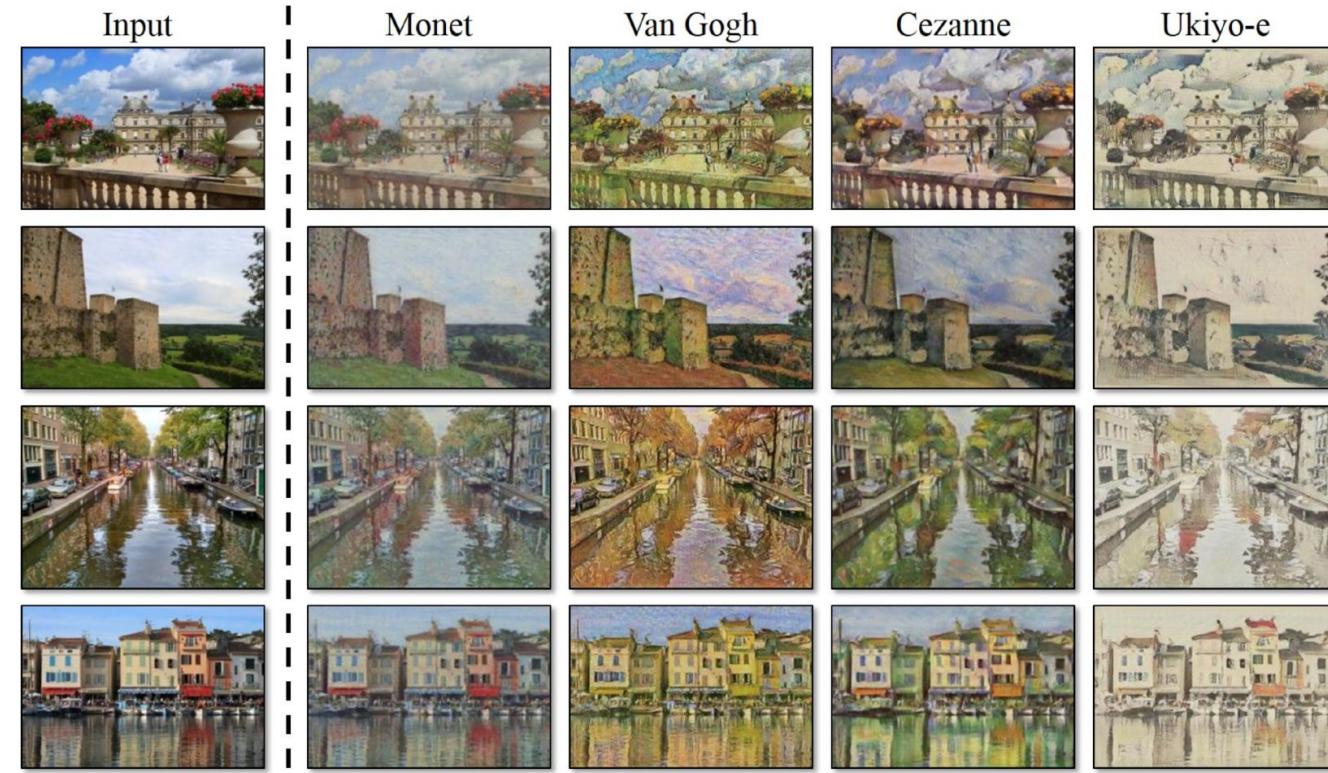
ReLU

Conv

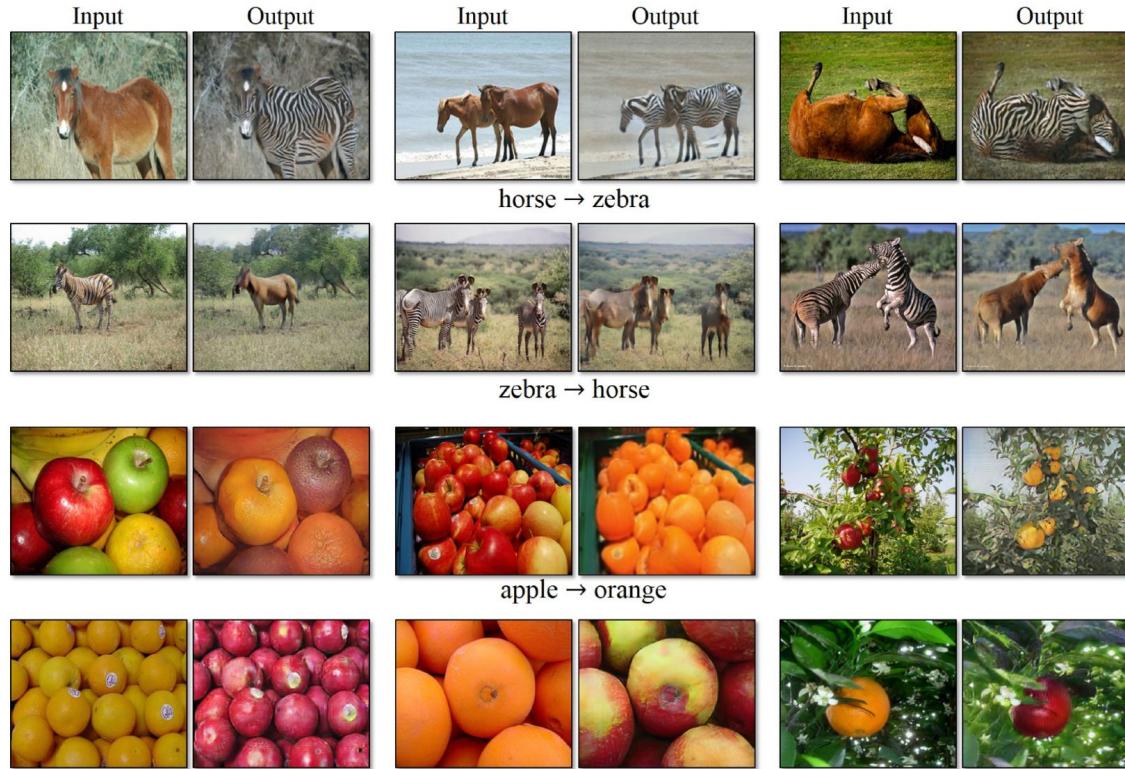
Batch-norm

add input

CycleGAN



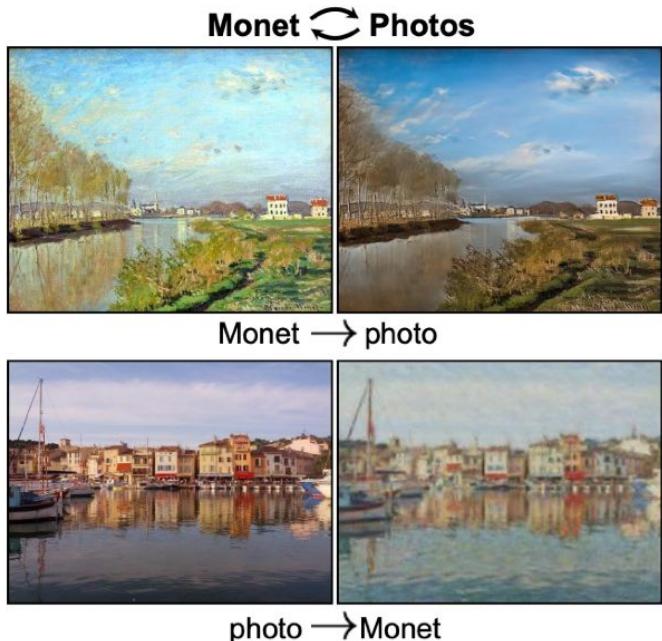
CycleGAN



Unpaired Image-to-image Translation

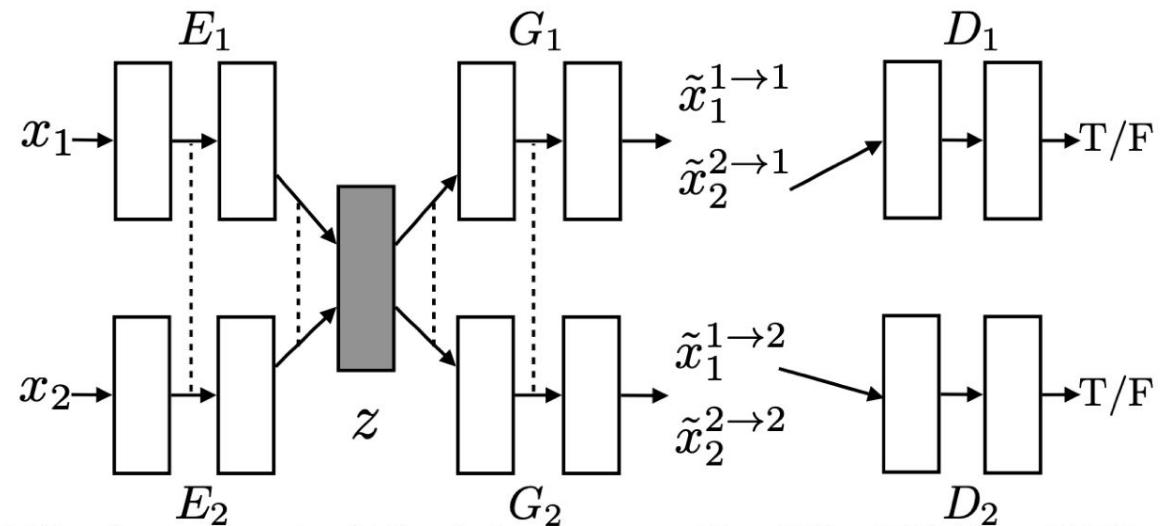
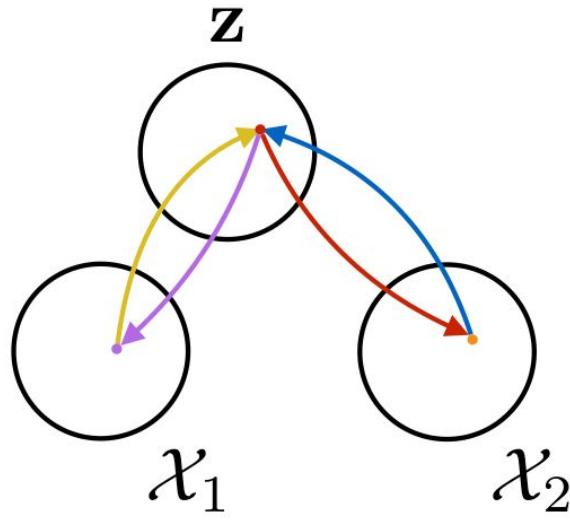
To solve it, constraints are necessary:

- Cycle-consistency constraint
- **Weight-sharing constraint**
- Geometry-consistency constraint



Unpaired Translation with weight sharing

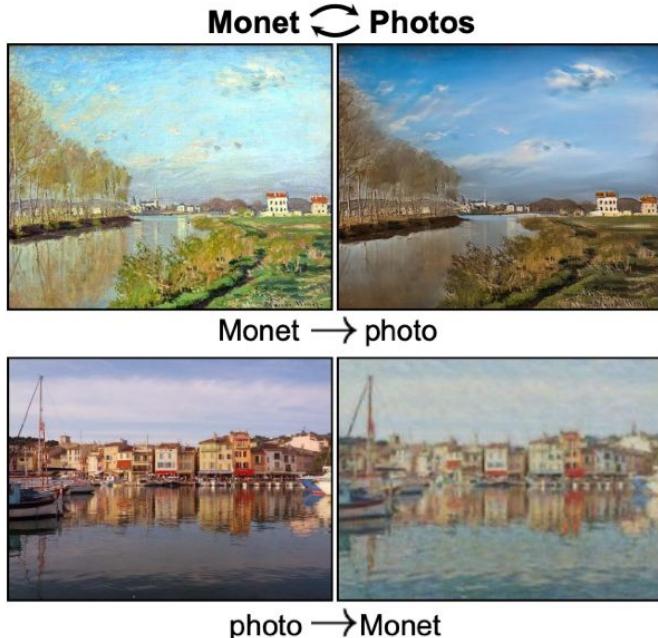
Train encoders with shared weights to go from image space to shared latent space



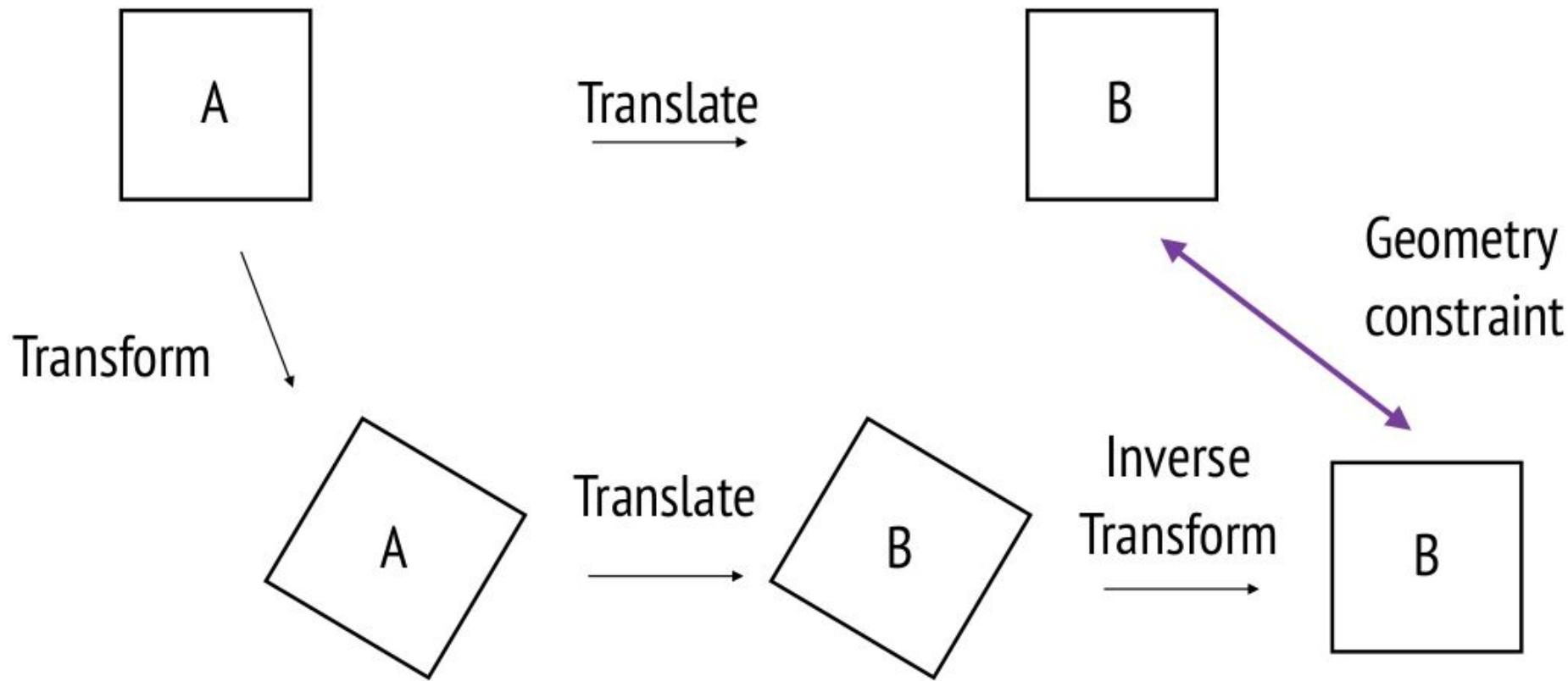
Unpaired Image-to-image Translation

To solve it, constraints are necessary:

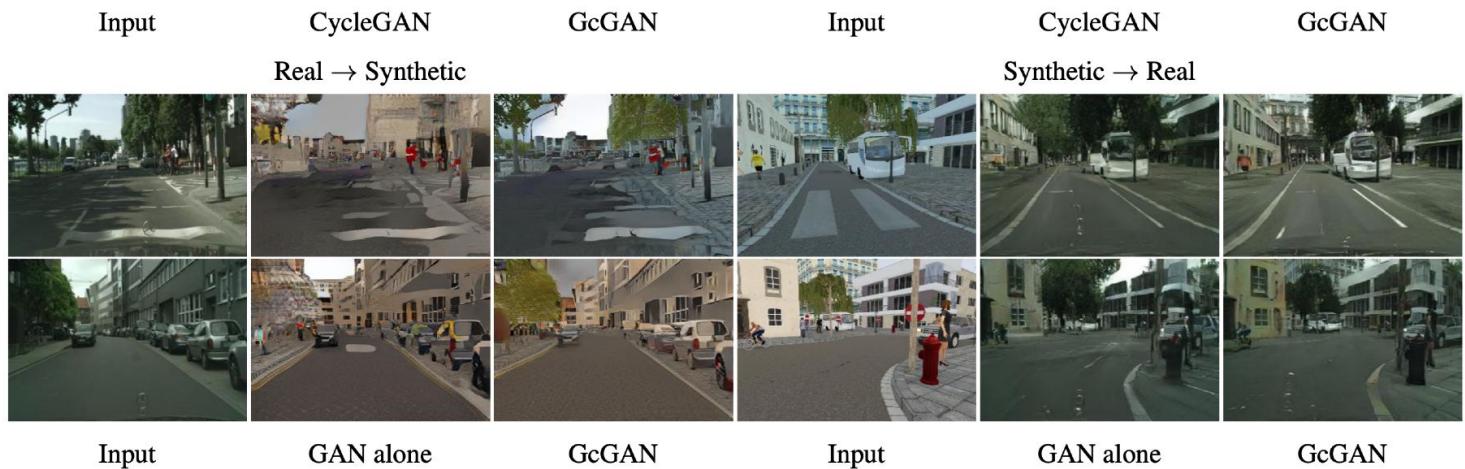
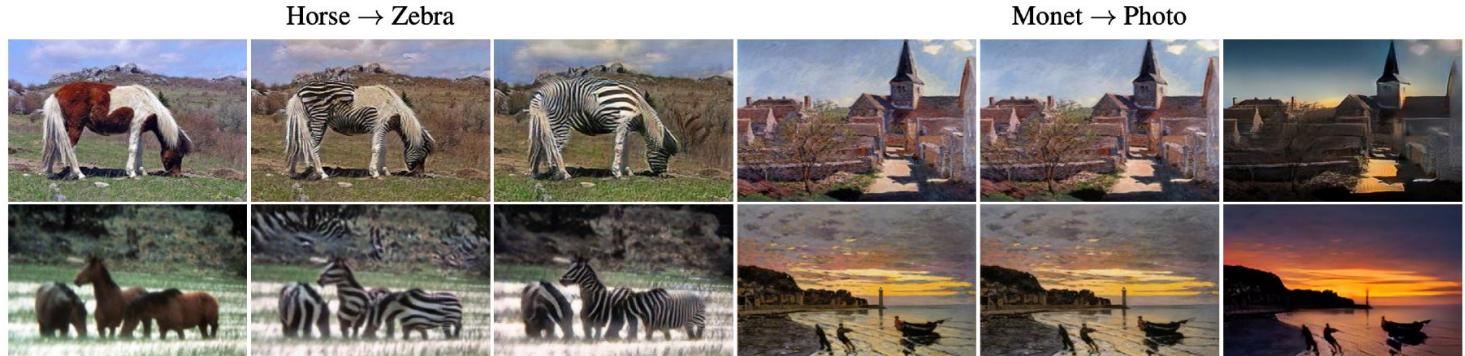
- Cycle-consistency constraint
- Weight-sharing constraint
- **Geometry-consistency constraint**



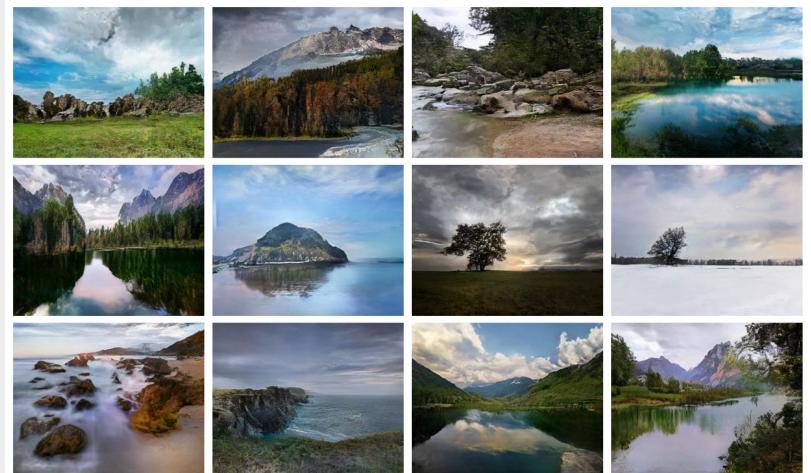
Geometry-consistency Constraint



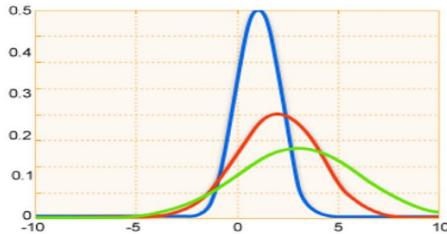
Geometry-consistency Constraint



Advanced methods



Style transfer

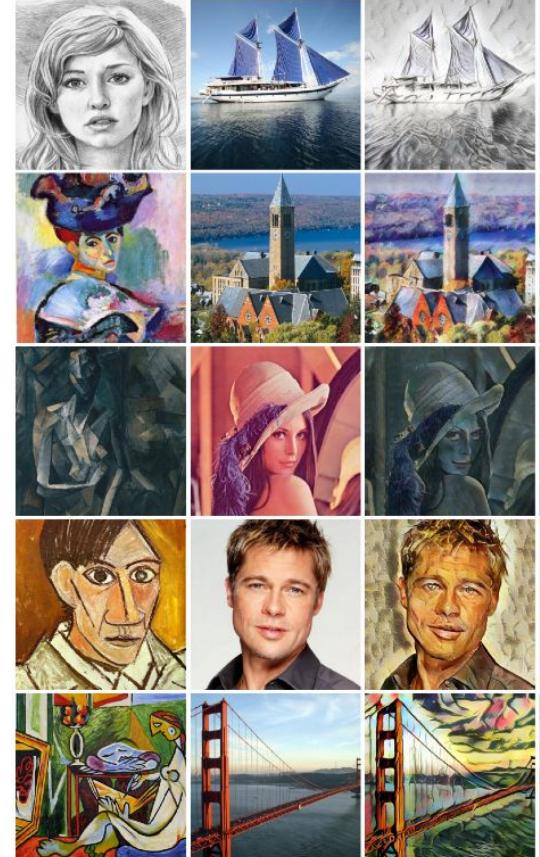


(a) Content image.

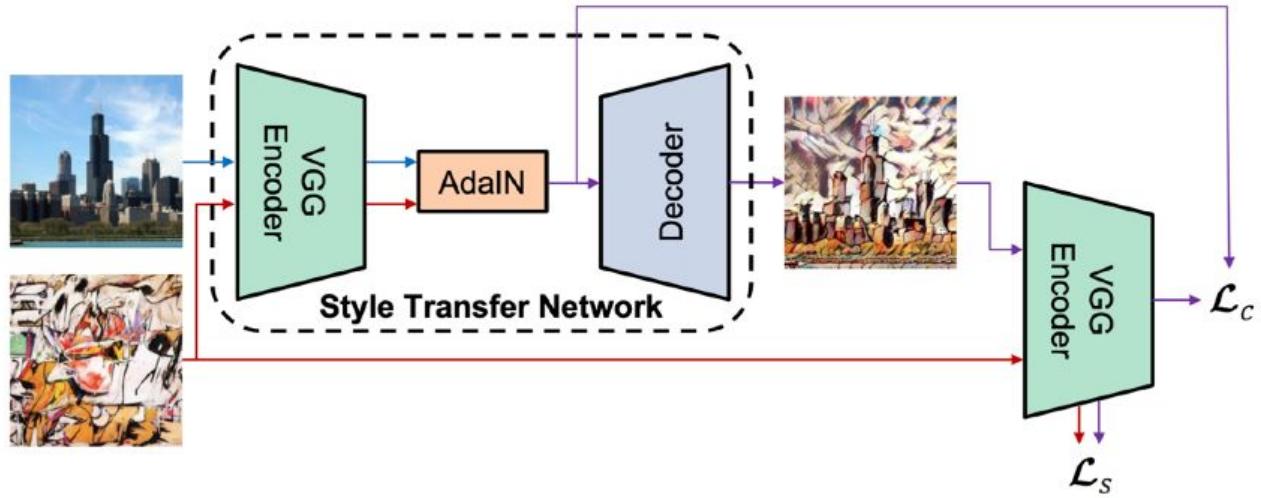


(c) Low contrast content image.

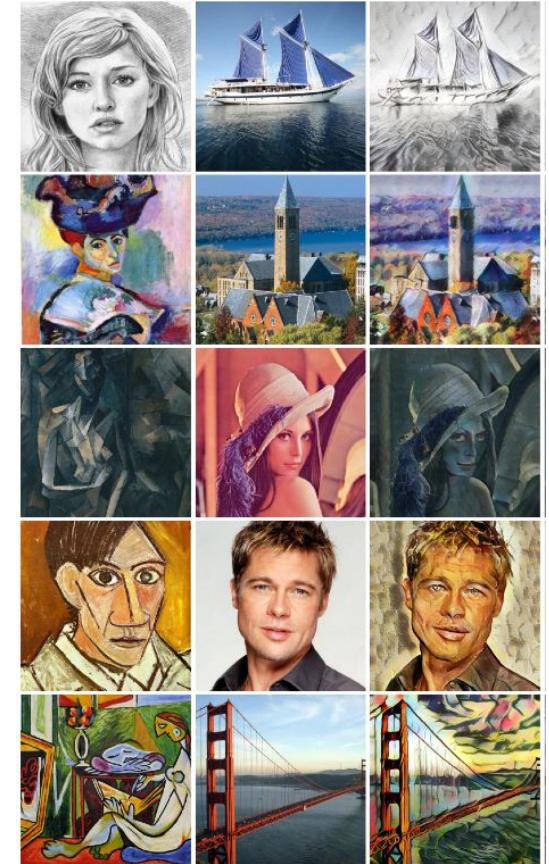
The two images are identical after instance normalization.
We obtain style invariant representations.



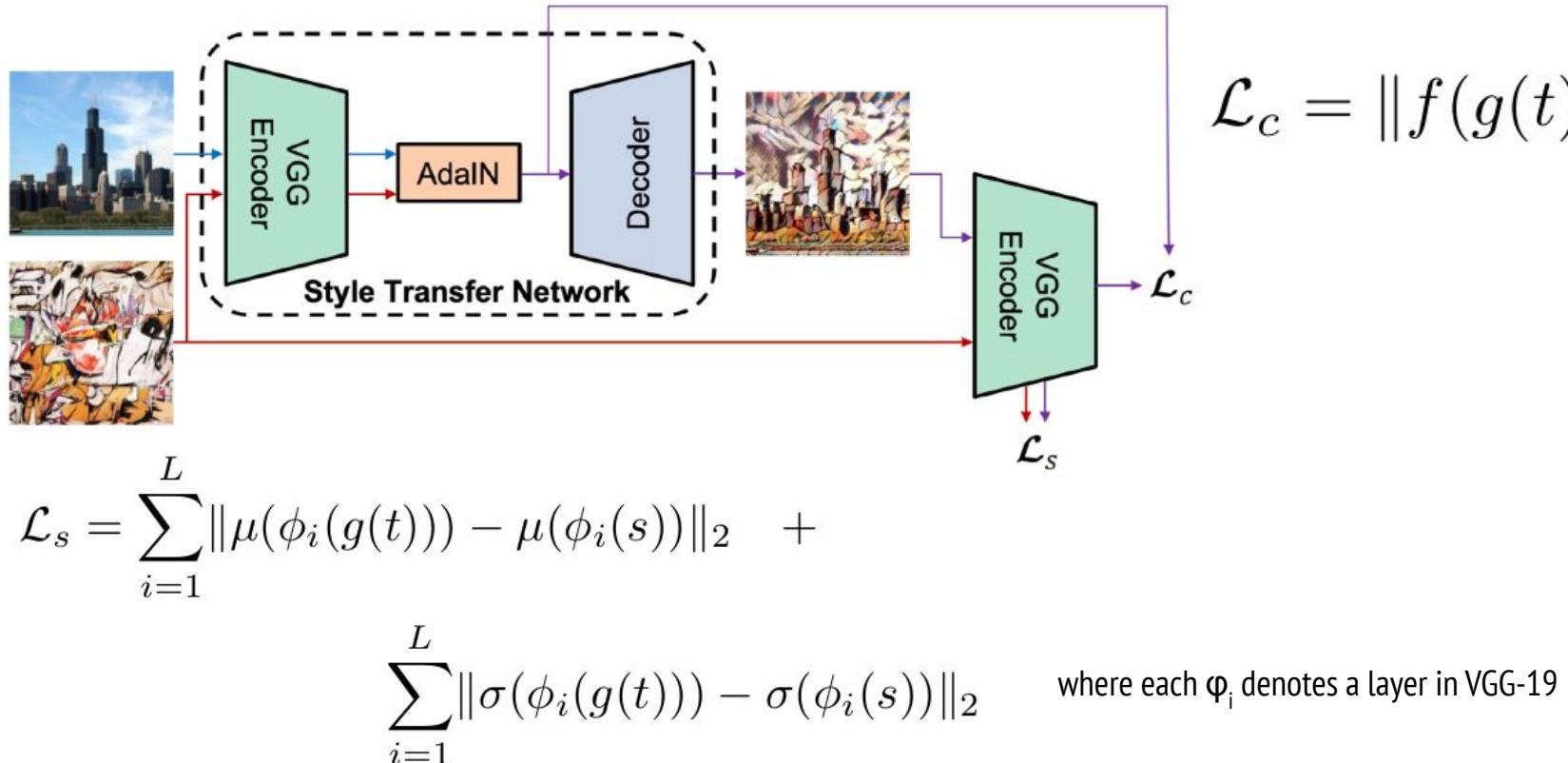
Style transfer



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

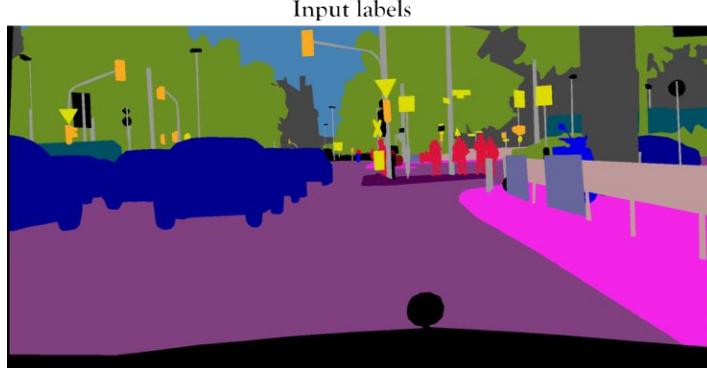
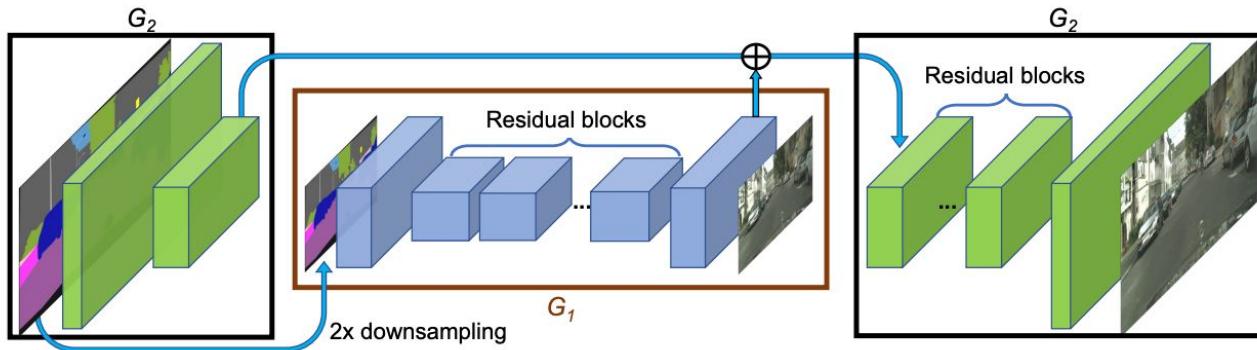


Style transfer

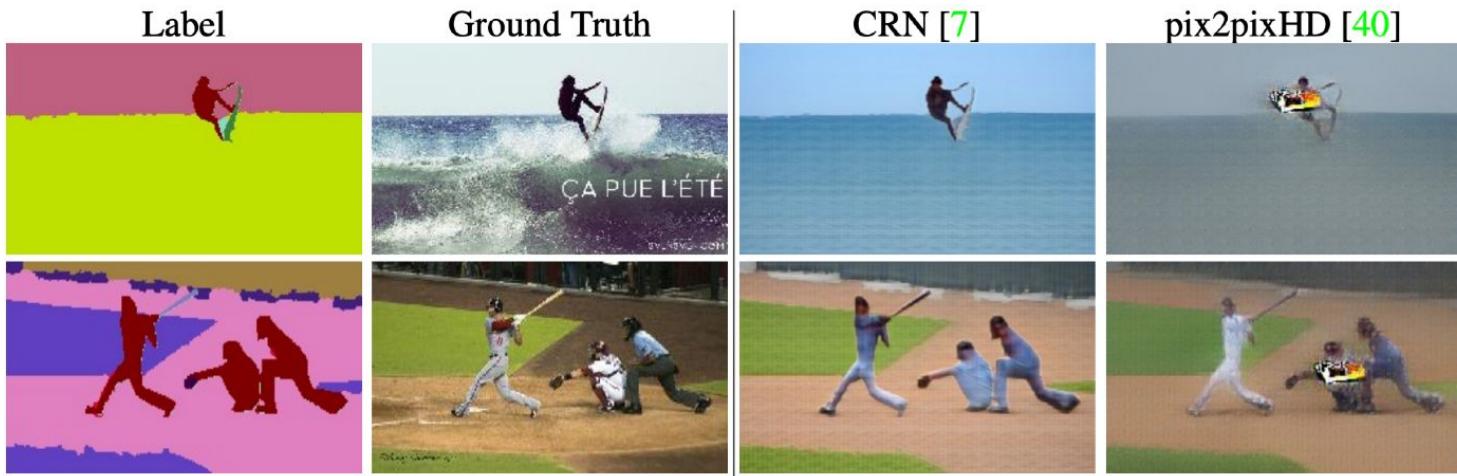


Pix2PixHD

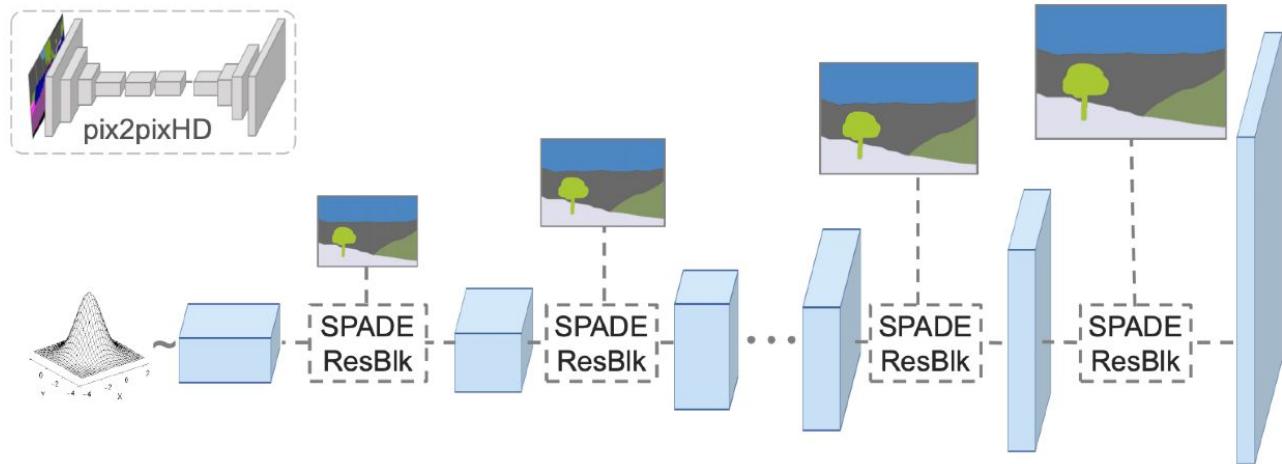
Idea: two-stage coarse-to-fine generation of HD images



Pix2PixHD

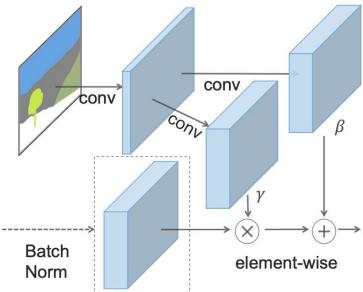
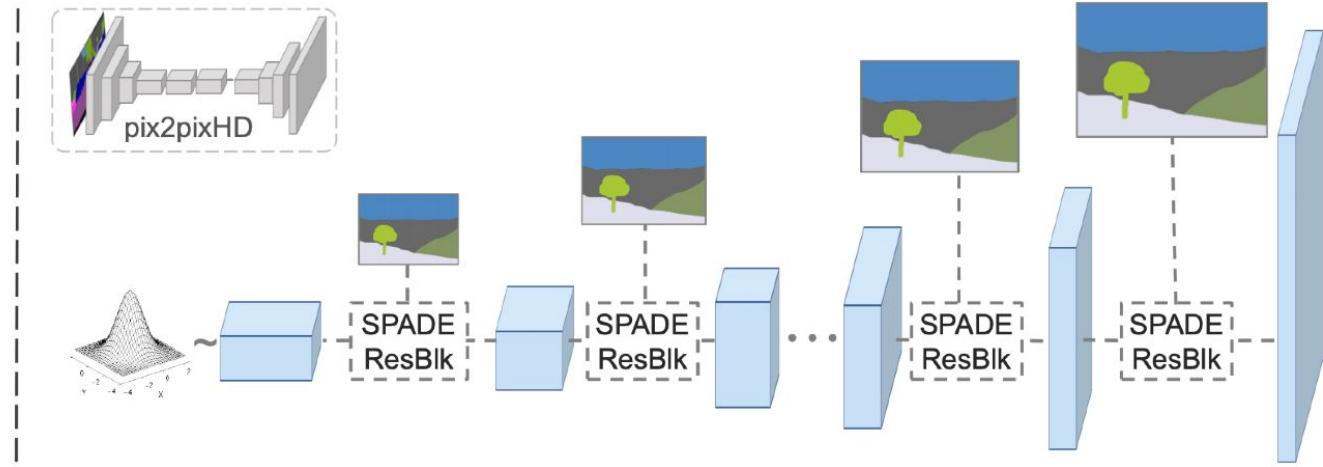
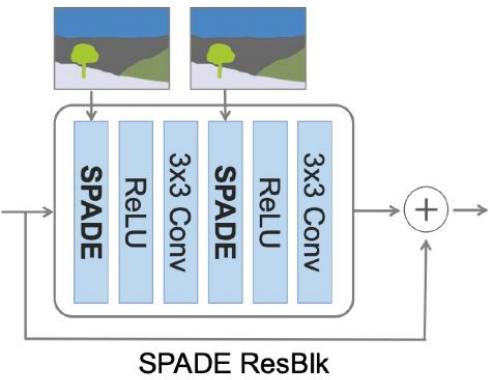


Spade



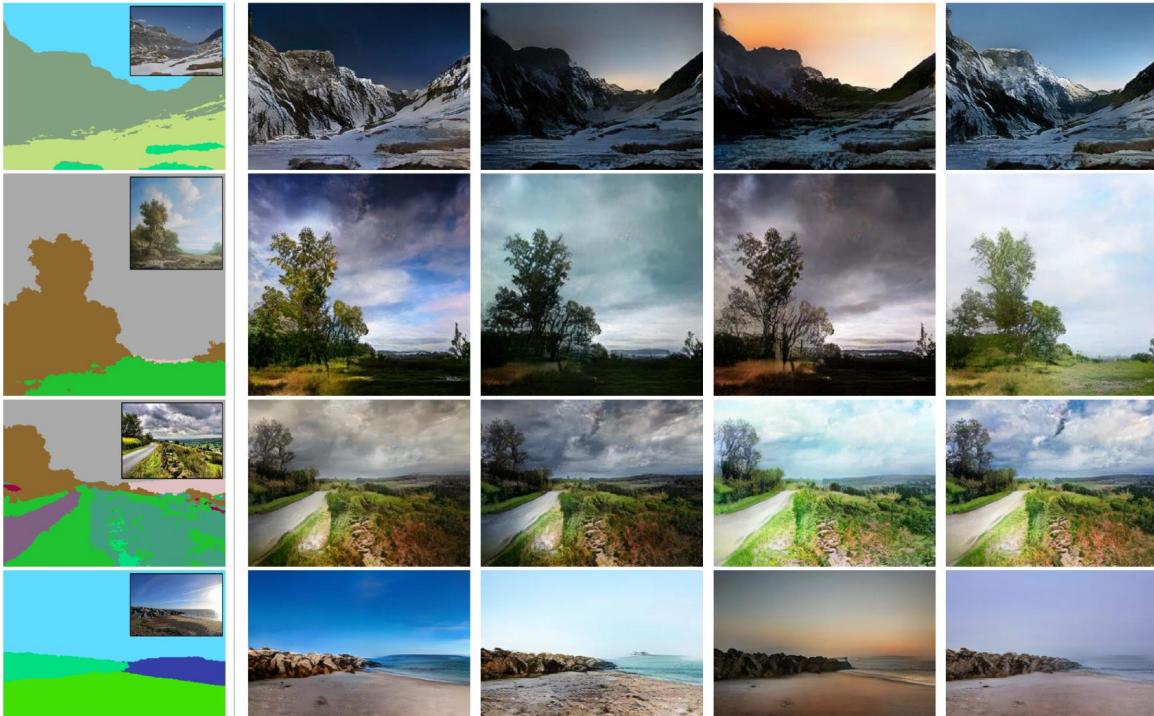
Park, Taesung, et al. "Semantic image synthesis with spatially-adaptive normalization." CVPR'2019

Spade



Park, Taesung, et al. "Semantic image synthesis with spatially-adaptive normalization." CVPR'2019

Spade

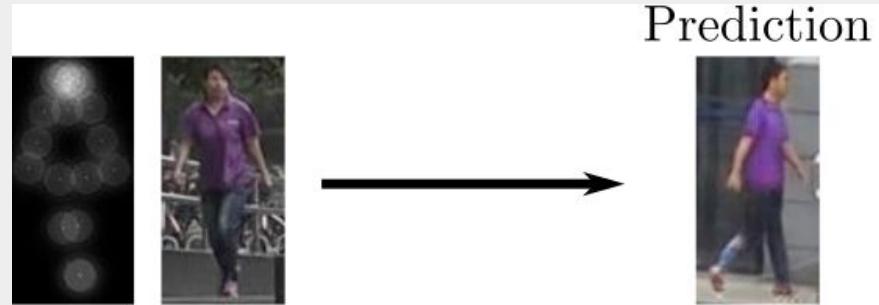


[Live demo](#)

Park, Taesung, et al. "Semantic image synthesis with spatially-adaptive normalization." CVPR'2019

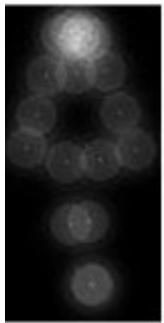
Pose-guided Image Generation

Stéphane Lathuilière

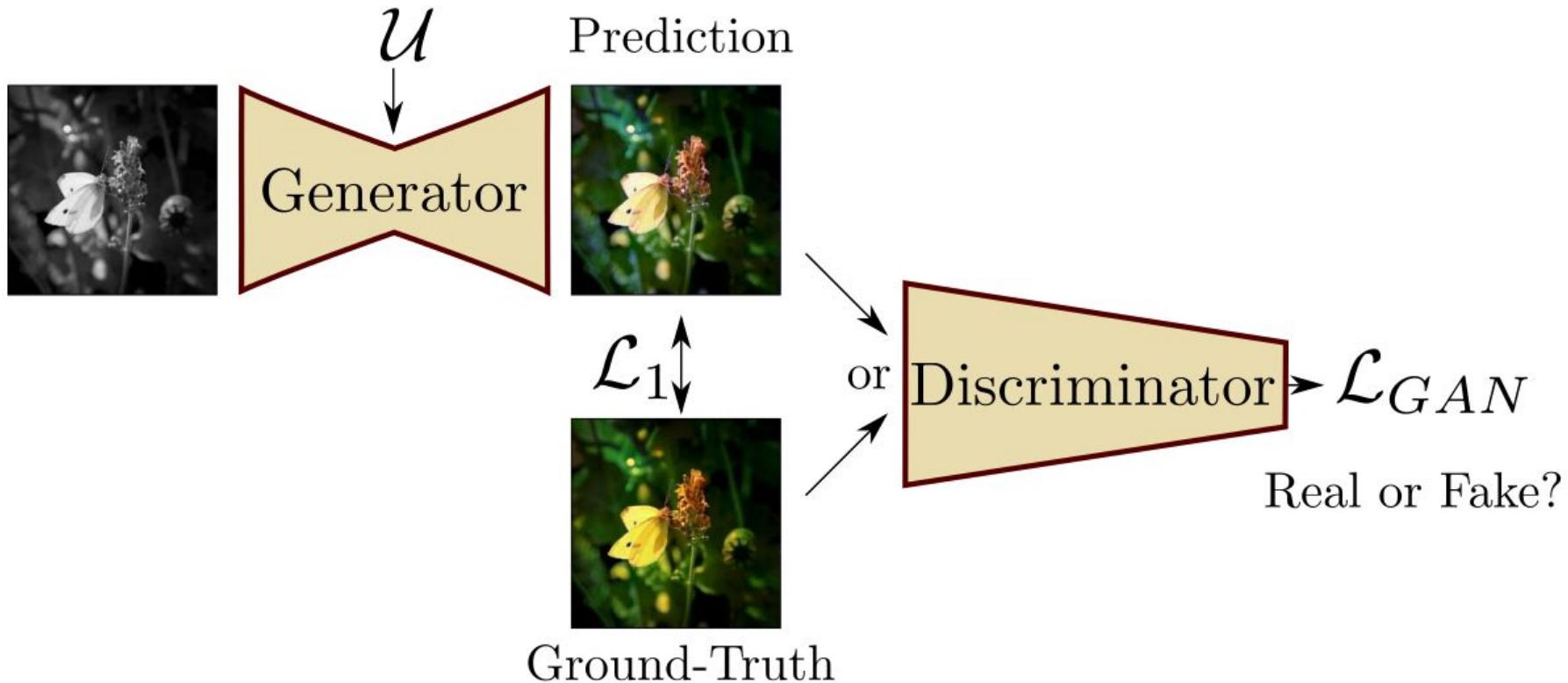


Introduction

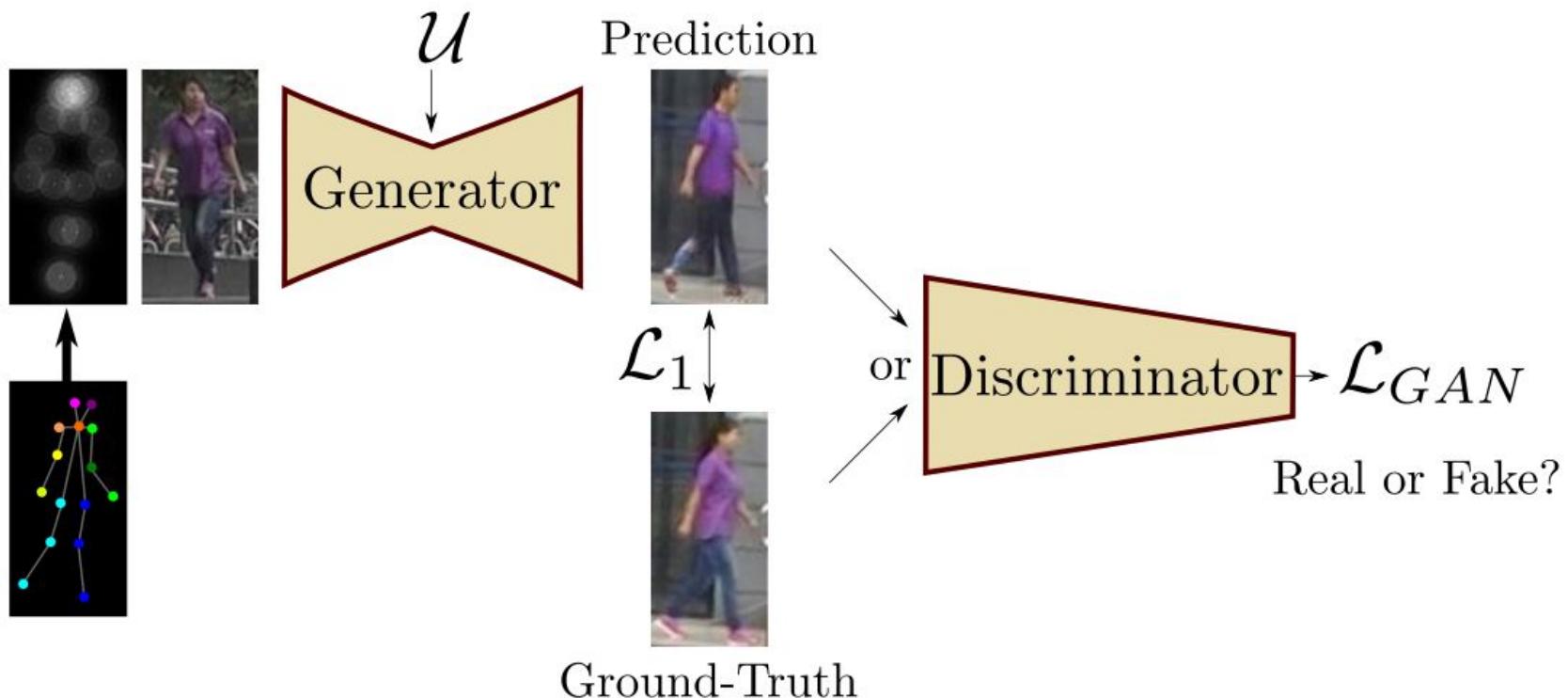
Prediction



Pix2Pix



Pix2Pix for pose guided



Pix2Pix for pose guided



Pose guided

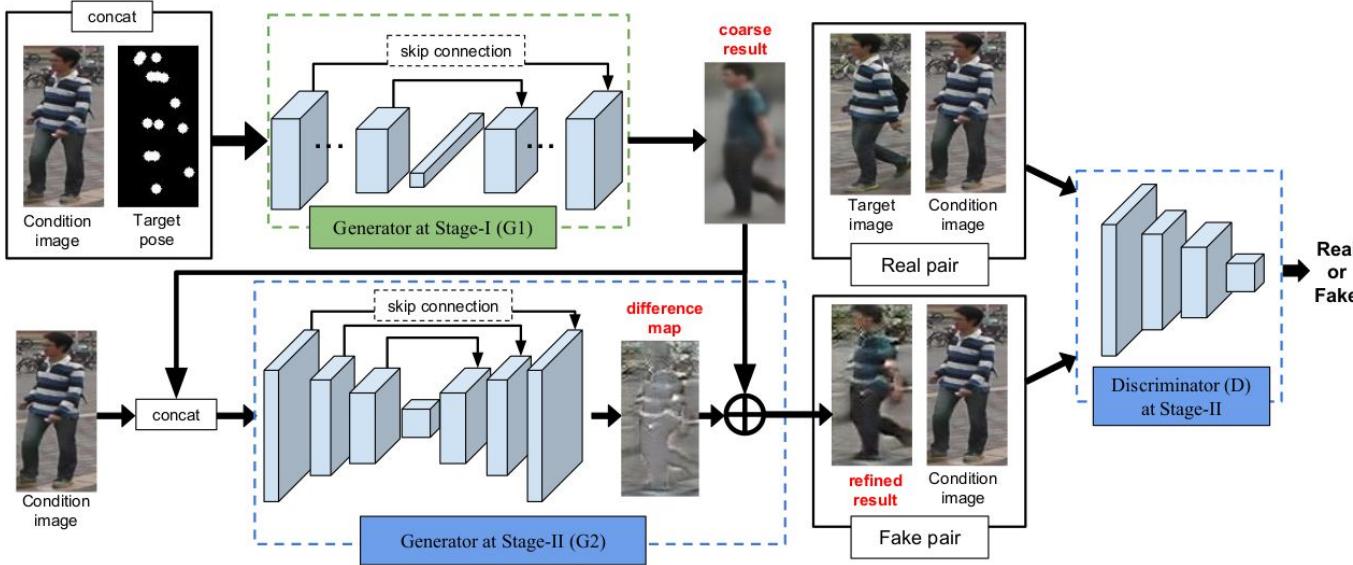
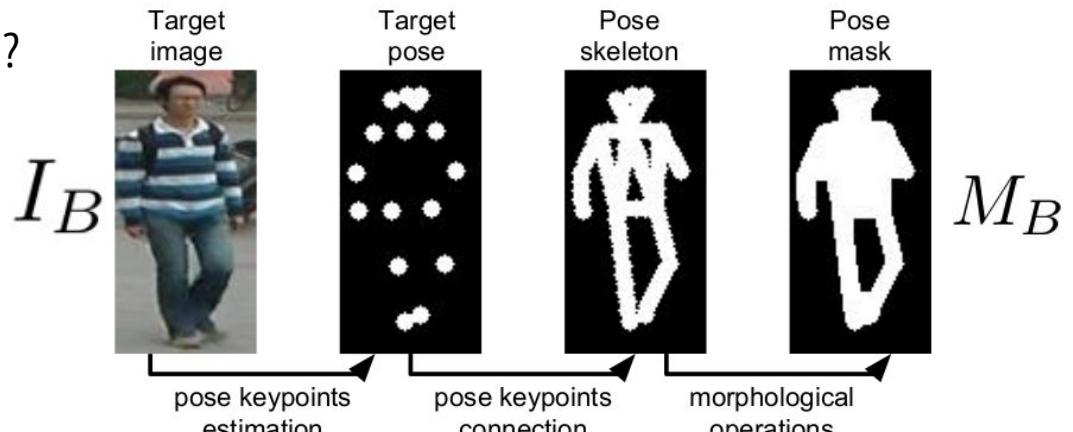


Figure 2: The overall framework of our Pose Guided Person Generation Network (PG²). It contains two stages. Stage-I focuses on pose integration and generates an initial result that captures the global structure of the human. Stage-II focuses on refining the initial result via adversarial training and generates sharper images.

Pose guided

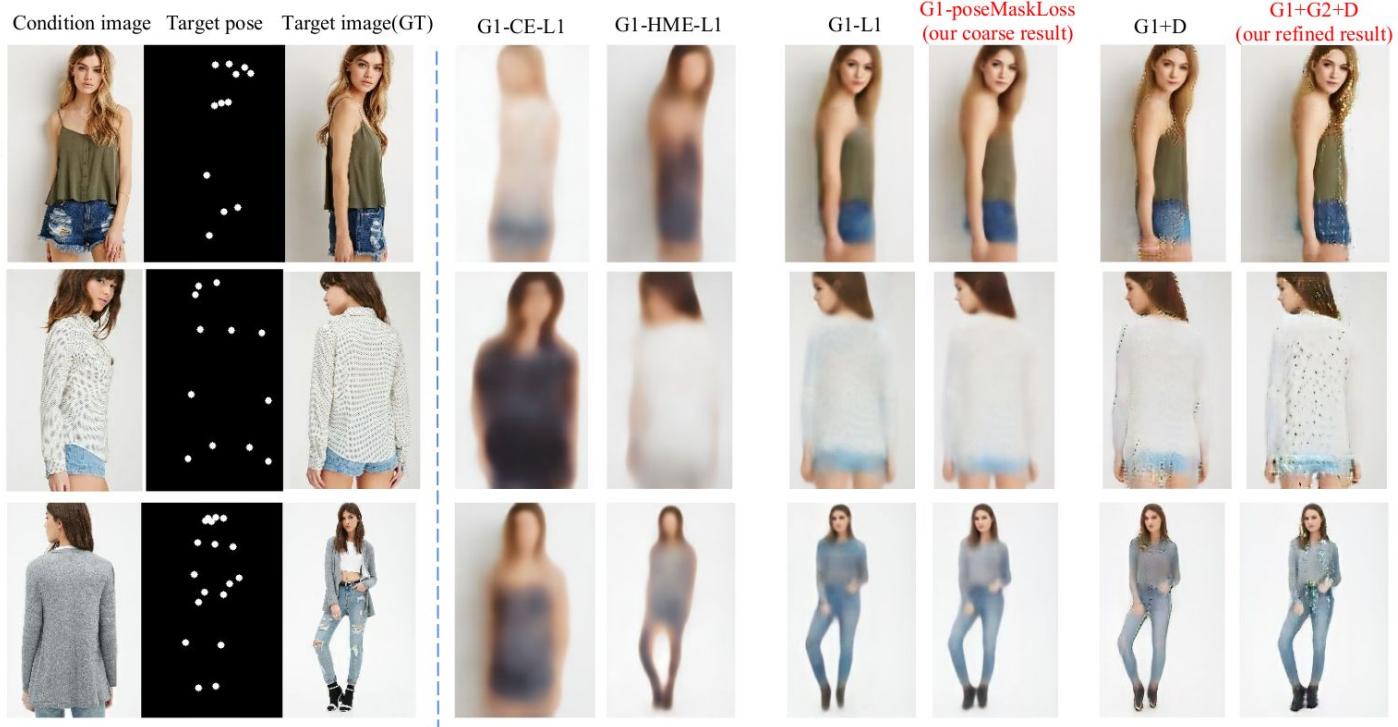
What about the background?



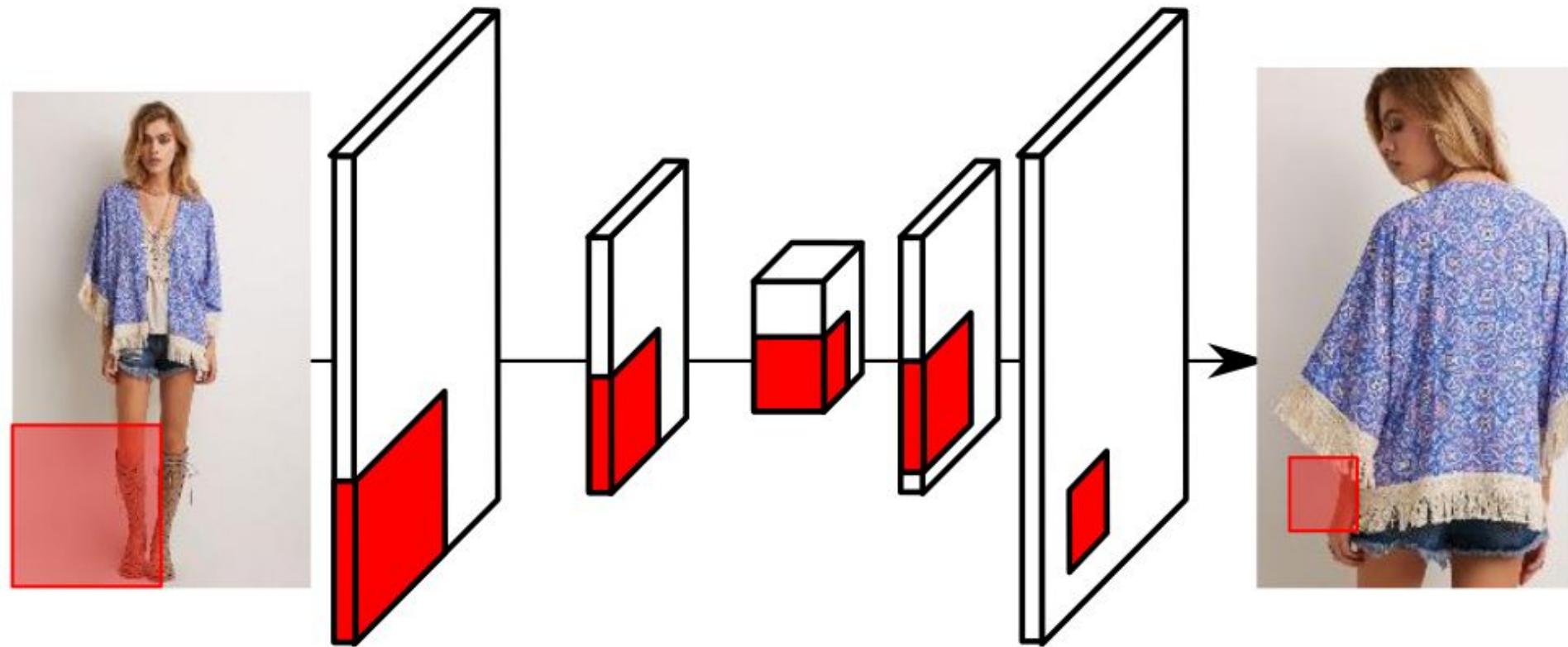
They give more weight to the foreground:

$$\mathcal{L}_{G1} = \|(\mathbf{G1}(I_A, P_B) - I_B) \odot (1 + M_B)\|_1$$

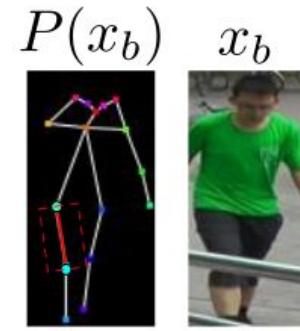
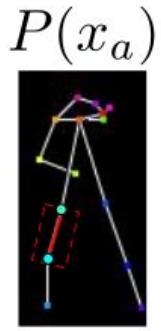
Pose guided



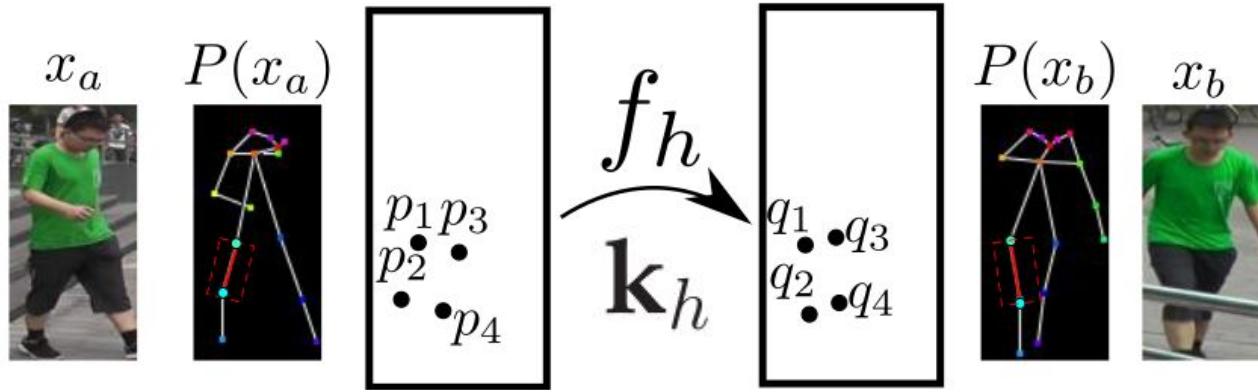
Pix2Pix for pose guided: problem



Pose guided: deformable GAN

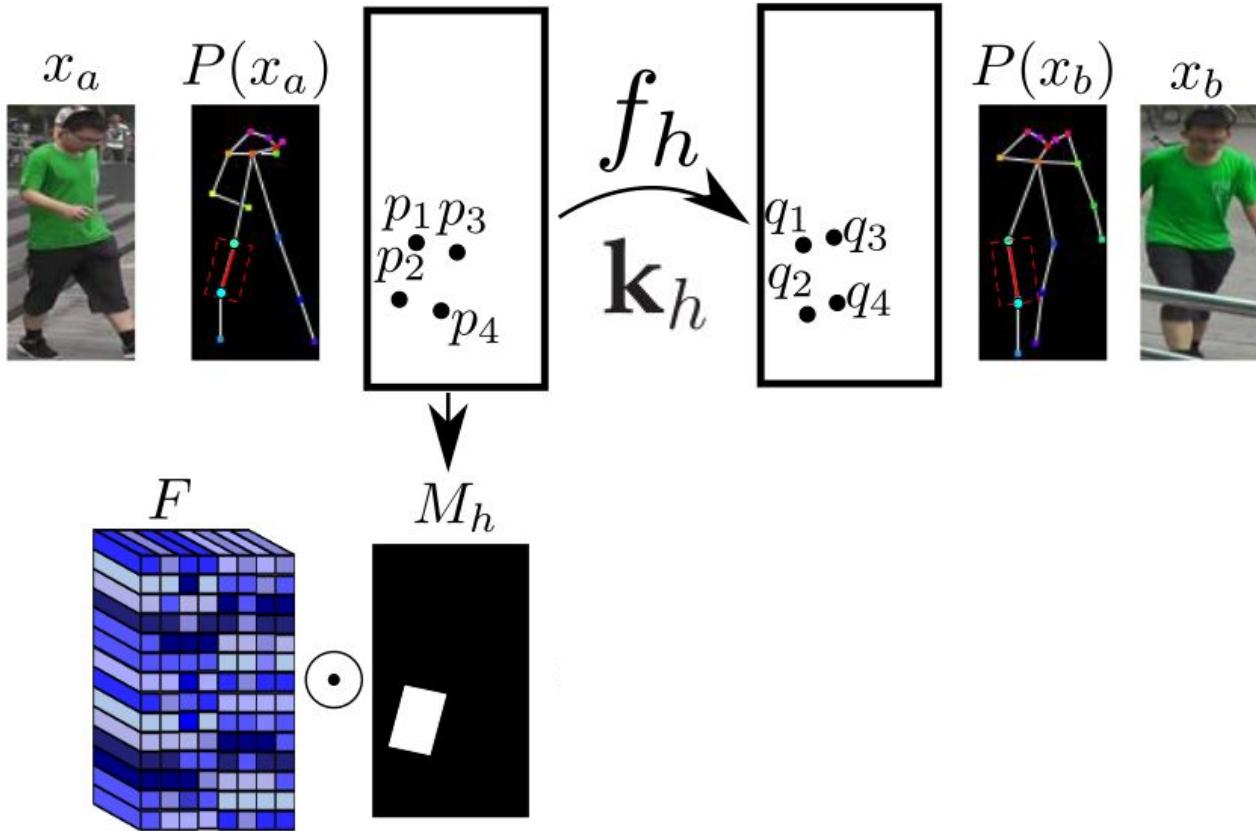


Pose guided: deformable GAN

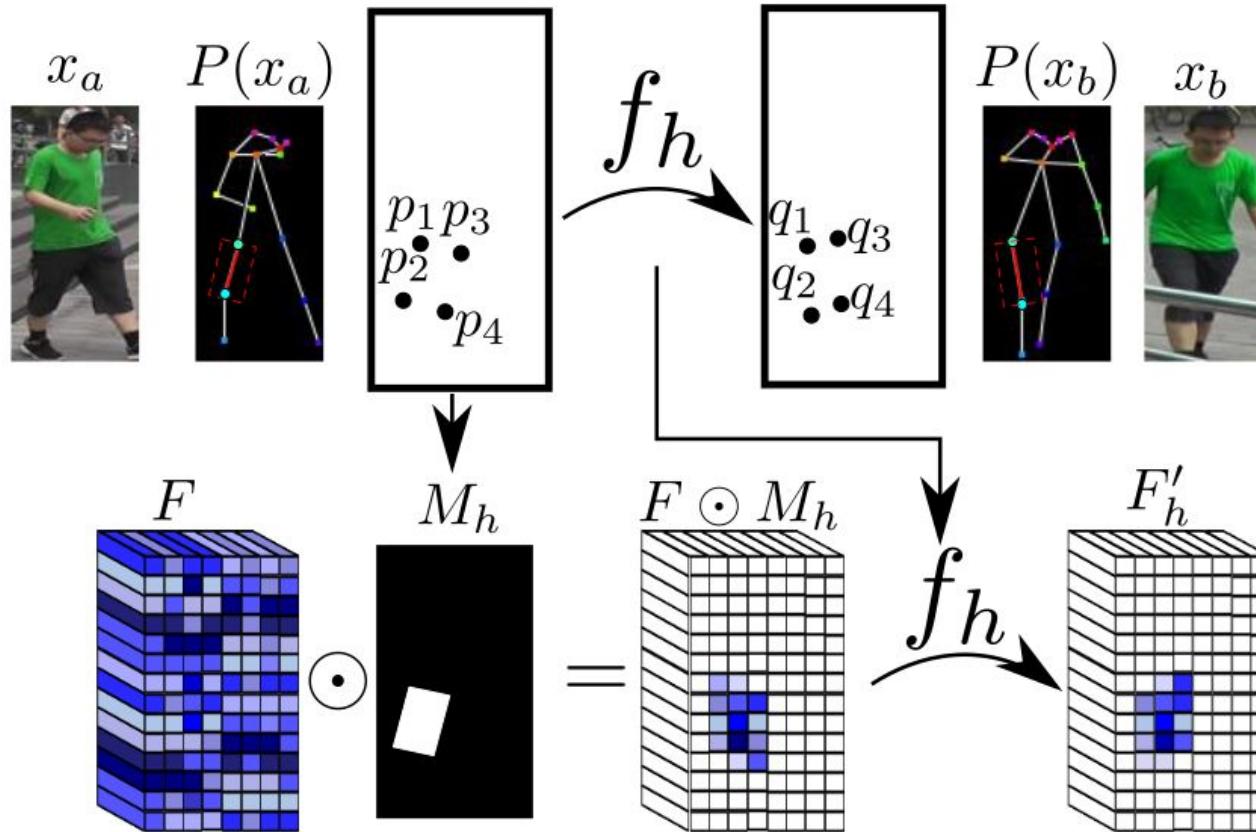


$$\min_{\mathbf{k}_h} \sum_{\mathbf{p}_j \in R_h^a, \mathbf{q}_j \in R_h^b} \|\mathbf{q}_j - f_h(\mathbf{p}_j; \mathbf{k}_h)\|_2^2.$$

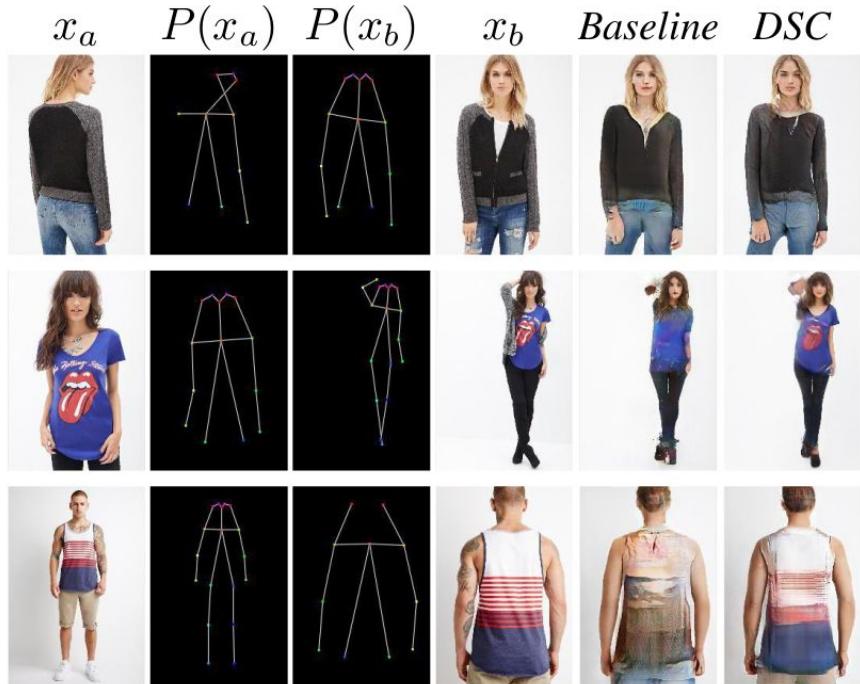
Pose guided: deformable GAN



Pose guided: deformable GAN



Pose guided: deformable GAN



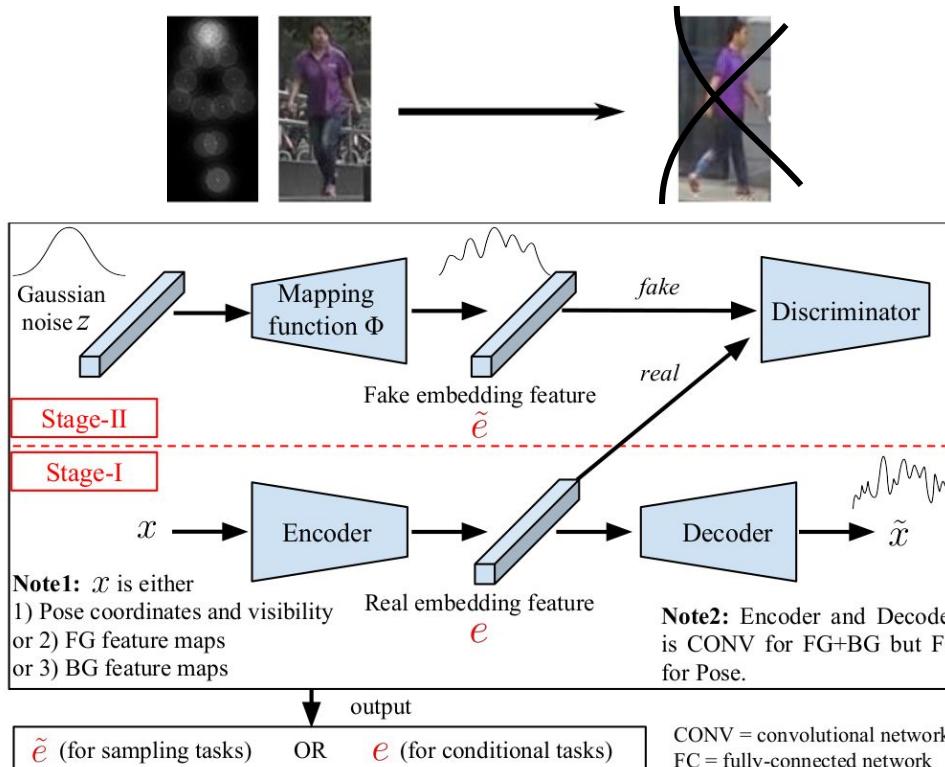
Pose guided: reconstruction vs GAN

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{NN}(G)$$

Market-1501					
λ	<i>SSIM</i>	<i>IS</i>	<i>mask-SSIM</i>	<i>mask-IS</i>	<i>DS</i>
0.1	0.292	2.621	0.808	3.168	0.697
0.01	0.290	3.185	0.805	3.502	0.720
0.001	0.245	3.566	0.779	3.634	0.609



Pose guided: Disentangled Generation



Pose guided: Disentangled Generation

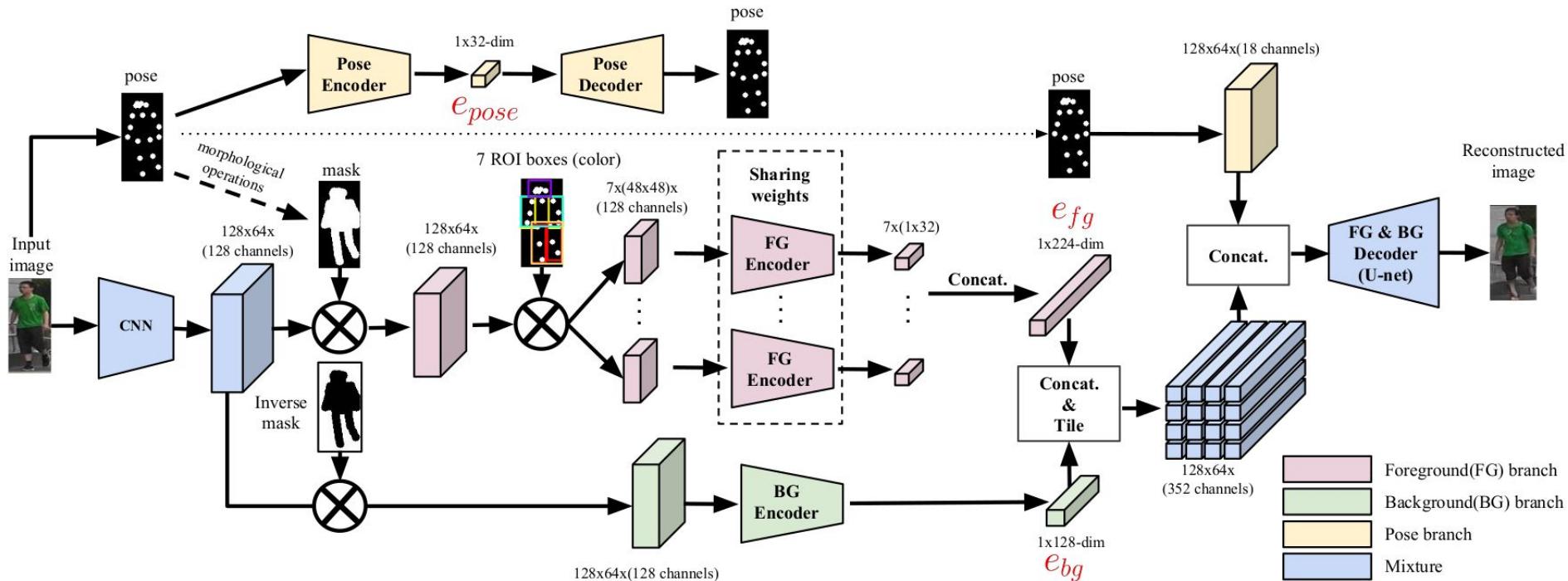


Figure 3: Stage-I: disentangled image reconstruction. This framework is composed of three branches: foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features to obtain the heatmaps at the sampling phase.

Pose guided: Disentangled Generation



Foreground (FG) sampling (fixed BG and Pose)



Background (BG) sampling (fixed FG and Pose)



Pose sampling (fixed FG and BG)

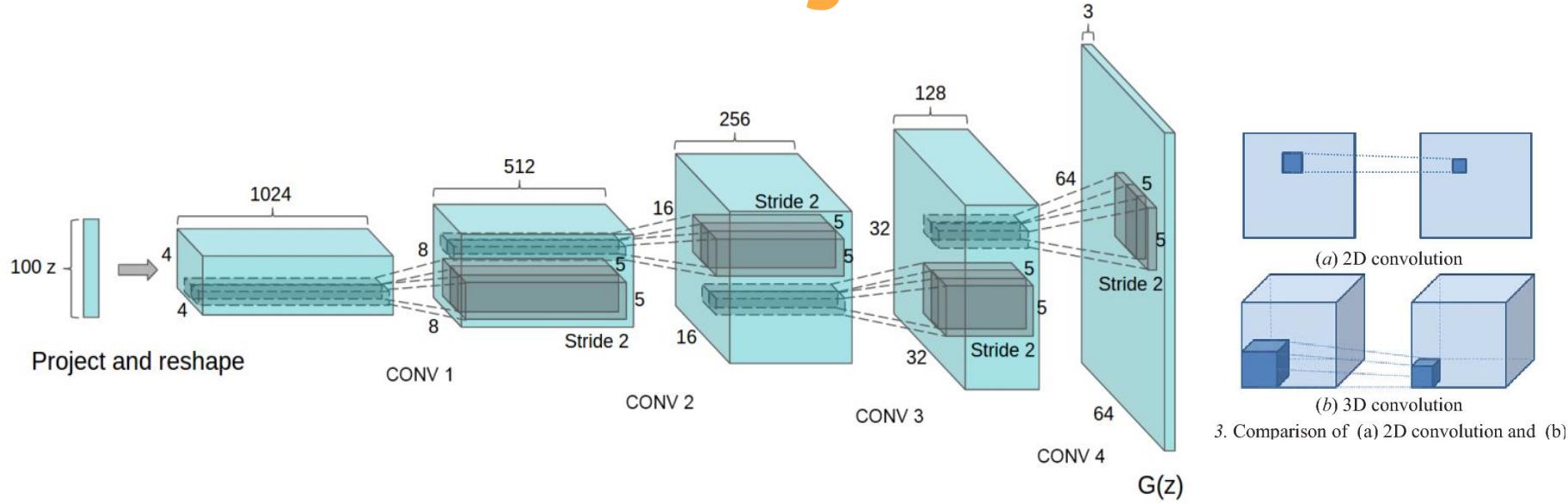


FG, BG and Pose sampling

Video generation

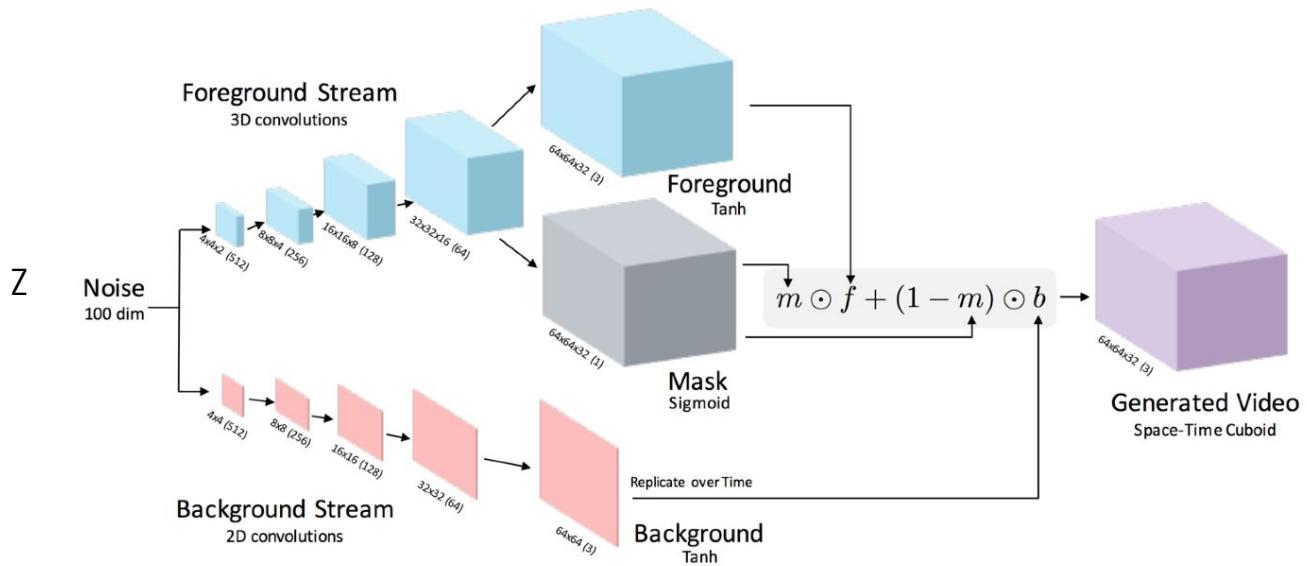


From DCGAN to video generation



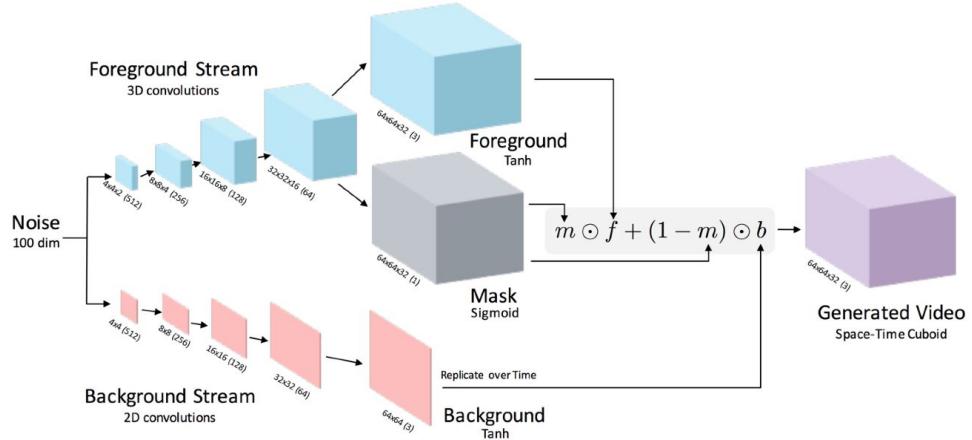
A. Radford, L. Metz, S. Chintala, [Unsupervised representation learning with deep convolutional generative adversarial networks](#), ICLR 2016
Figure from Lung nodule detection based on 3D convolutional neural networks Lei Fan et al.

Video Generation



Saito, Matsumoto, and Saito. "Temporal generative adversarial nets with singular value clipping." ICCV'2017.

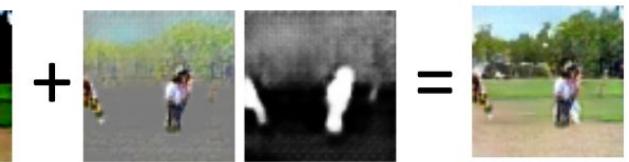
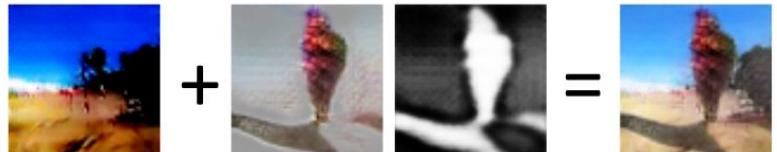
Video Generation



Background Foreground Generation

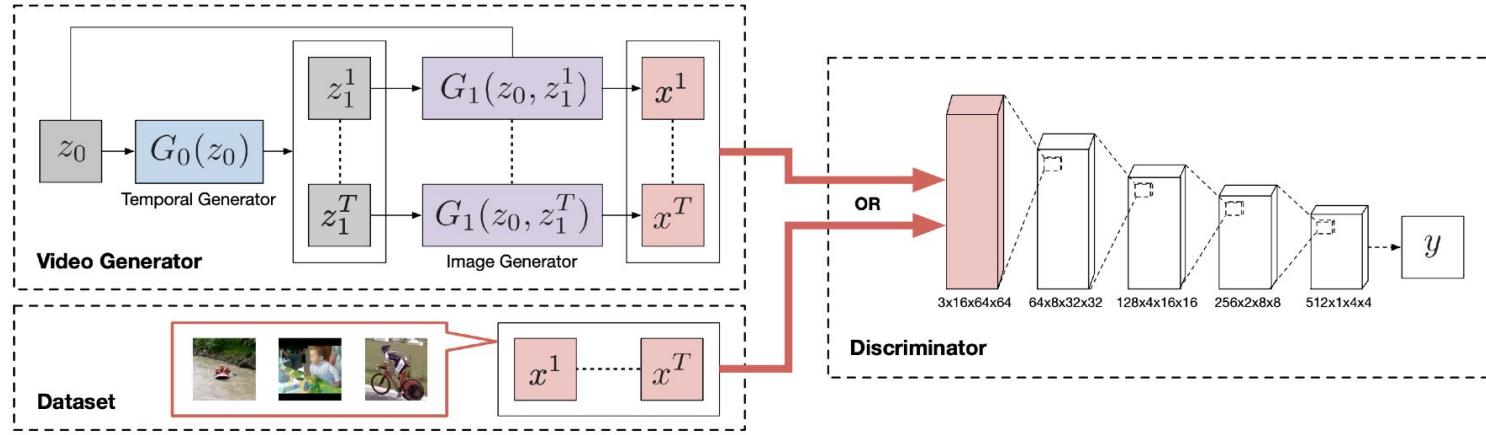


Background Foreground Generation



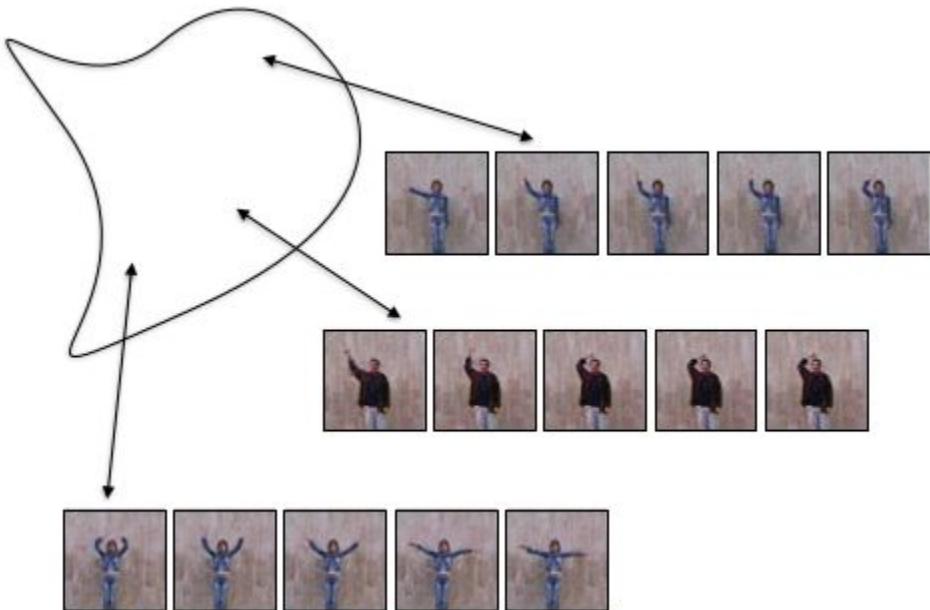
Saito, Matsumoto, and Saito. "Temporal generative adversarial nets with singular value clipping." ICCV'2017.

Video Generation



Saito, Matsumoto, and Saito. "Temporal generative adversarial nets with singular value clipping." ICCV'2017.

Video Generation



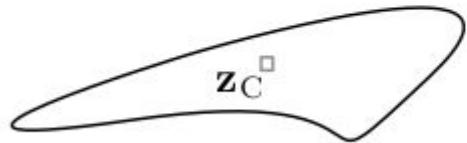
Limitations:

- Fixed length videos only
- No control over motion and content

Video Generation: MoCoGAN

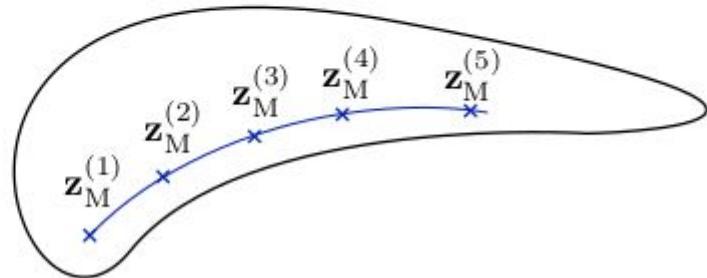
Sampled content

$$\mathbf{z}_C = [\mathbf{z}_C, \mathbf{z}_C, \dots, \mathbf{z}_C]$$

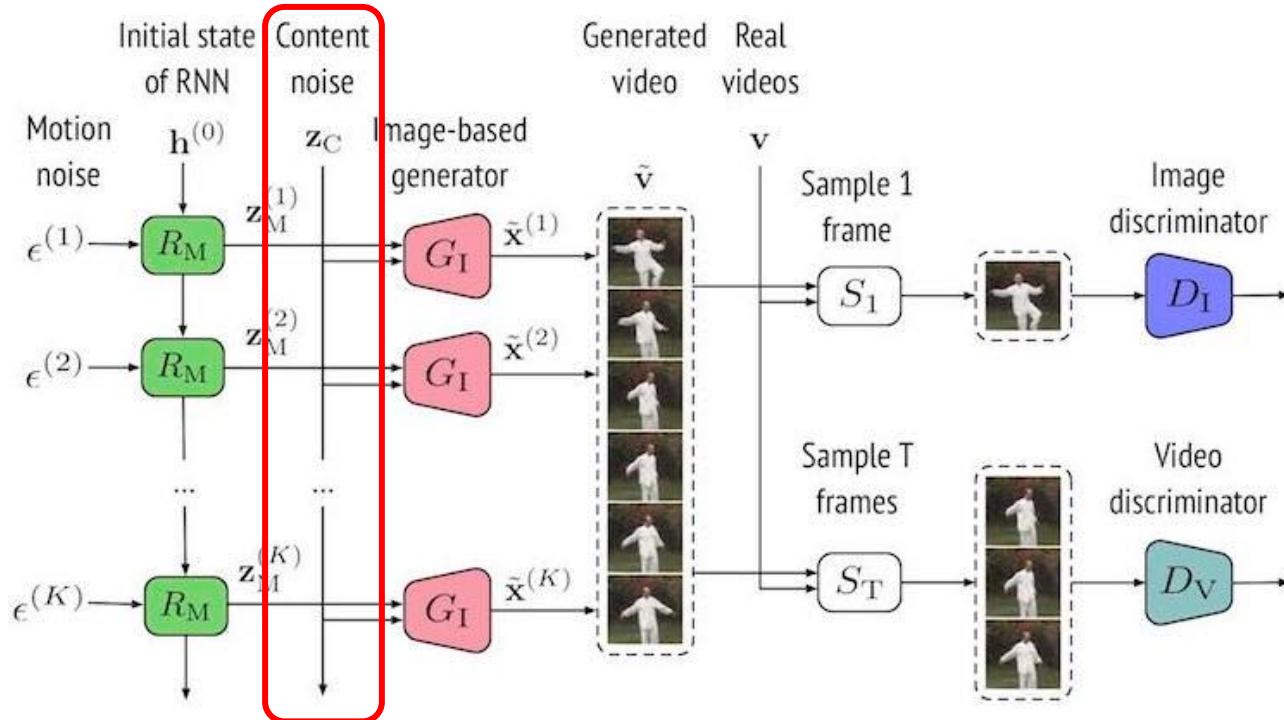


Motion trajectory

$$\mathbf{z}_M = [\mathbf{z}_M^{(1)}, \mathbf{z}_M^{(2)}, \dots, \mathbf{z}_M^{(K)}]$$

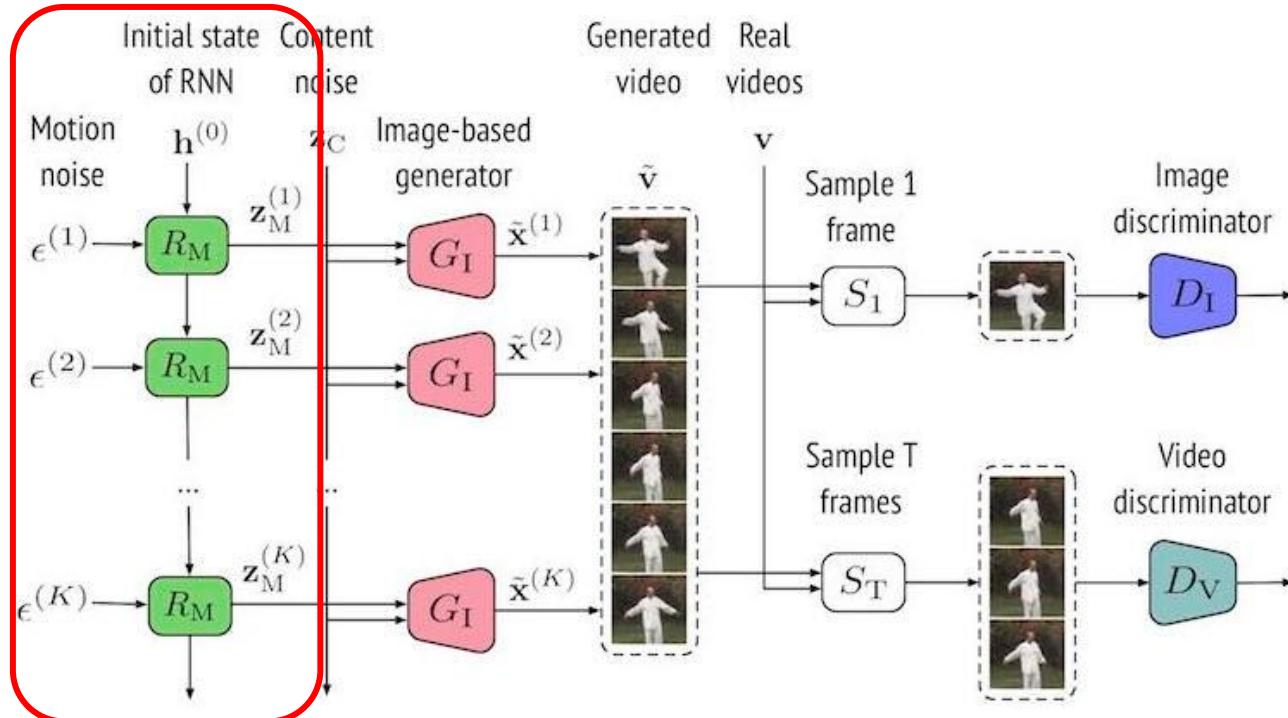


Video Generation: MoCoGAN



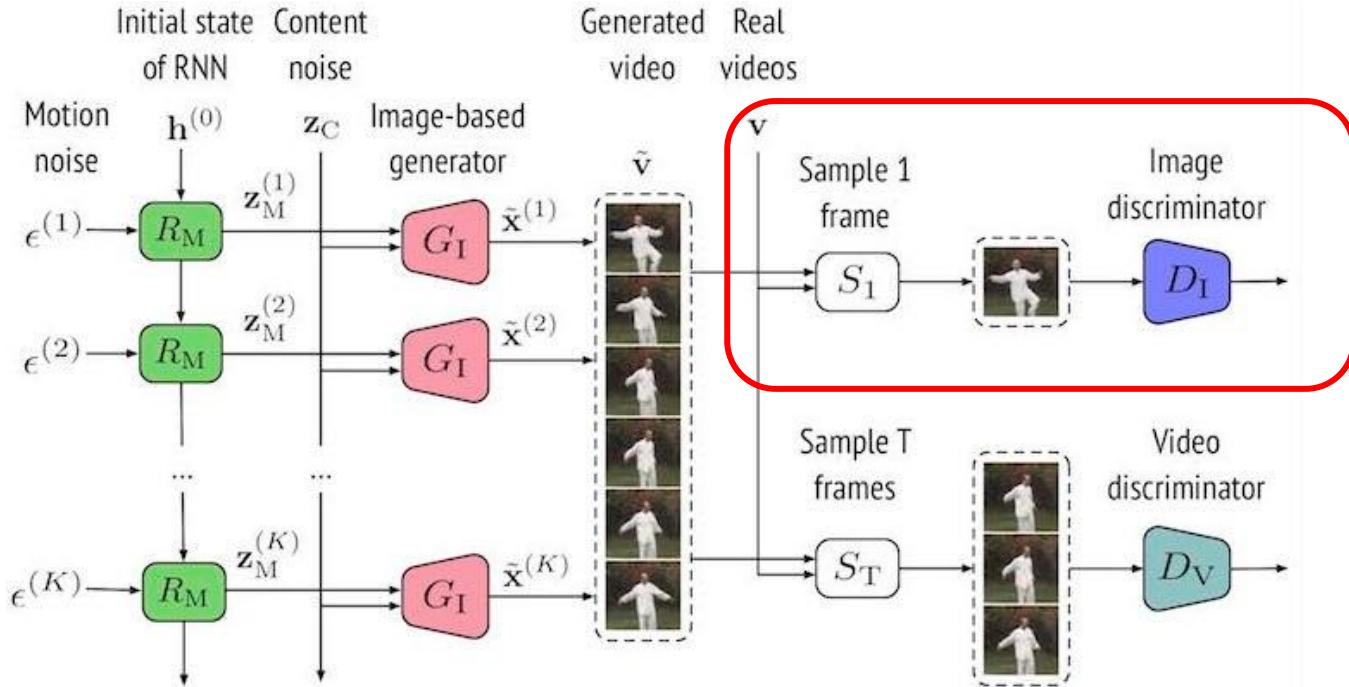
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN



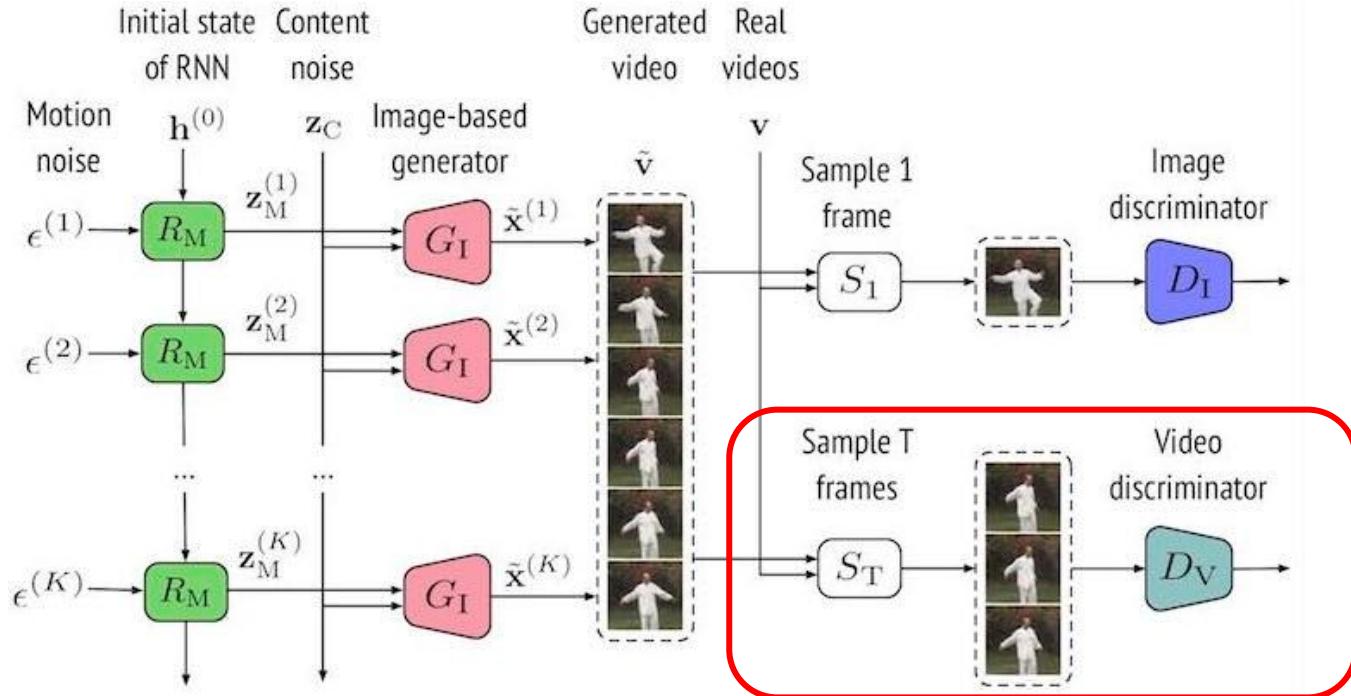
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN



MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN



MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN



MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN

We can change the motion without changing the content:



MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN

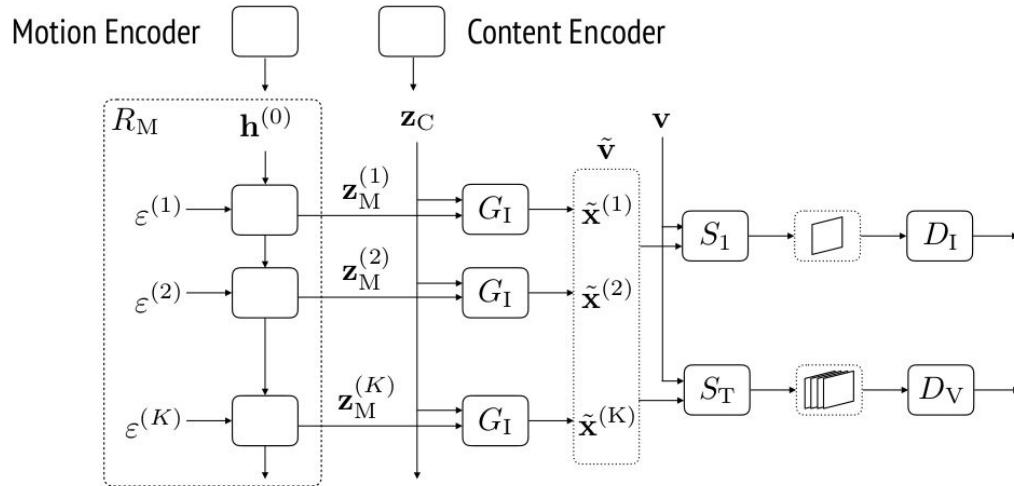
Video prediction:



MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: MoCoGAN

Video prediction:



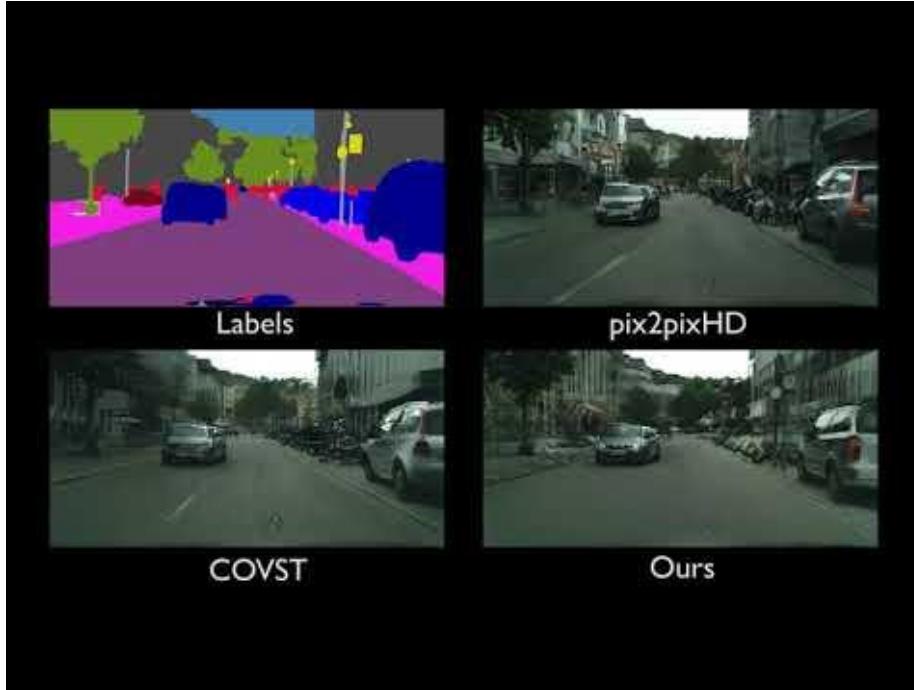
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

Video Generation: DVGAN



Clark, Aidan, Jeff Donahue, and Karen Simonyan. "Efficient video generation on complex datasets." ArXiv preprint

Video-to-video translation



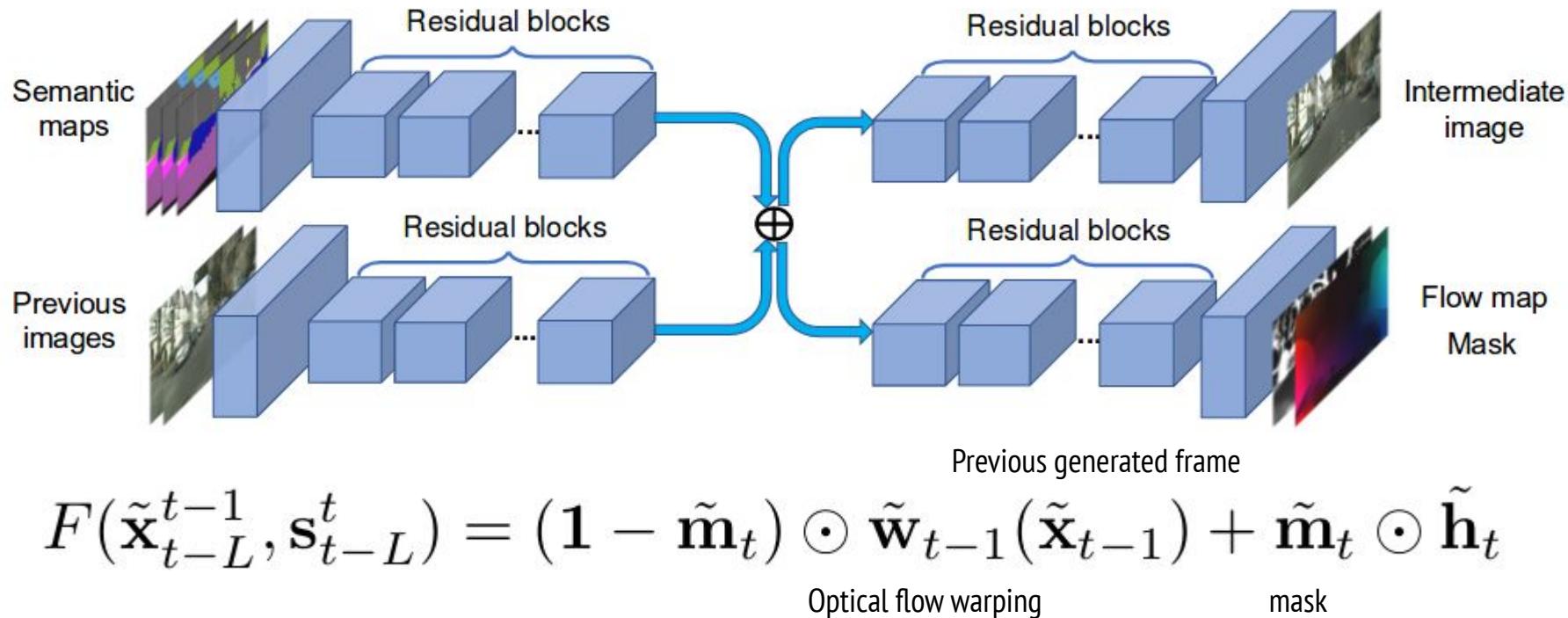
Wang, Ting-Chun, et al. "Video-to-video synthesis." NIPS'2018

Video-to-video translation



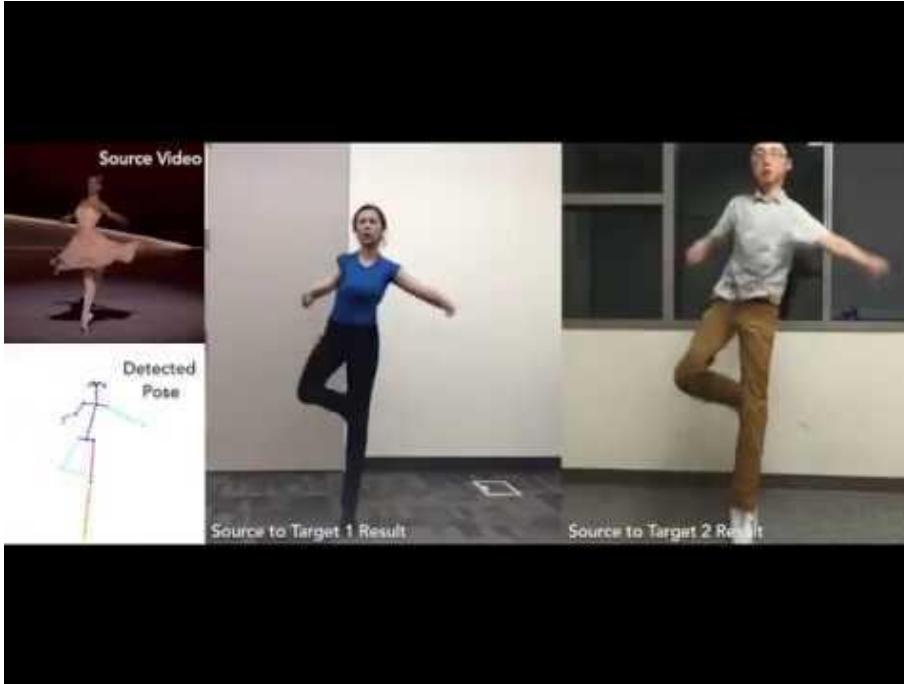
Wang, Ting-Chun, et al. "Video-to-video synthesis." NIPS'2018

Video-to-video translation



Wang, Ting-Chun, et al. "Video-to-video synthesis." NIPS'2018

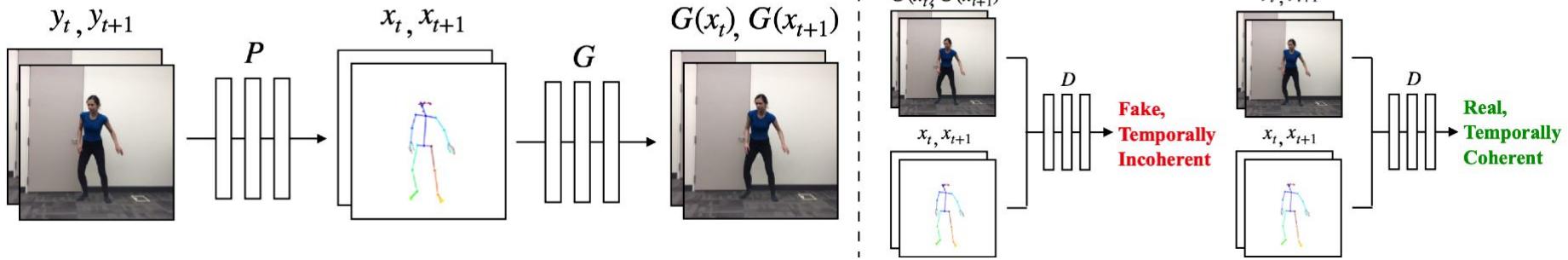
Animating Single Subject



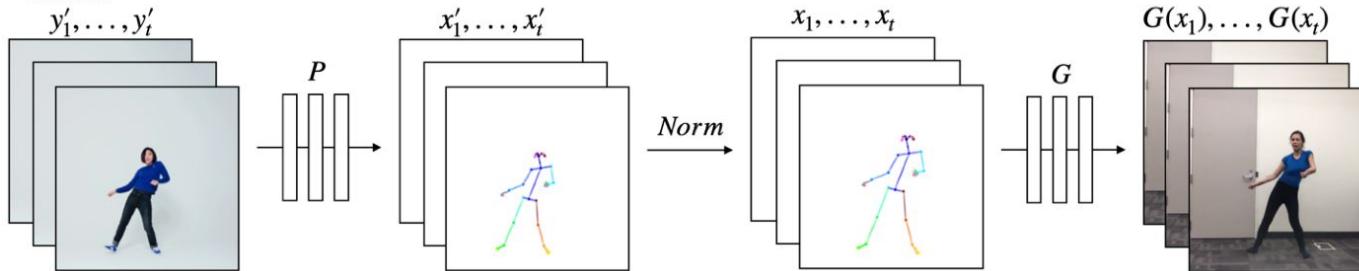
Chan, Caroline, et al. "Everybody dance now." ICCV'2019.

Animating Single Subject

Training

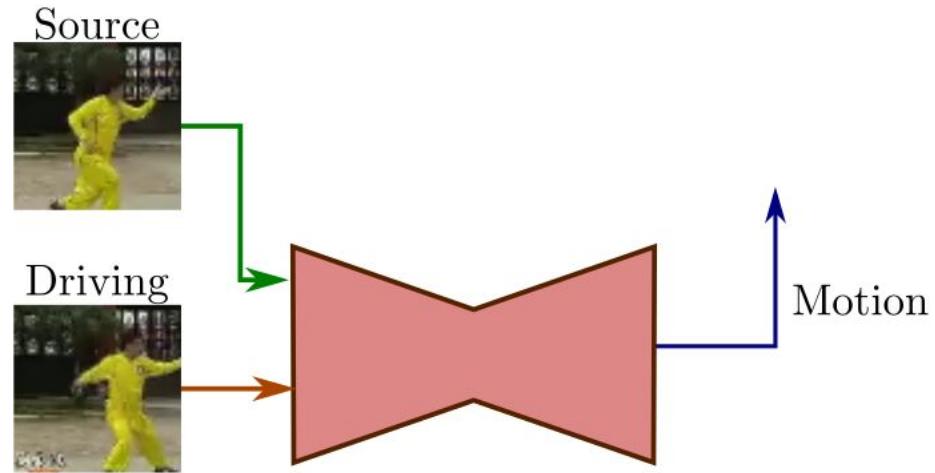


Transfer

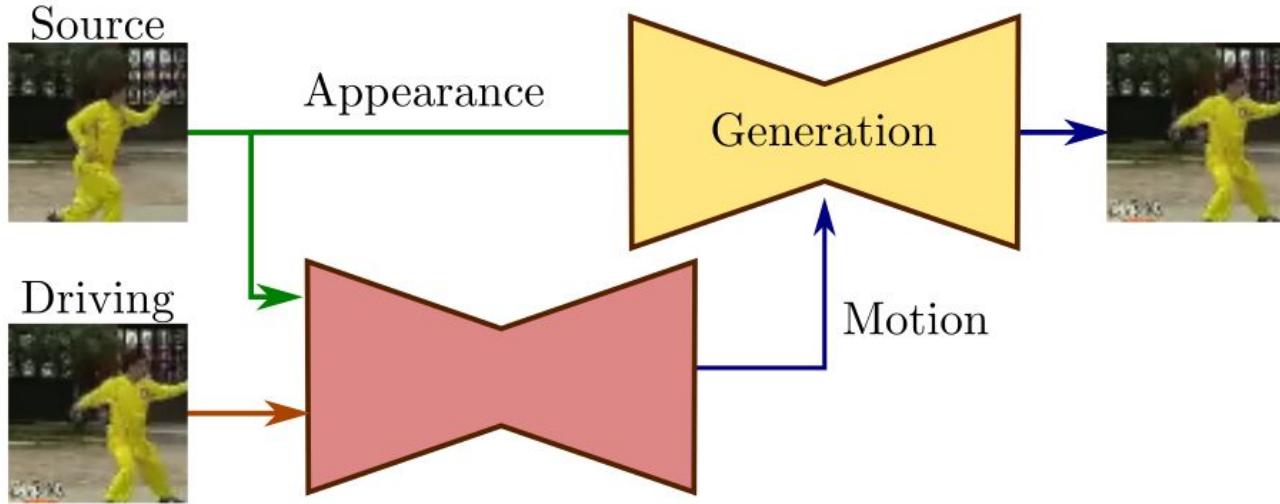


Chan, Caroline, et al. "Everybody dance now." ICCV'2019.

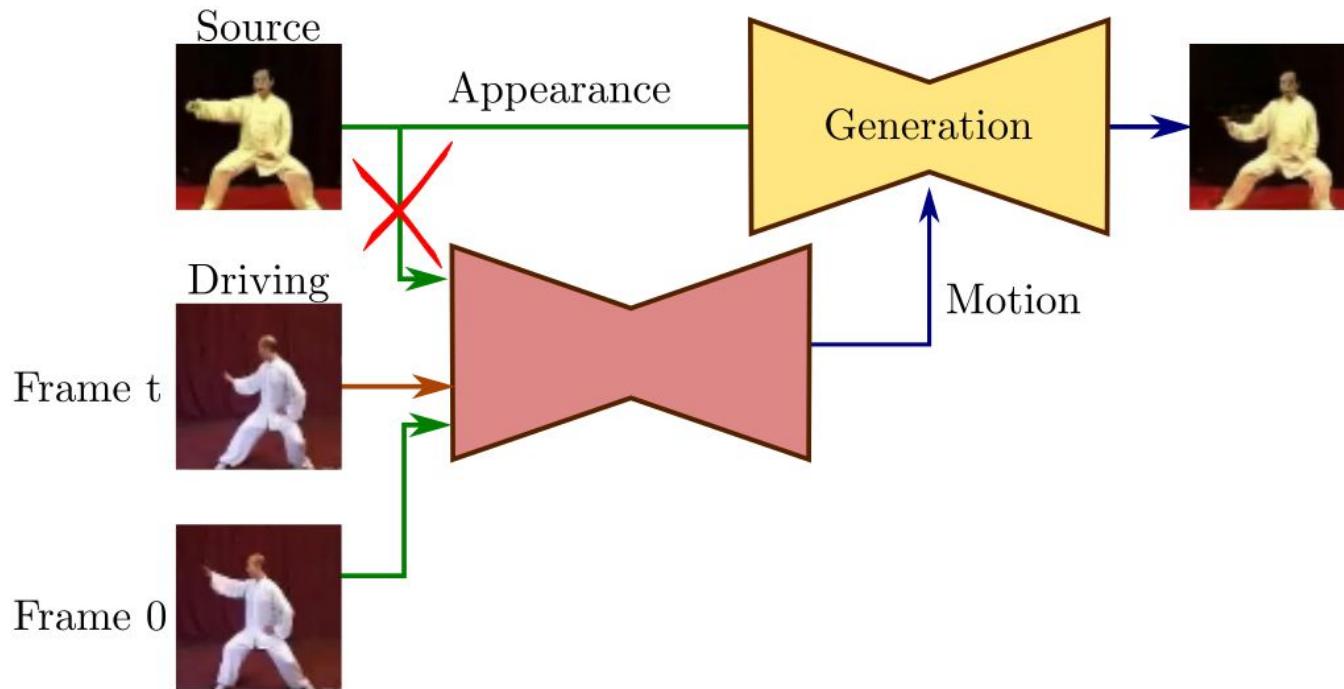
Retargeting



Retargeting



Retargeting



Retargeting



A. Siarohin, S. Lathuilière, S. Tulyakov , E. Ricci, N. Sebe, Animating Arbitrary Objects via Deep Motion Transfer , CVPR 2019