

# Image-to-image Translation

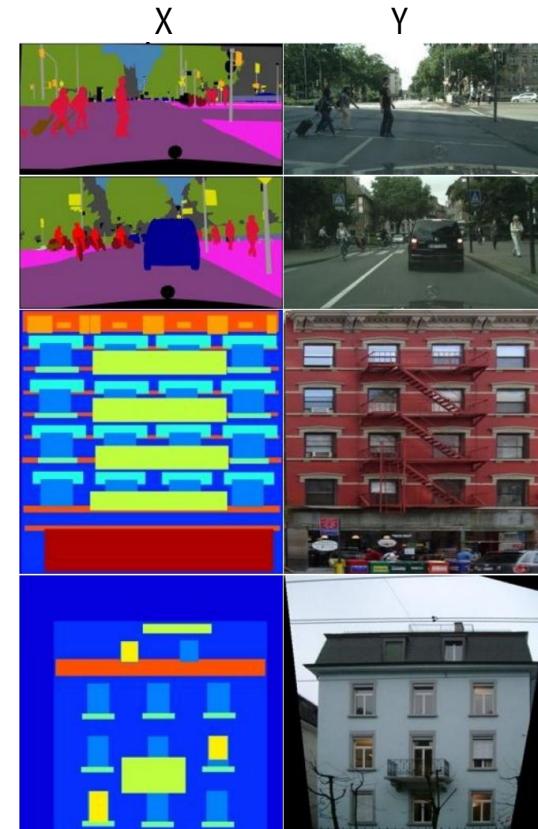
Given two domain the goal is to translate image from one possible representation to another.

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

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

Paired image-to-image translation

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



# Image-to-image Translation

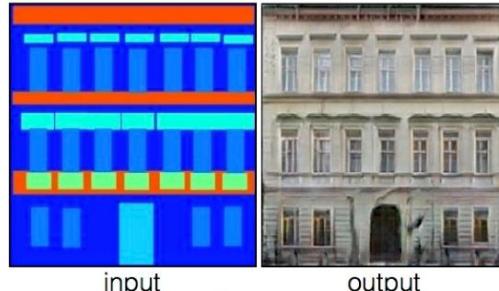
Labels to Street Scene



input

output

Labels to Facade



input

output

BW to Color



input

output

Aerial to Map



input

output

Day to Night



input

output

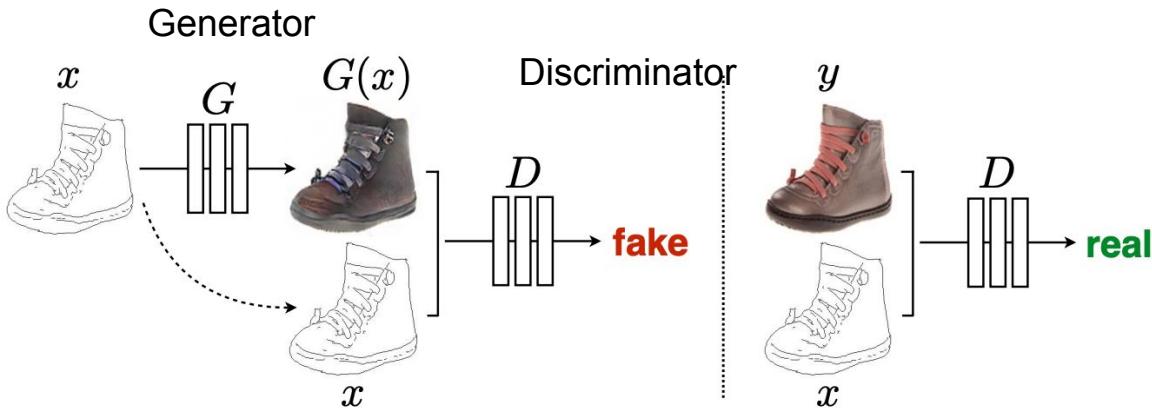
Edges to Photo



input

output

# Image-to-image Translation: Pix2Pix



Combined GAN-loss and reconstruction loss:

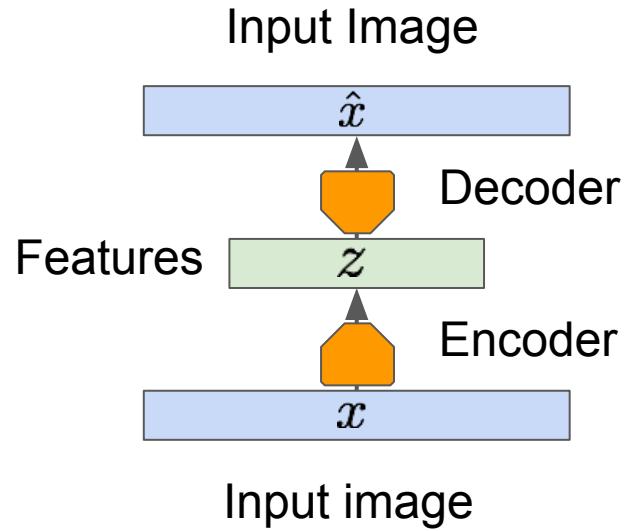
$$\begin{aligned}\mathcal{L}_{cGAN}(G, D) = & \mathbb{E}_{x,y}[\log D(x,y)] + \\ & \mathbb{E}_{x,z}[\log(1 - D(x, G(x, z)))]\end{aligned}$$

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

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

# Encoder-Decoder

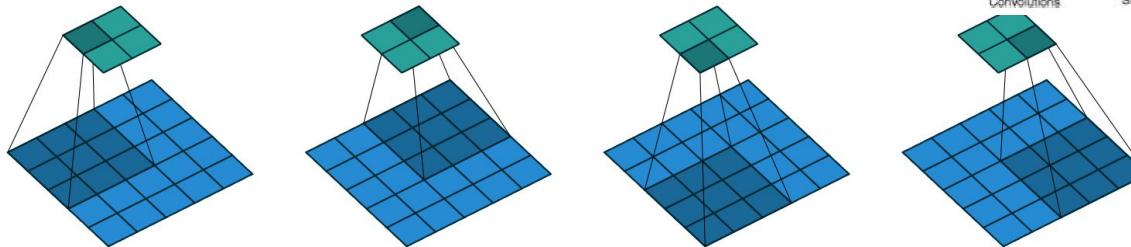
What is the architecture of G



# “De-convolution”

Architecture:

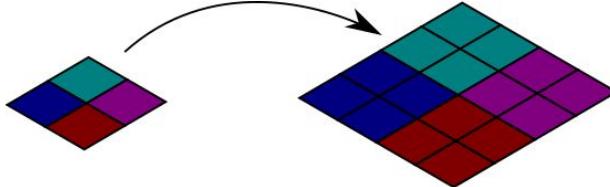
 Encoder: ConvNet architecture



Convolving a  $3 \times 3$  kernel over a  $5 \times 5$  input using  $2 \times 2$  strides

 Decoder: How to go from a vector to an image? How to “deconvolve”?

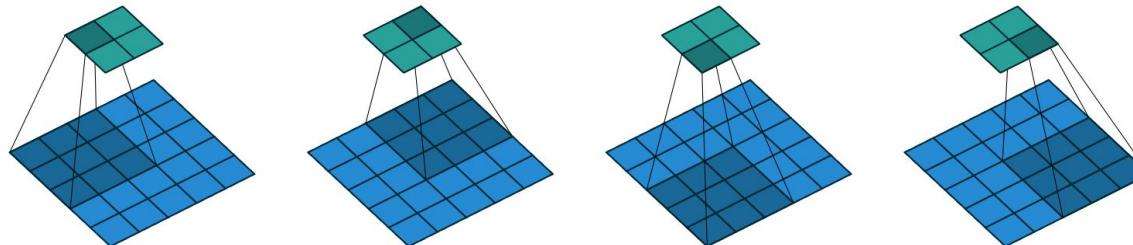
Up-sampling



# “De-convolution”

## Architecture:

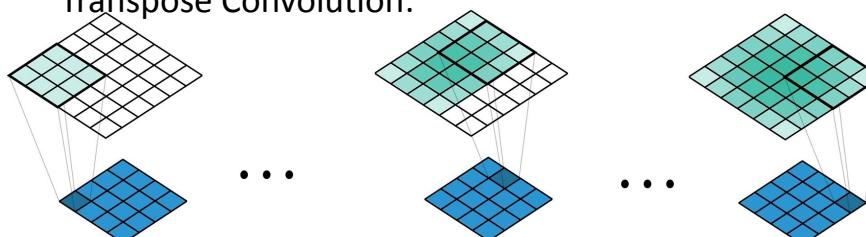
Encoder: ConvNet architecture



Convolving a  $3 \times 3$  kernel over a  $5 \times 5$  input using  $2 \times 2$  strides

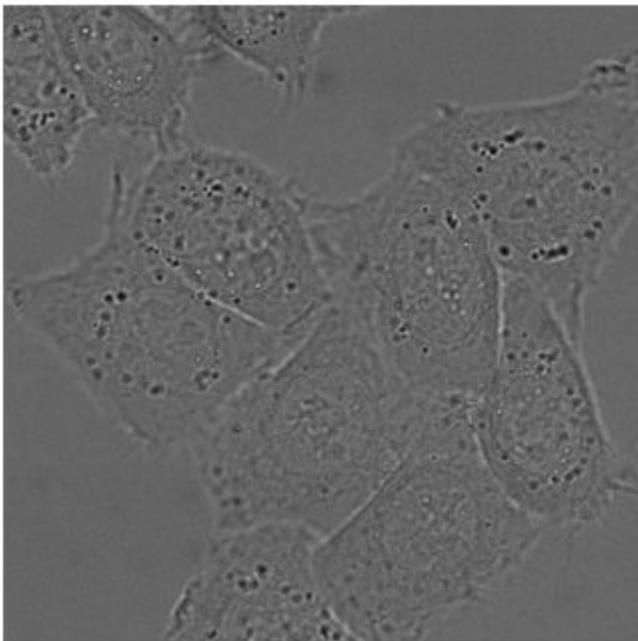
Decoder: How to go from a vector to an image? How to “deconvolve”?

Transpose Convolution:

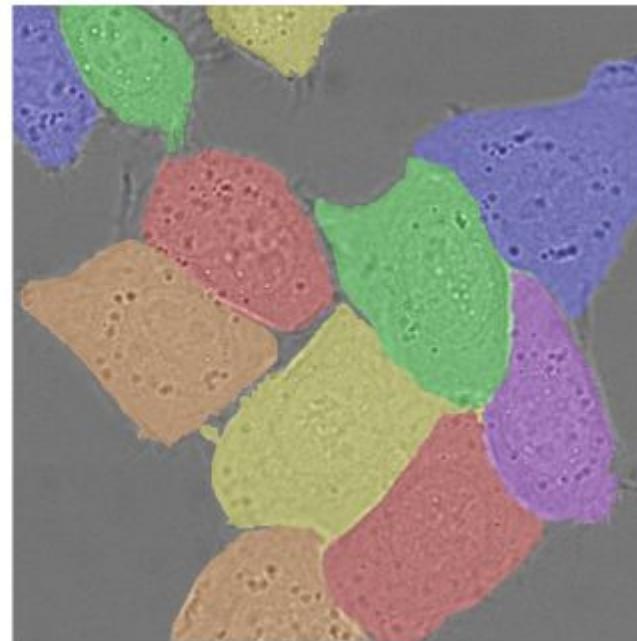


# Image-to-image Translation: Pix2Pix

a



b

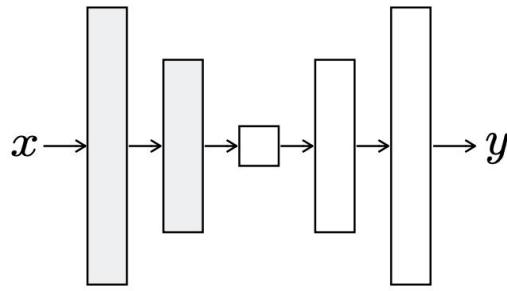


U-Net: Convolutional Networks for Biomedical Image Segmentation Olaf Ronneberger, Philipp Fischer, Thomas Brox MICCAI 2015

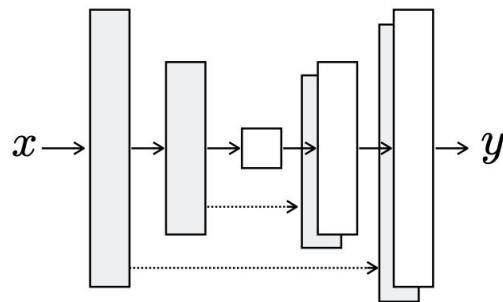
# Image-to-image Translation: Pix2Pix

Skip connections in generator

Encoder-decoder

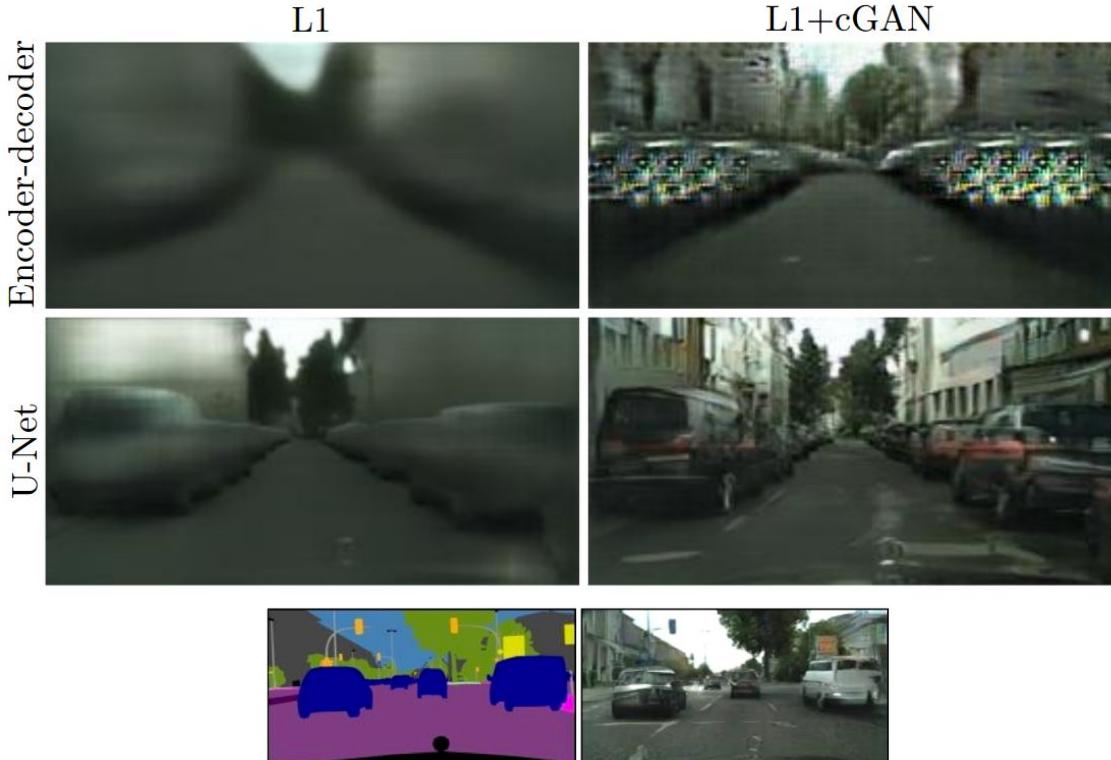


U-Net



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

# Image-to-image Translation: Pix2Pix

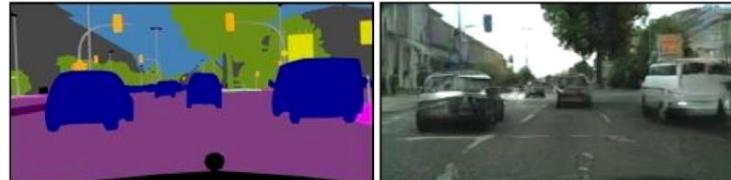


U-Net: Convolutional Networks for Biomedical Image Segmentation Olaf Ronneberger, Philipp Fischer, Thomas Brox MICCAI 2015

# Pix2Pix: Ablations

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

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

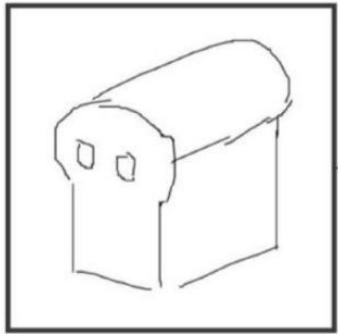
Discriminator receptive field	Per-pixel acc.	Per-class acc.	Class IOU
<b>1×1</b>	0.39	0.15	0.10
<b>16×16</b>	0.65	0.21	<b>0.17</b>
<b>70×70</b>	<b>0.66</b>	<b>0.23</b>	<b>0.17</b>
<b>286×286</b>	0.42	0.16	0.11

Table 3: FCN-scores for different receptive field sizes of the discriminator, evaluated on Cityscapes labels→photos. Note that input images are  $256 \times 256$  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”



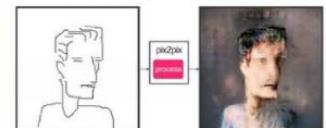
by Brannon Dorsey

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.

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

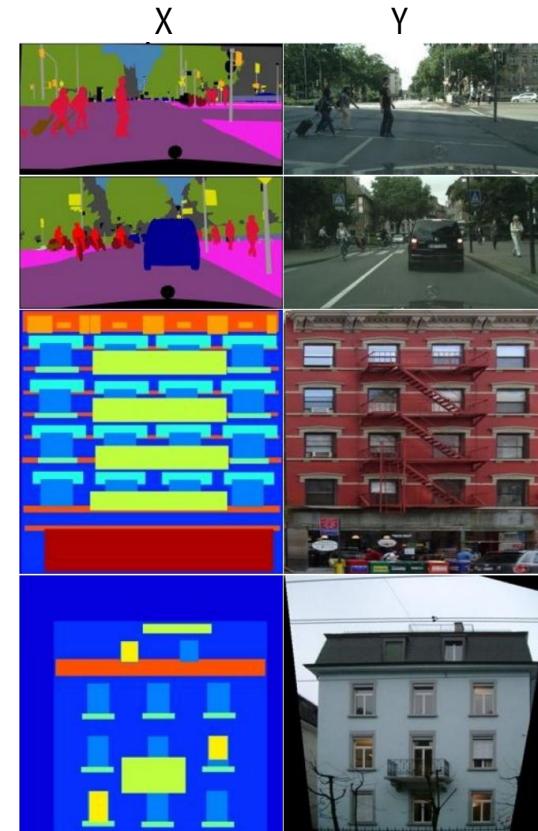
$$\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})$$

Paired image-to-image translation

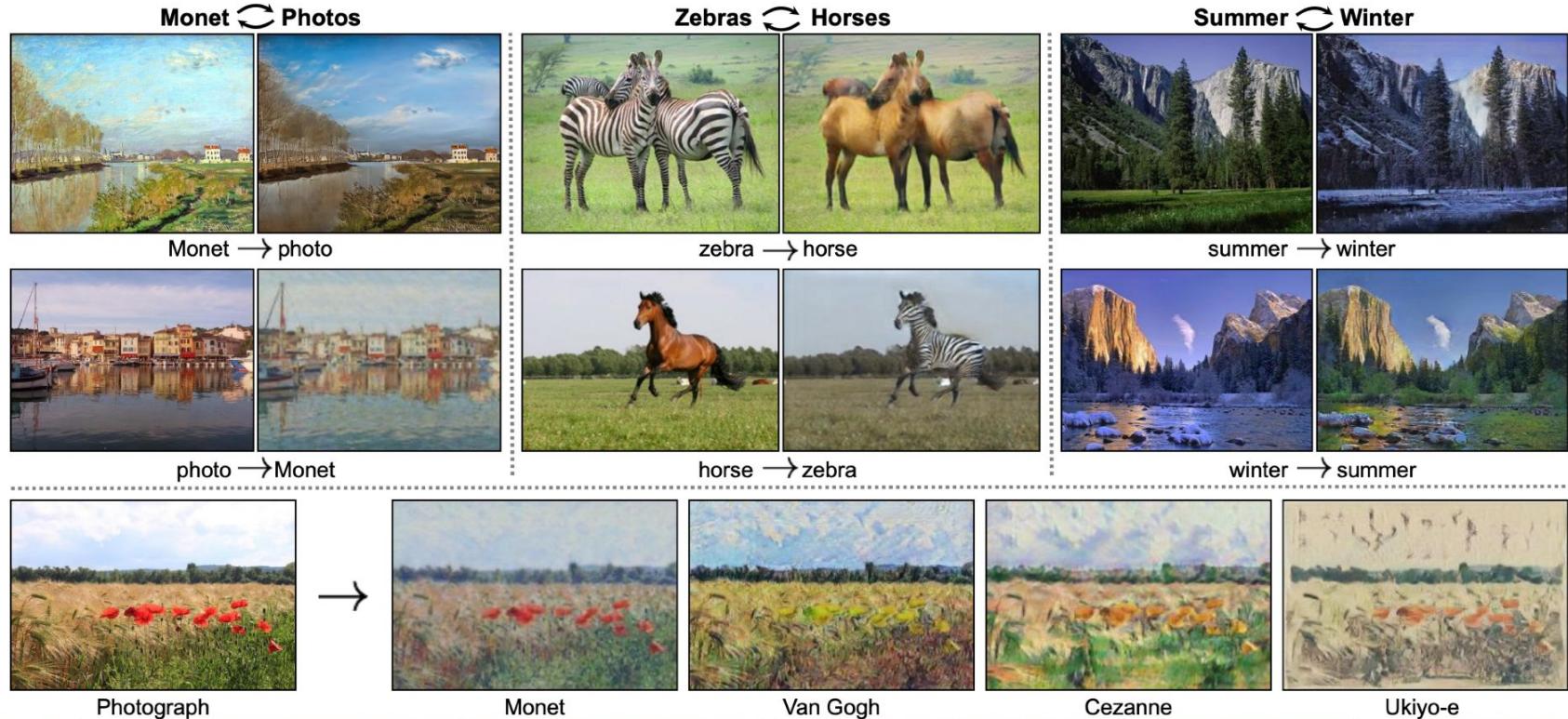
$$\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})$$

Unpaired

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



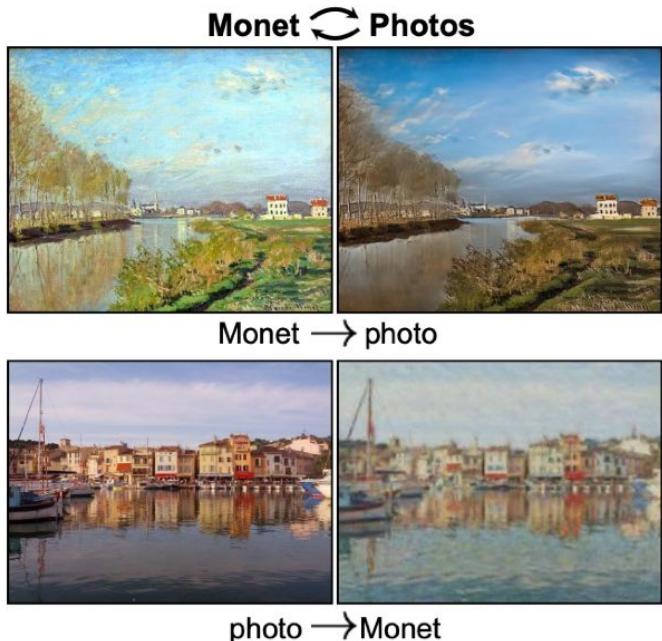
# Unpaired Image-to-image Translation



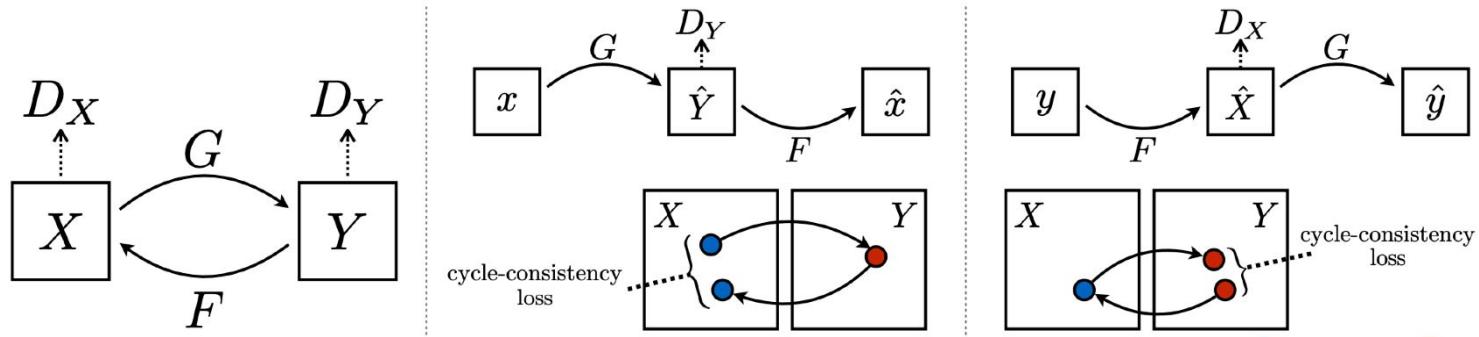
# Unpaired Image-to-image Translation

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:

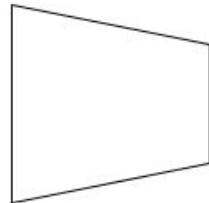
$$\begin{aligned}\mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]\end{aligned}$$

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

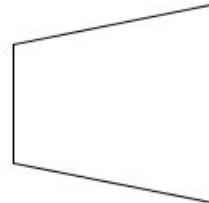
Downsampling



Strided conv

Batch-norm

Upsampling



ResBlocks

Conv

Batch-norm

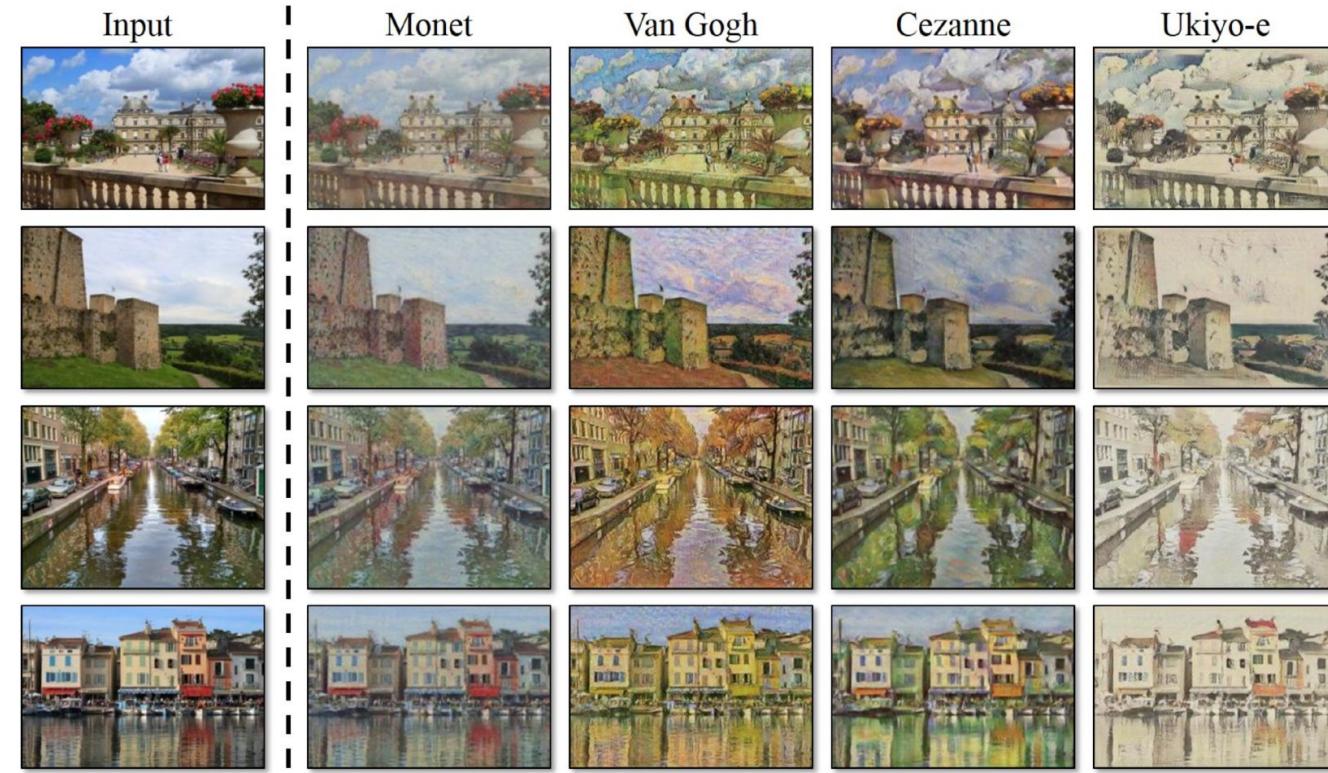
ReLU

Conv

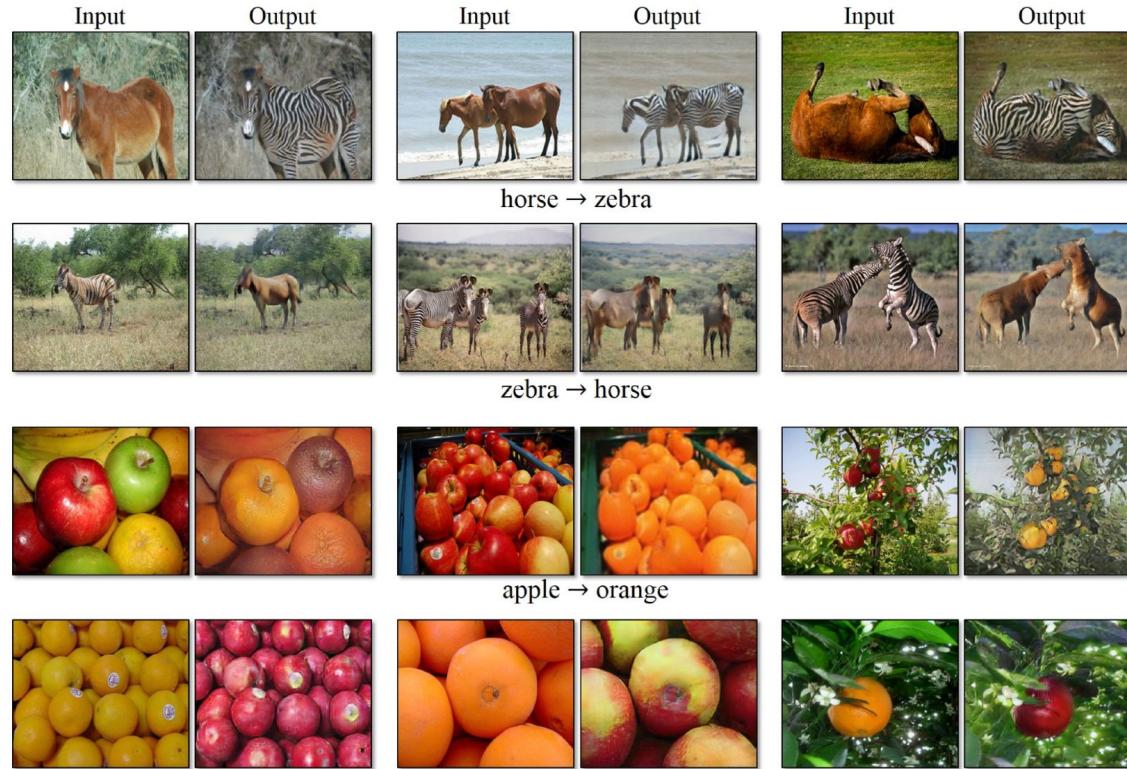
Batch-norm

add input

# CycleGAN



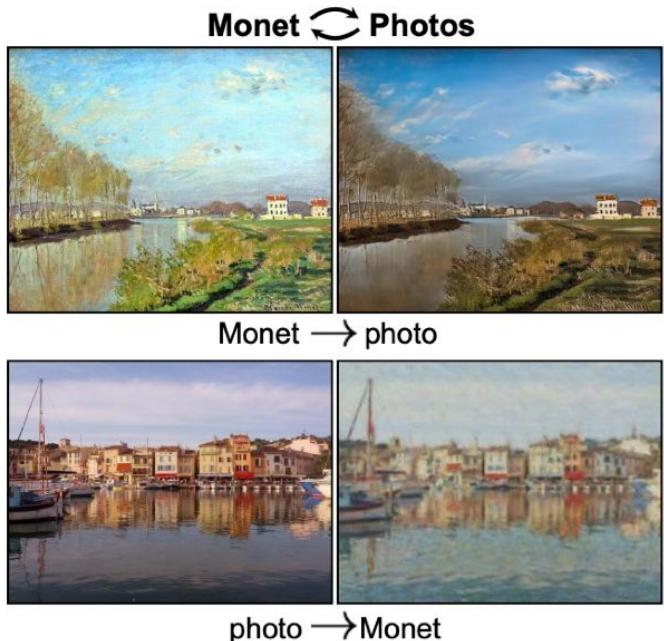
# CycleGAN



# Unpaired Image-to-image Translation

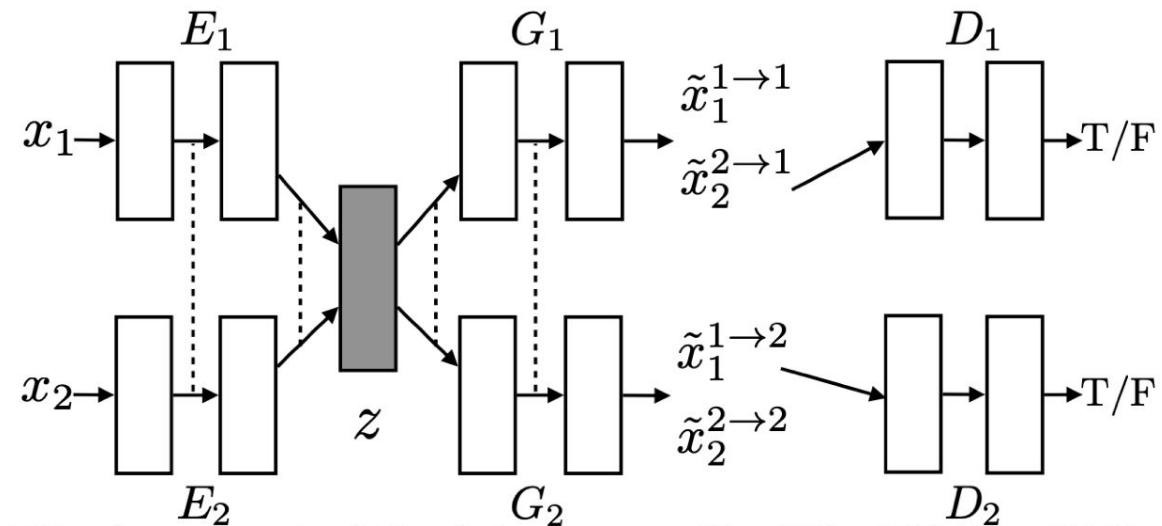
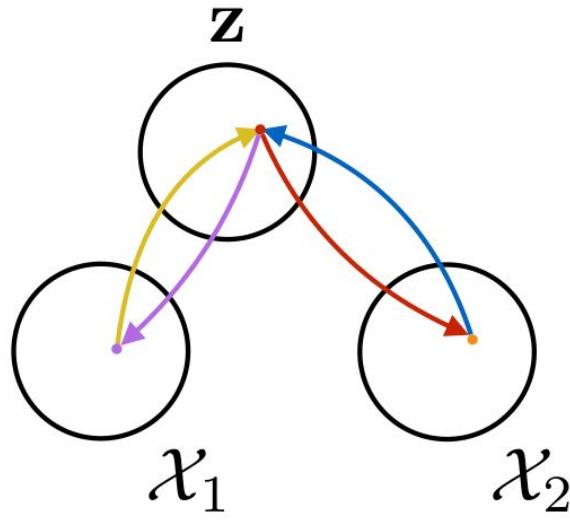
To solve it, constraints are necessary:

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



# Unpaired Translation with weight sharing

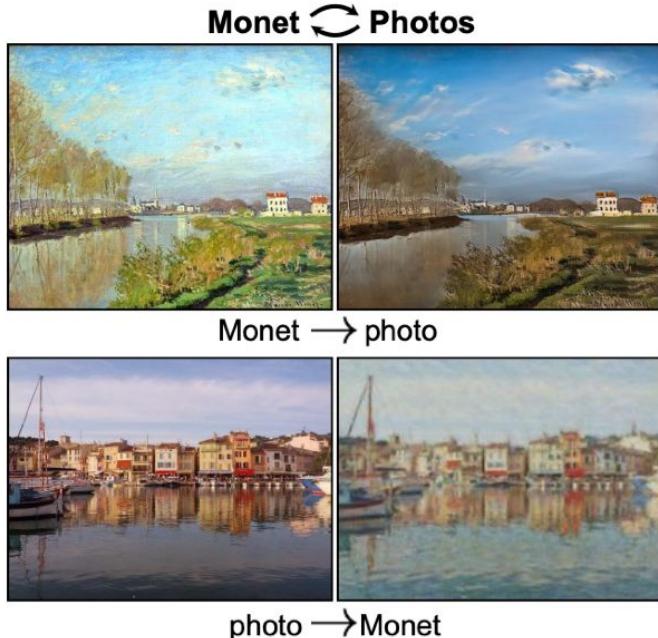
Train encoders with shared weights to go from image space to shared latent space



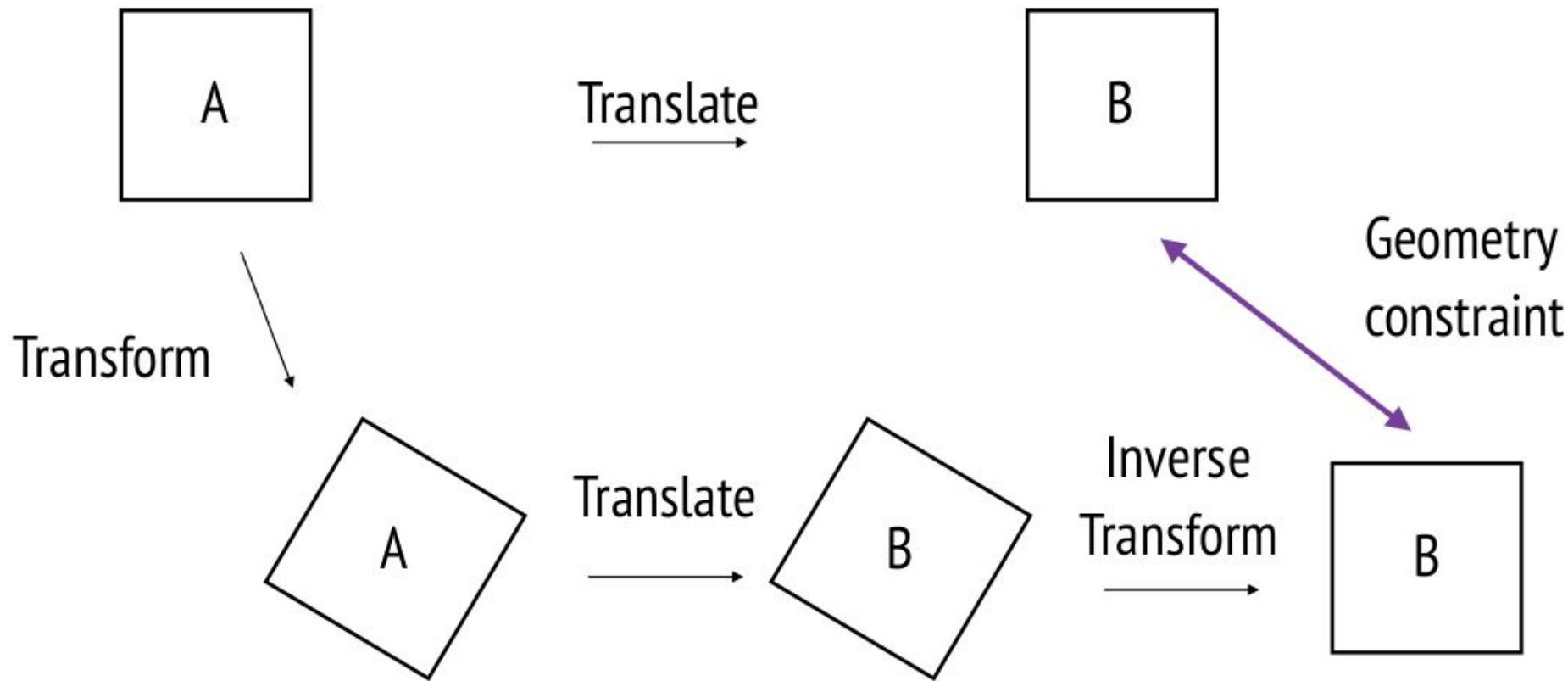
# Unpaired Image-to-image Translation

To solve it, constraints are necessary:

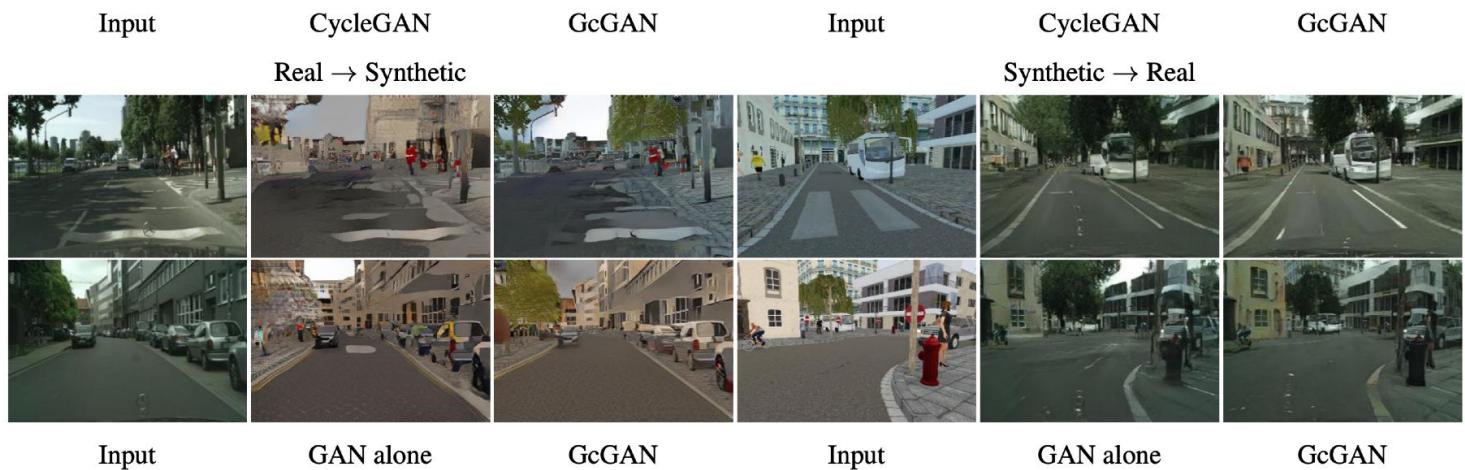
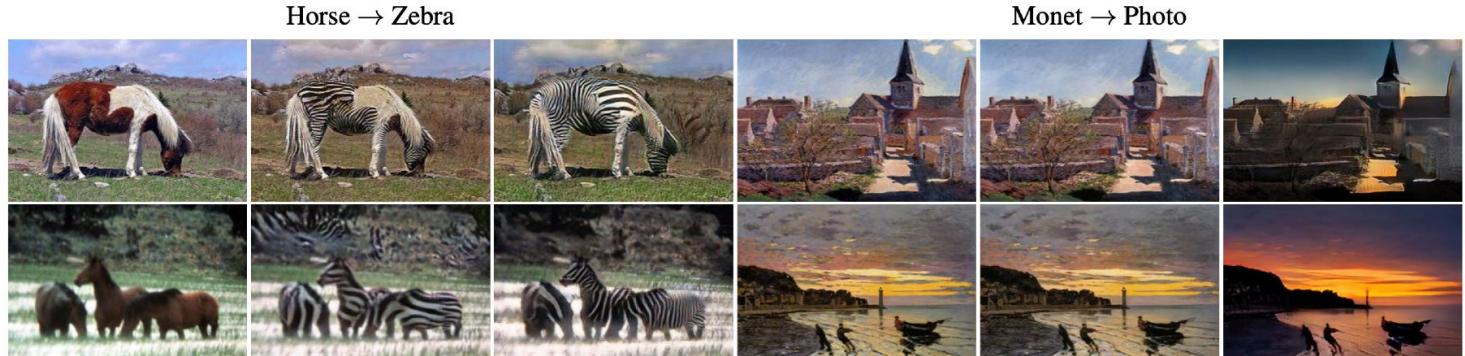
- Cycle-consistency constraint
- Weight-sharing constraint
- **Geometry-consistency constraint**



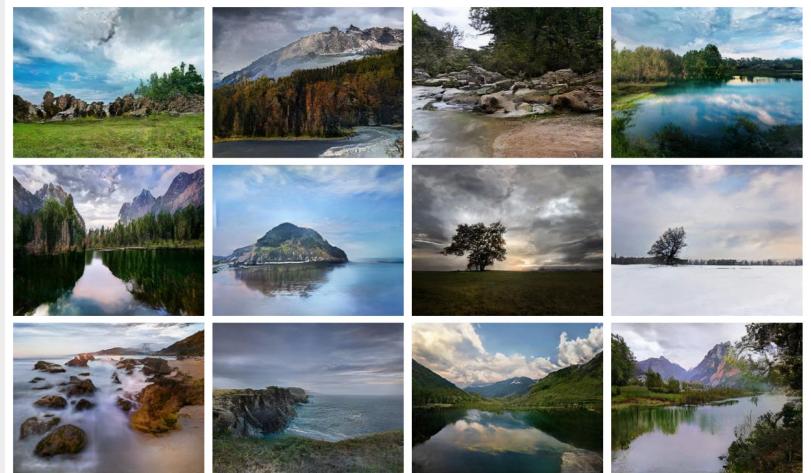
# Geometry-consistency Constraint



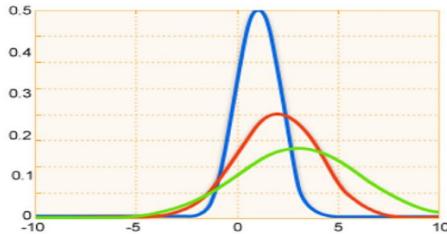
# Geometry-consistency Constraint



# Advanced methods



# Style transfer

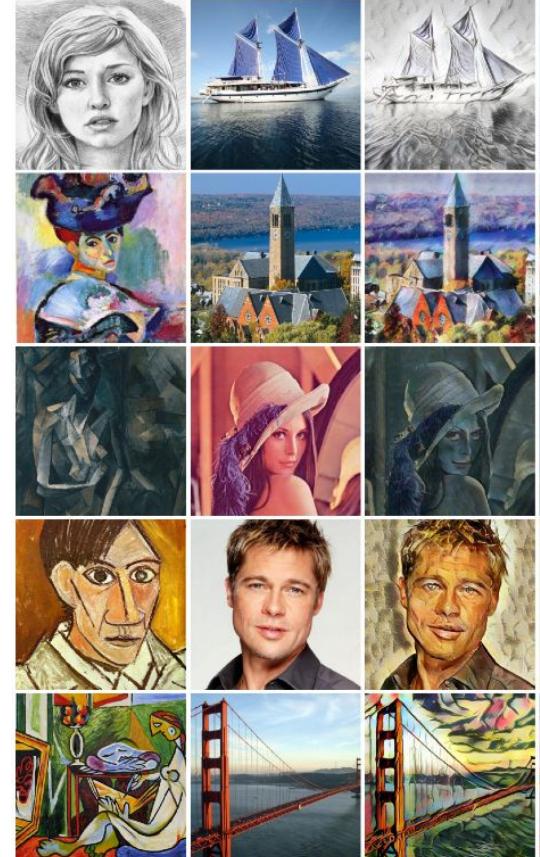


(a) Content image.



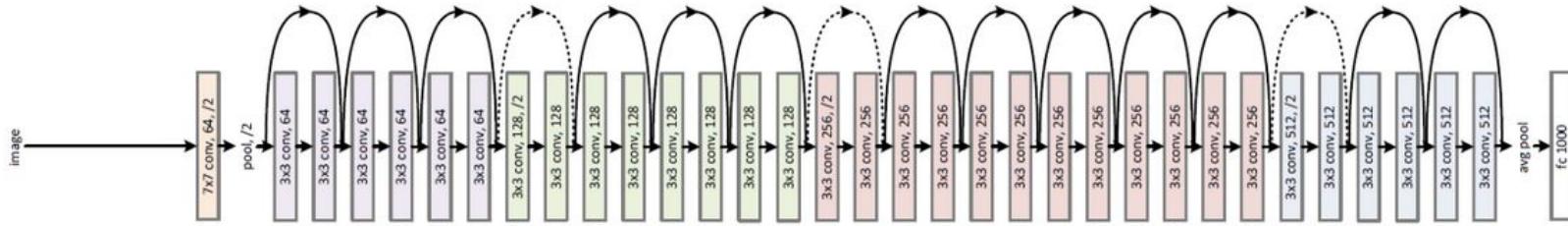
(c) Low contrast content image.

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

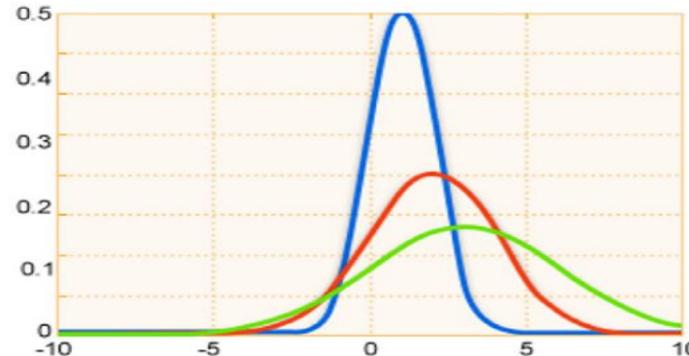


# Layers: Batch-normalization

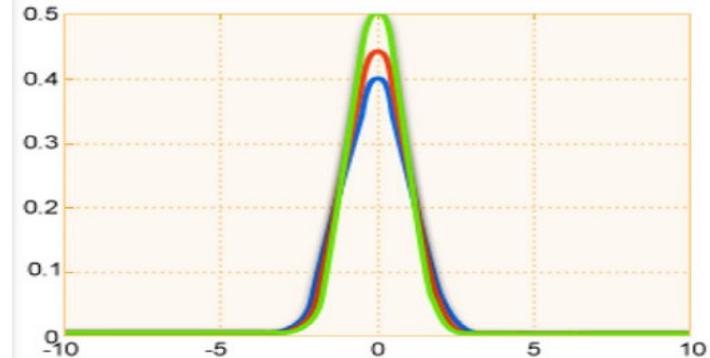
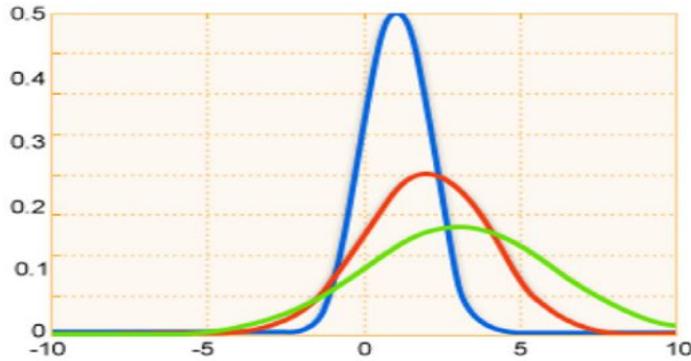
34-layer residual



Parameter update:  $\theta_{n+1} = \theta_n - \gamma \nabla L_b(\theta_n)$



# Layers: Batch-normalization



**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;  
Parameters to be learned:  $\gamma, \beta$   
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Batch Normalization: Accelerating Deep Network  
Training by Reducing Internal Covariate Shift  
Sergey Ioffe, Christian Szegedy

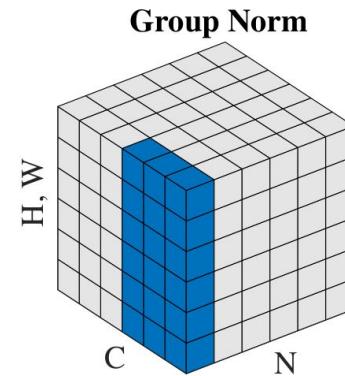
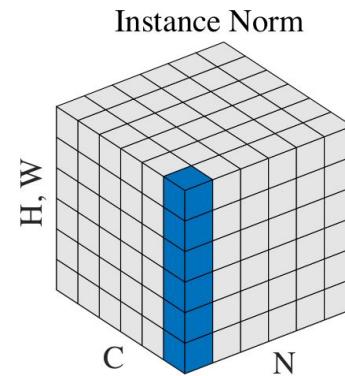
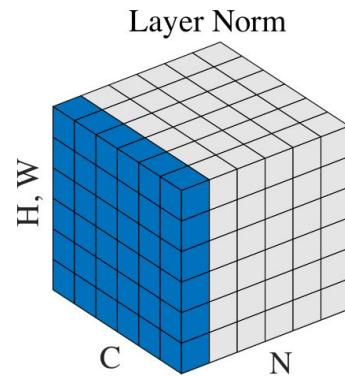
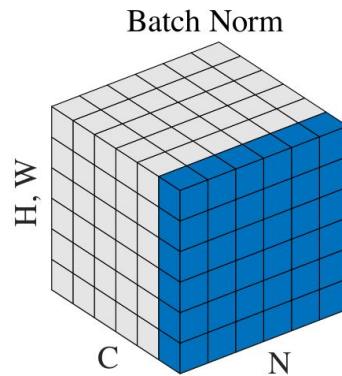
# Layers: Batch-normalization

Advantages:

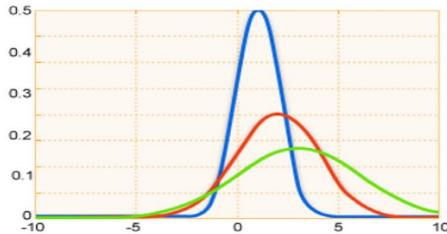
- Faster training
- Better optimization (lower final loss)
- Better generalization
- Specific usages in:
  - few-shot learning
  - domain adaptation
  - generation

# Layers: Other normalizations

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$



# Style transfer

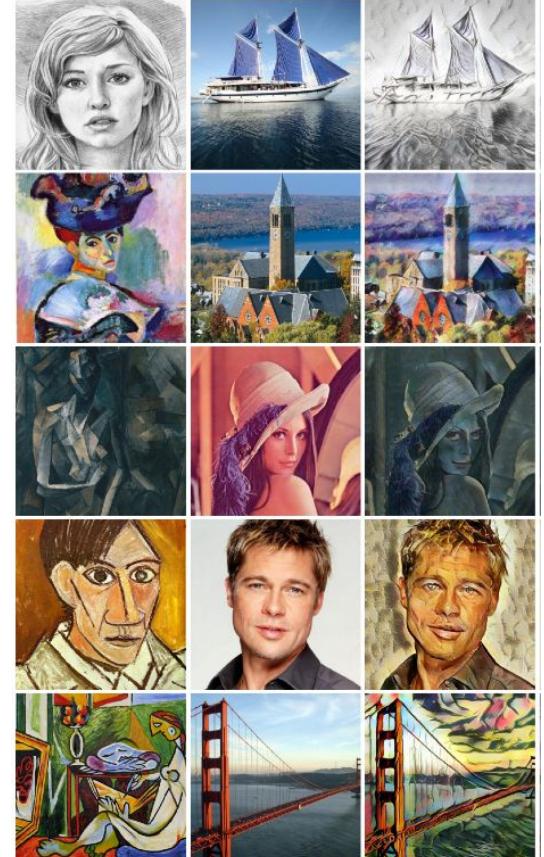


(a) Content image.

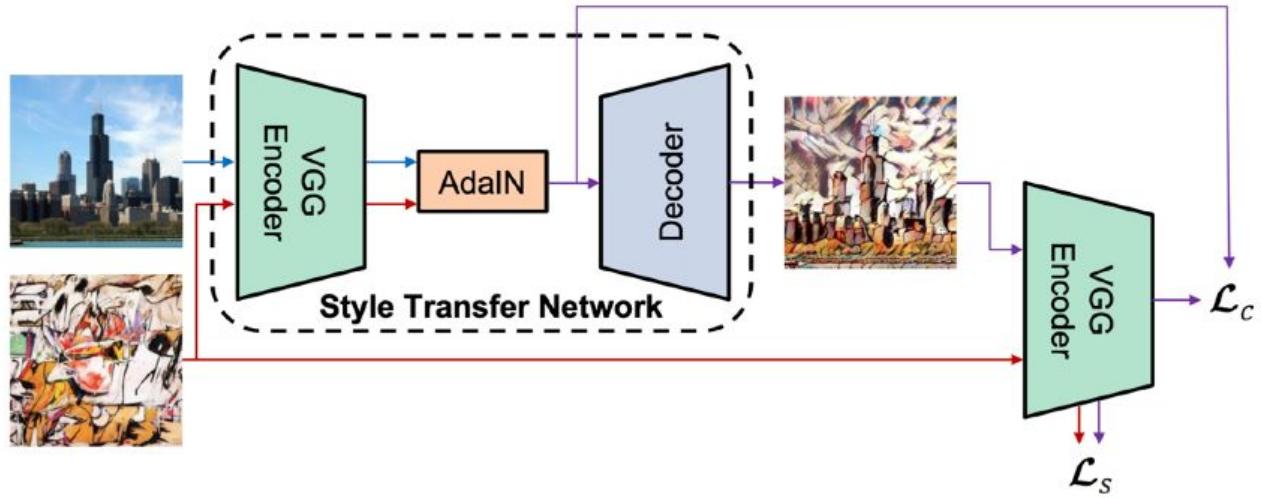


(c) Low contrast content image.

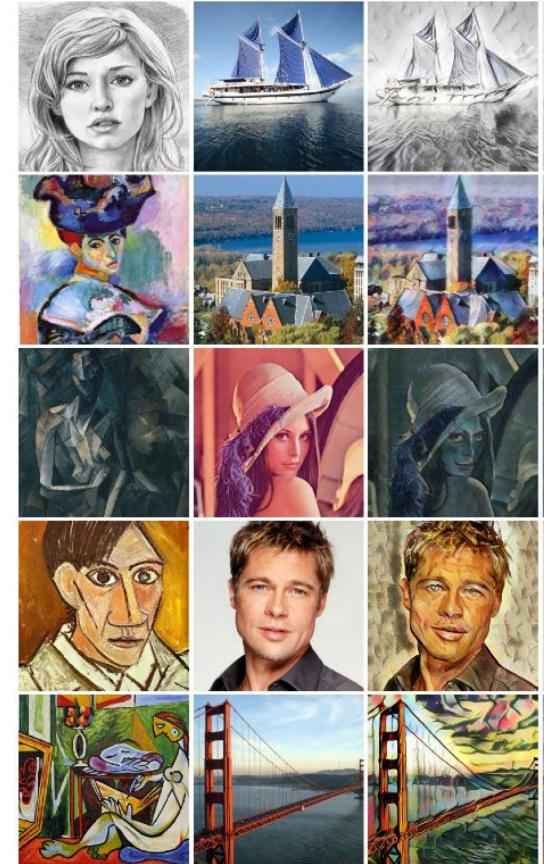
The two images are identical after instance normalization.  
We obtain style invariant representations.



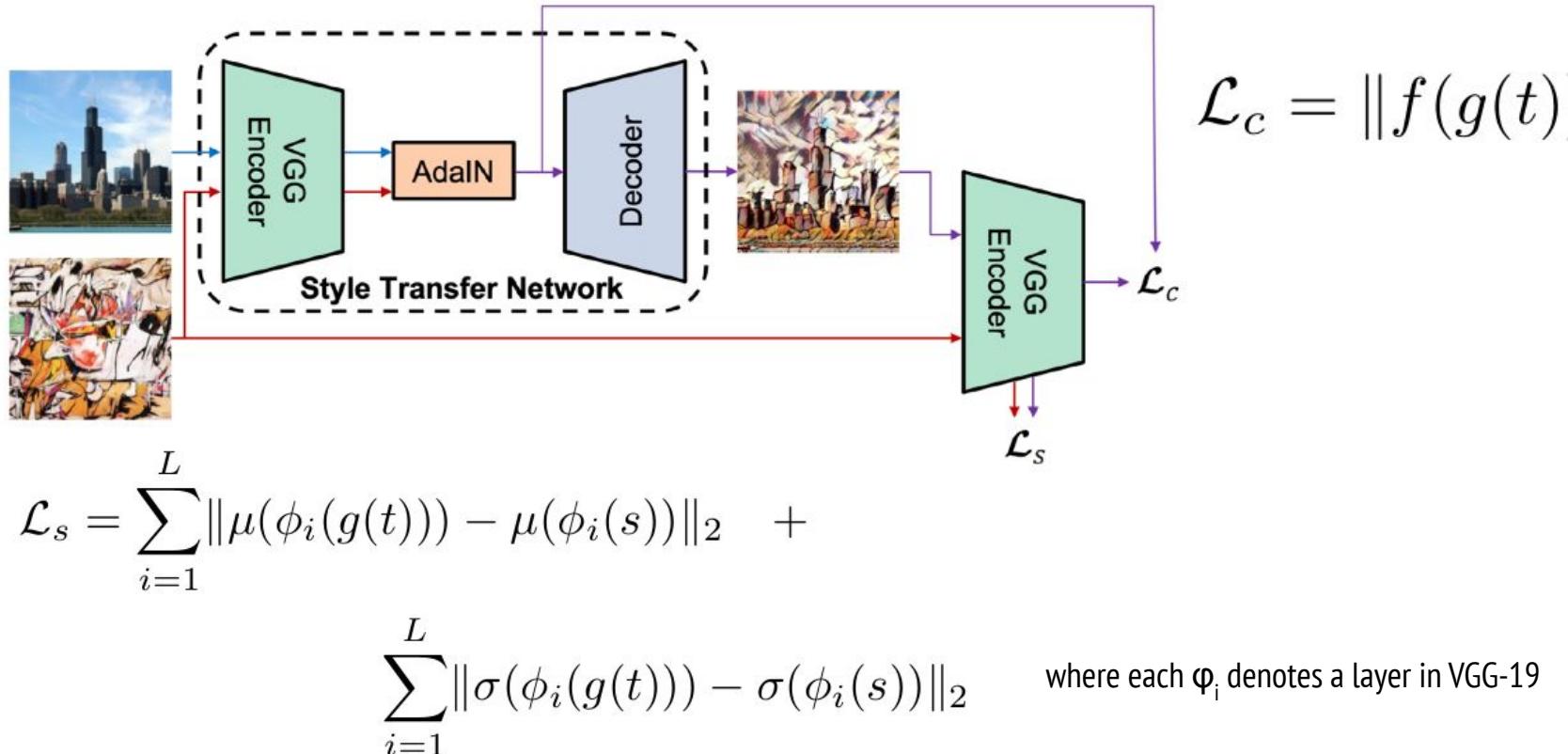
# Style transfer



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

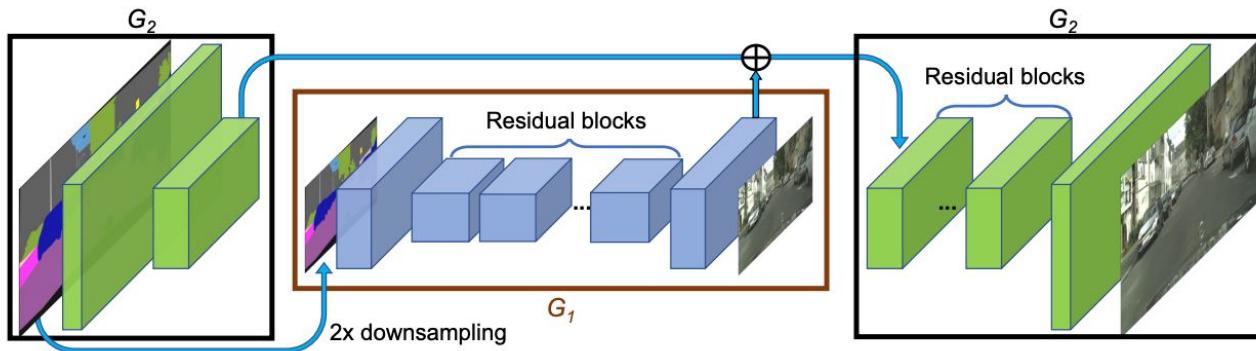


# Style transfer

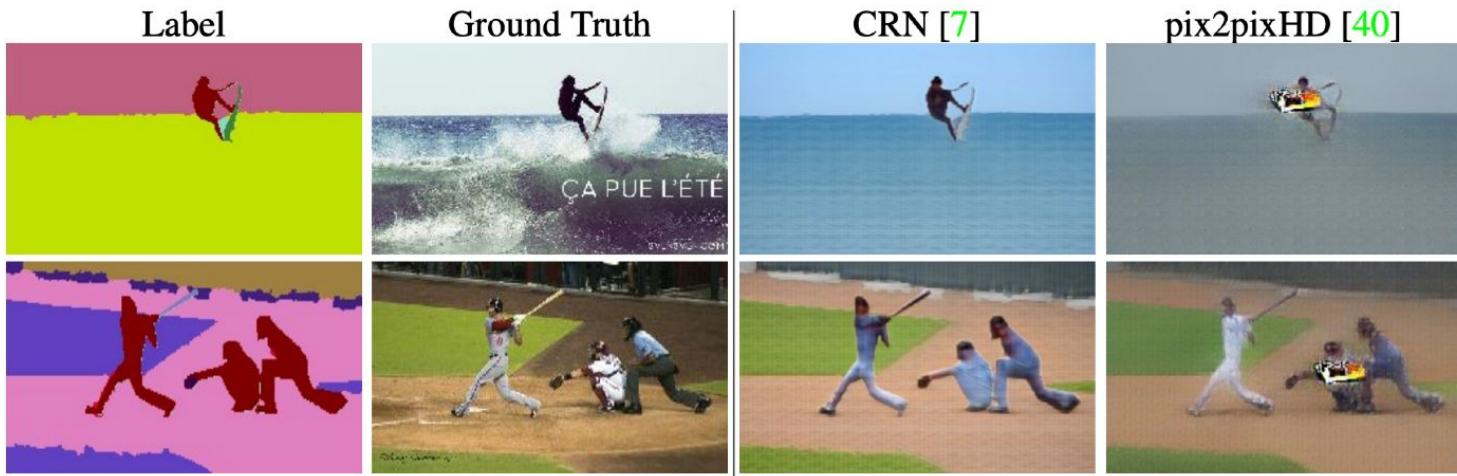


# Pix2PixHD

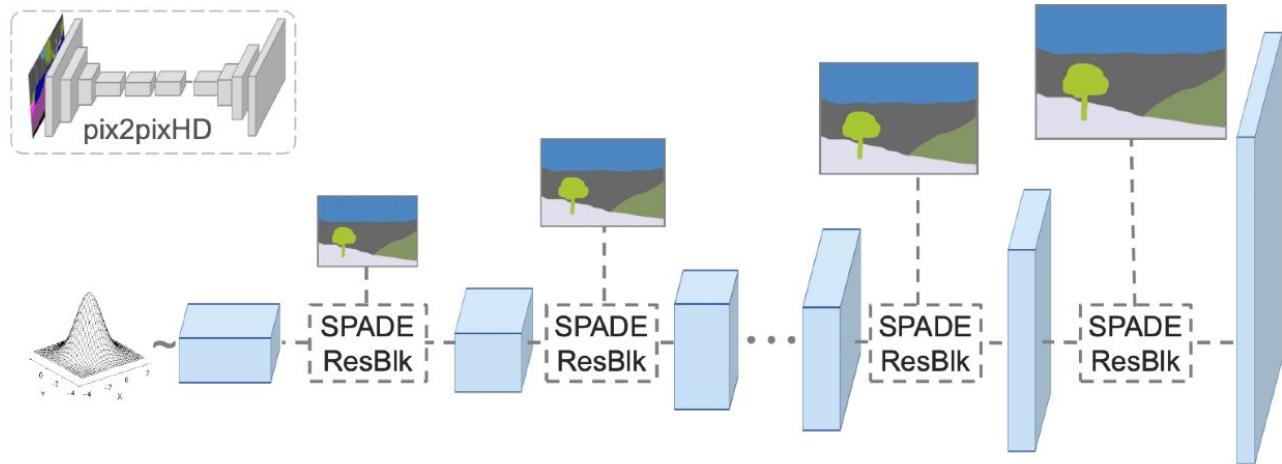
Idea: two-stage coarse-to-fine generation of HD images



# Pix2PixHD

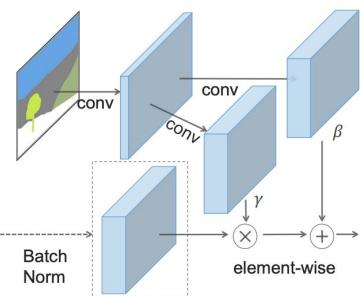
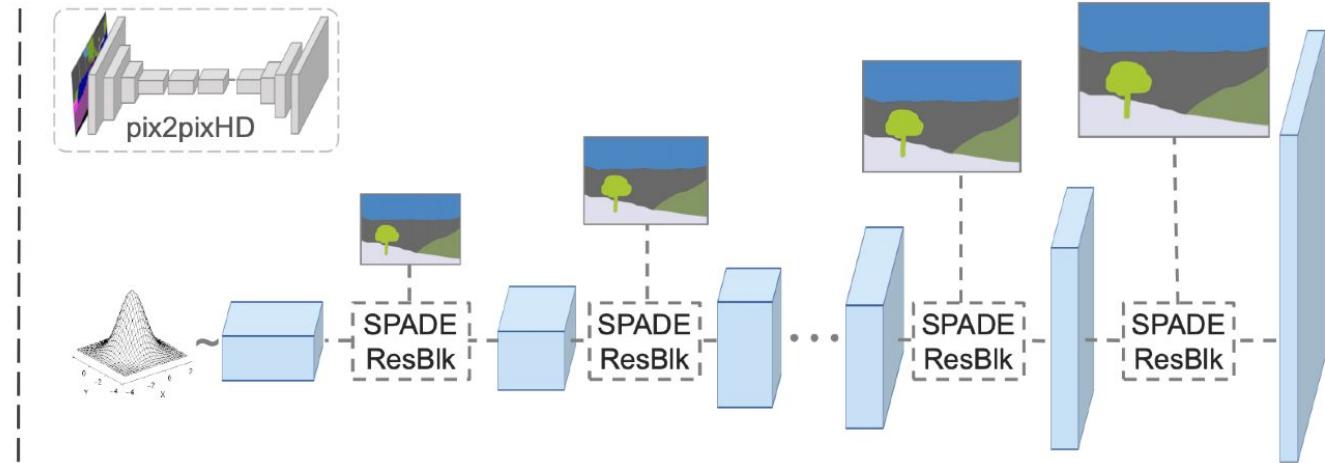
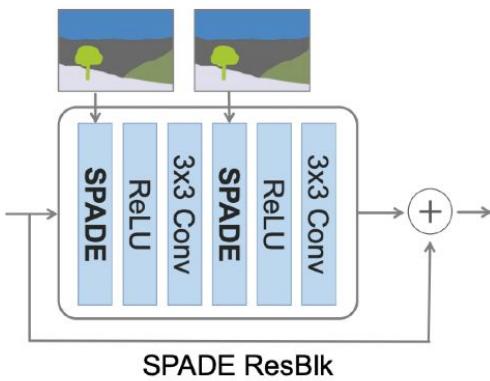


# Spade



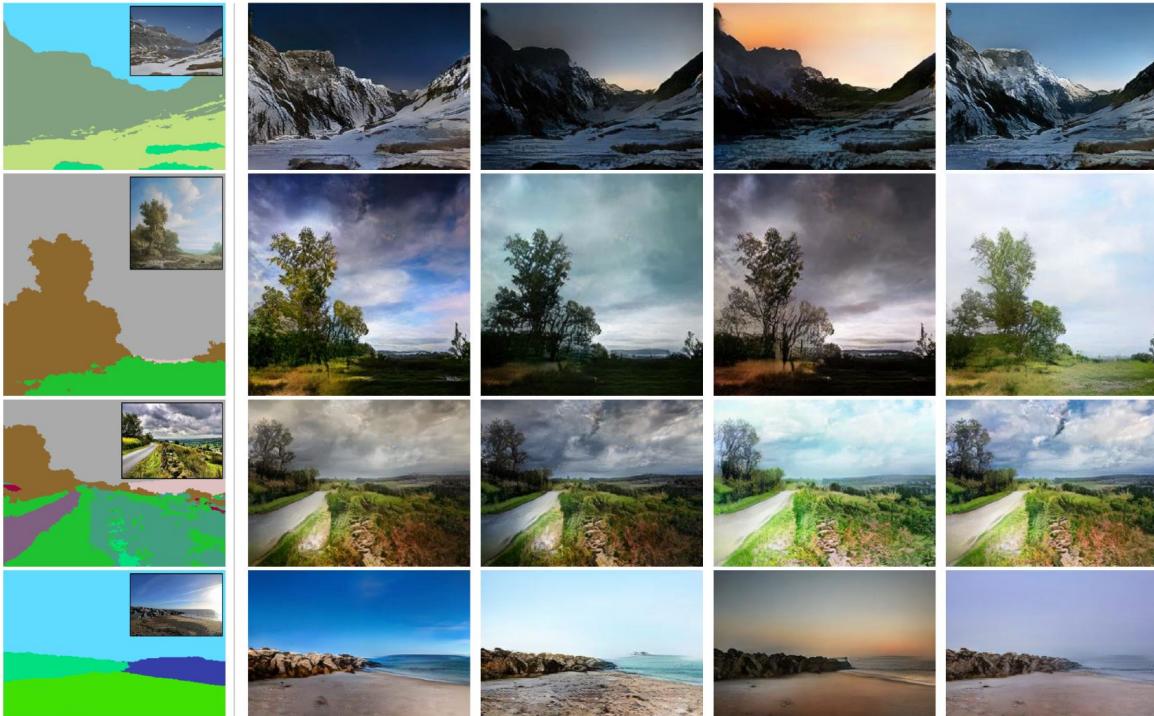
Park, Taesung, et al. "Semantic image synthesis with spatially-adaptive normalization." CVPR'2019

# Spade



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

# Spade

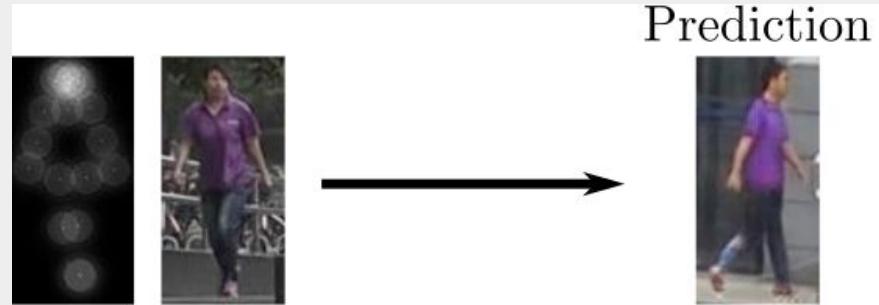


[Live demo](#)

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

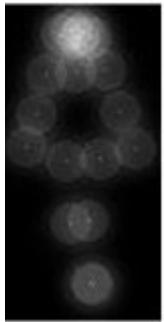
# Pose-guided Image Generation

Stéphane Lathuilière

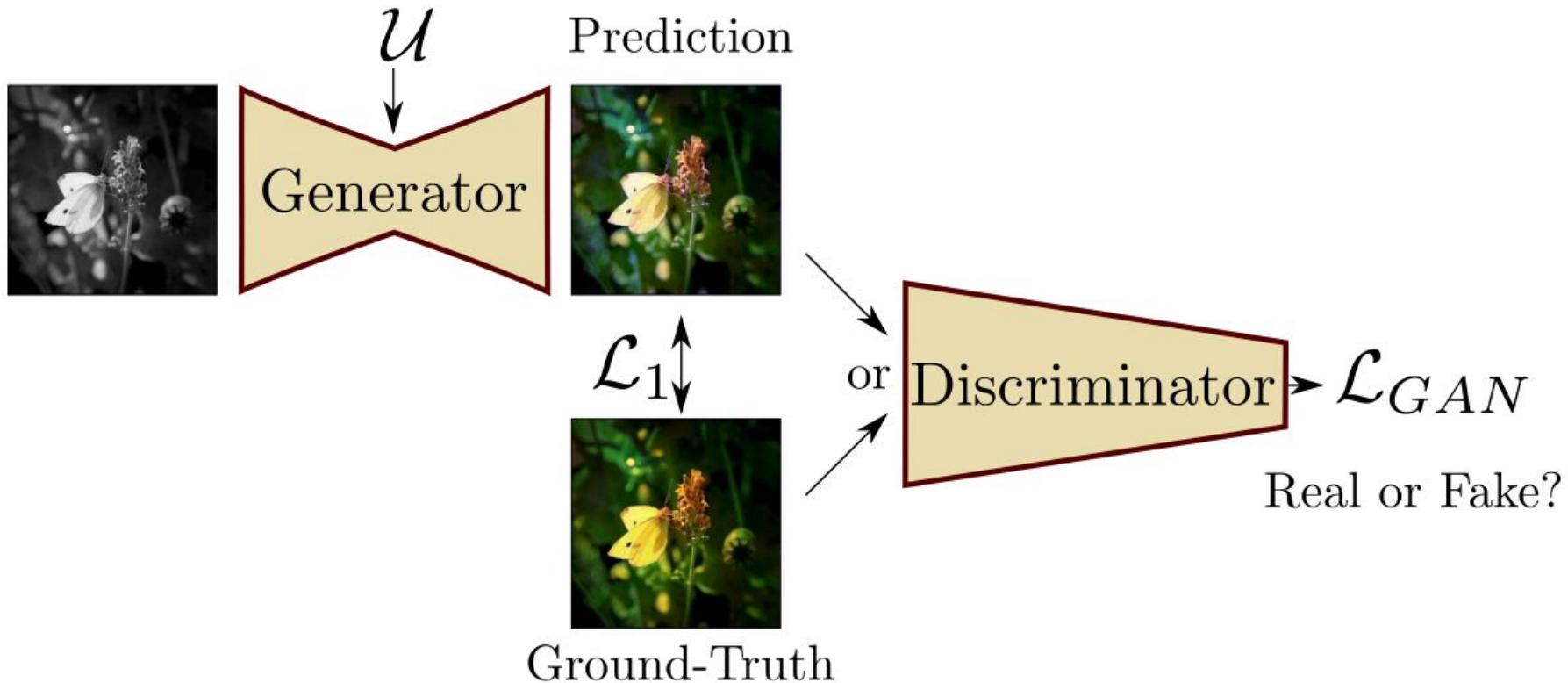


# Introduction

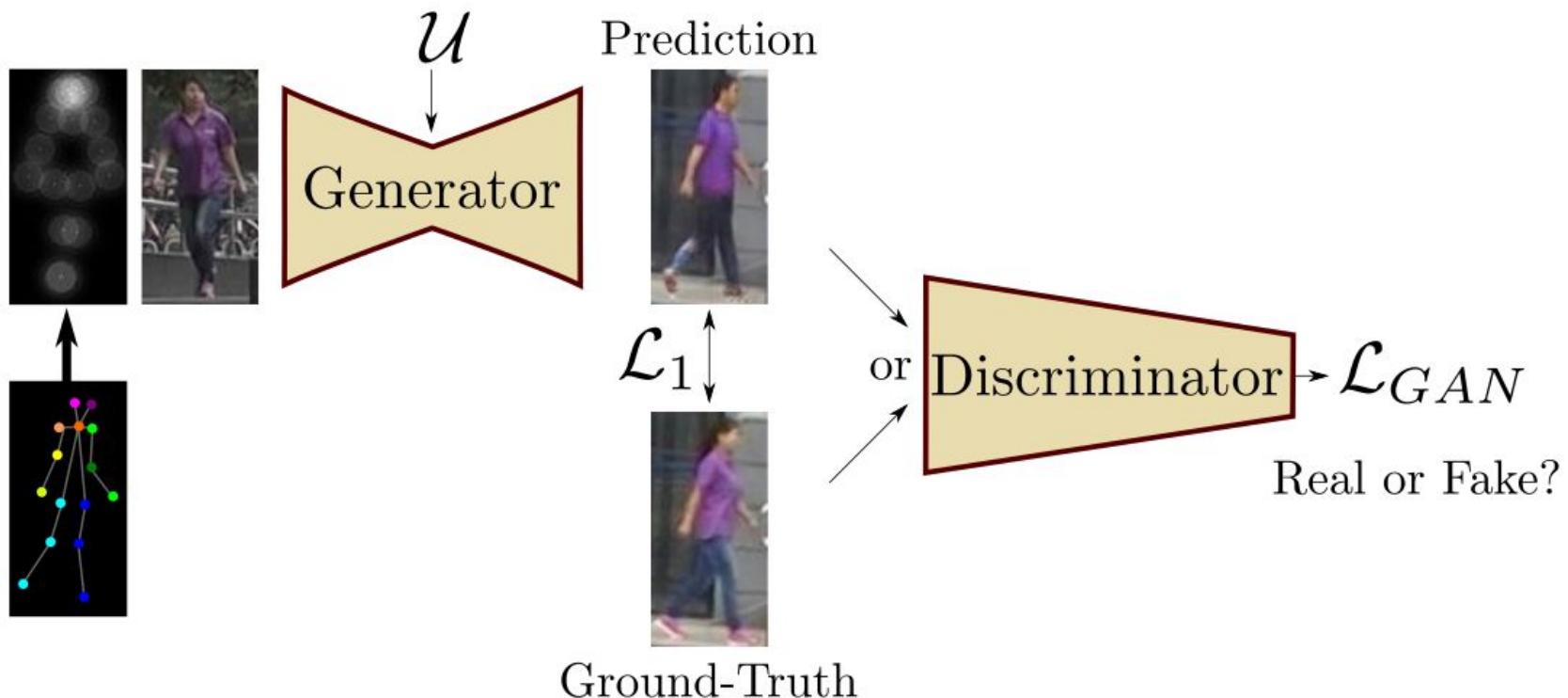
Prediction



# Pix2Pix



# Pix2Pix for pose guided



# Pix2Pix for pose guided



# Pose guided

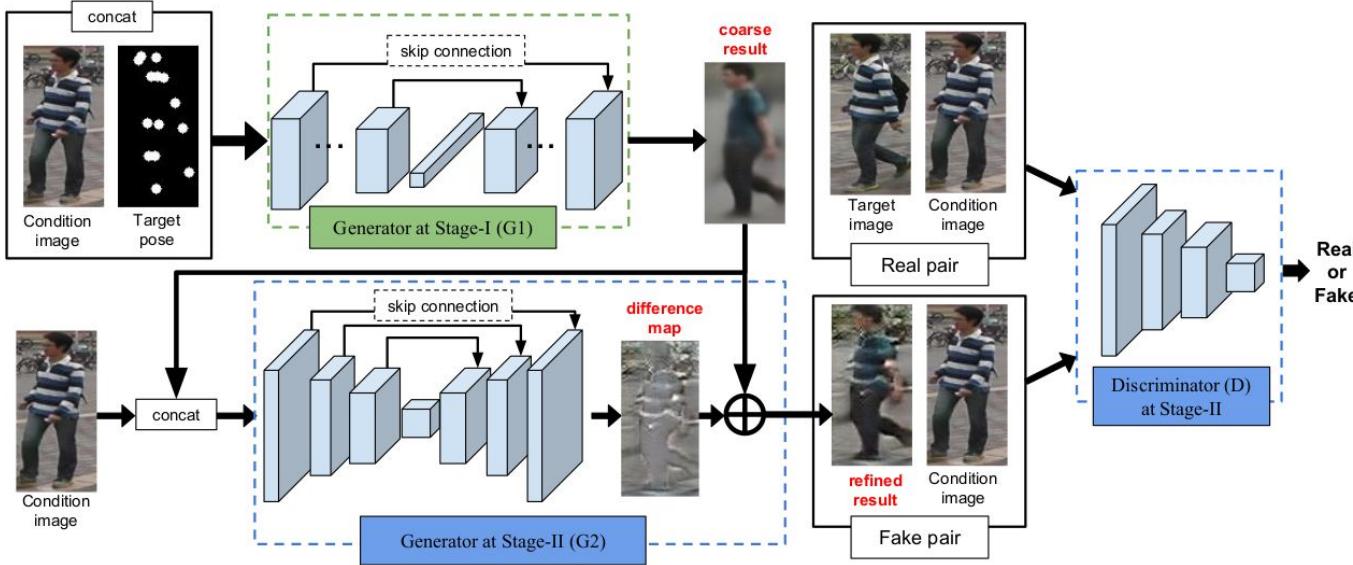
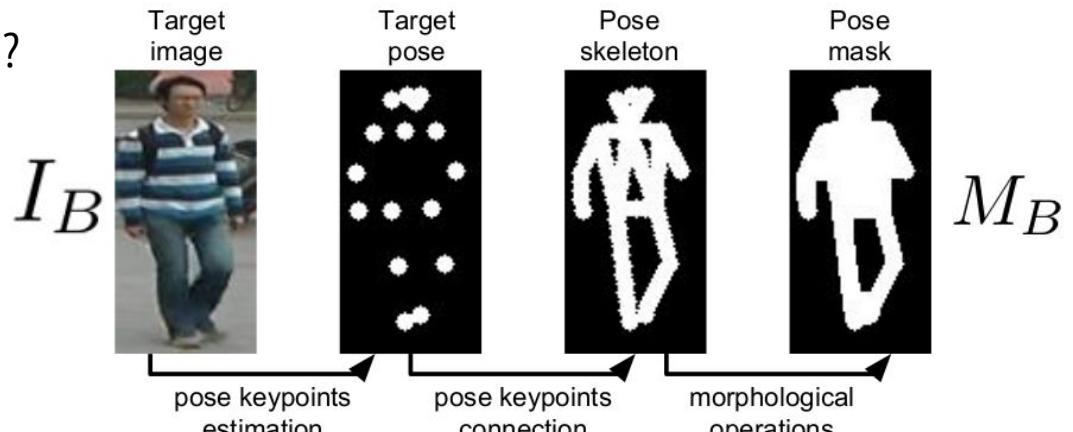


Figure 2: The overall framework of our Pose Guided Person Generation Network (PG<sup>2</sup>). It contains two stages. Stage-I focuses on pose integration and generates an initial result that captures the global structure of the human. Stage-II focuses on refining the initial result via adversarial training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

# Pose guided

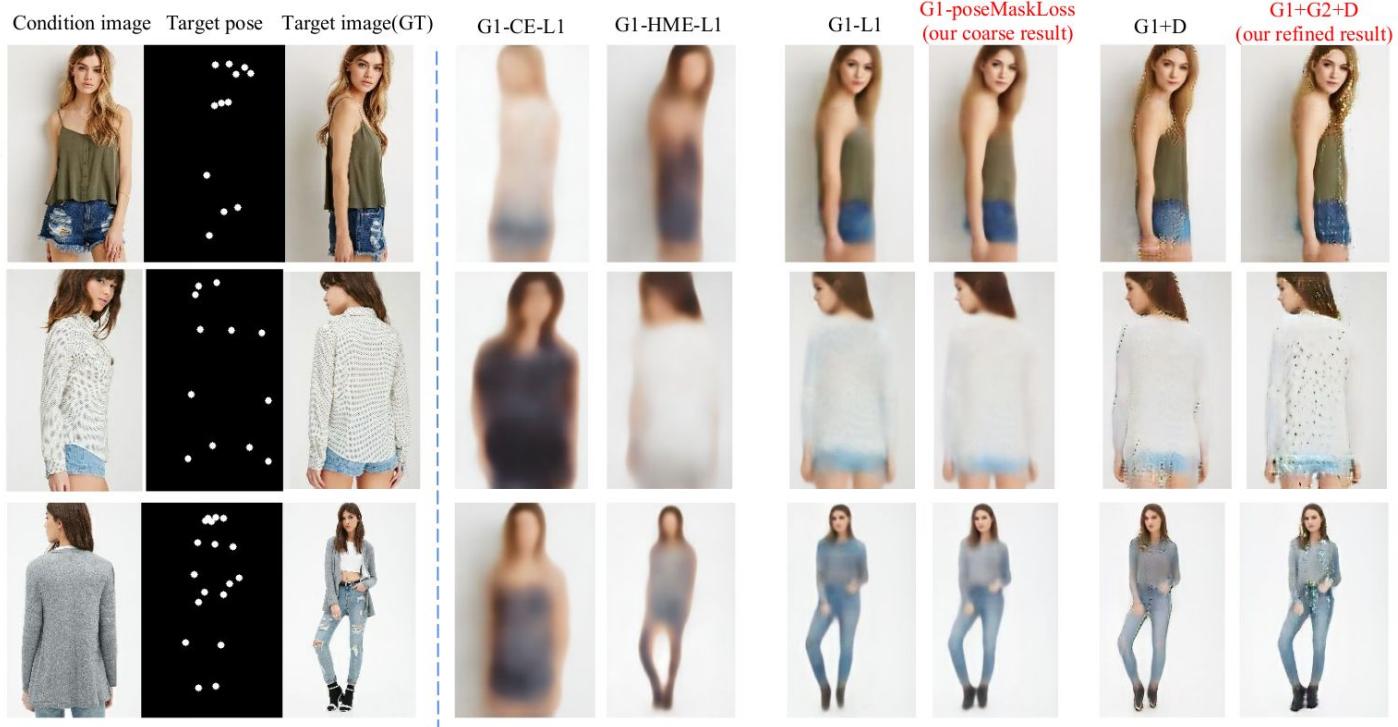
What about the background?



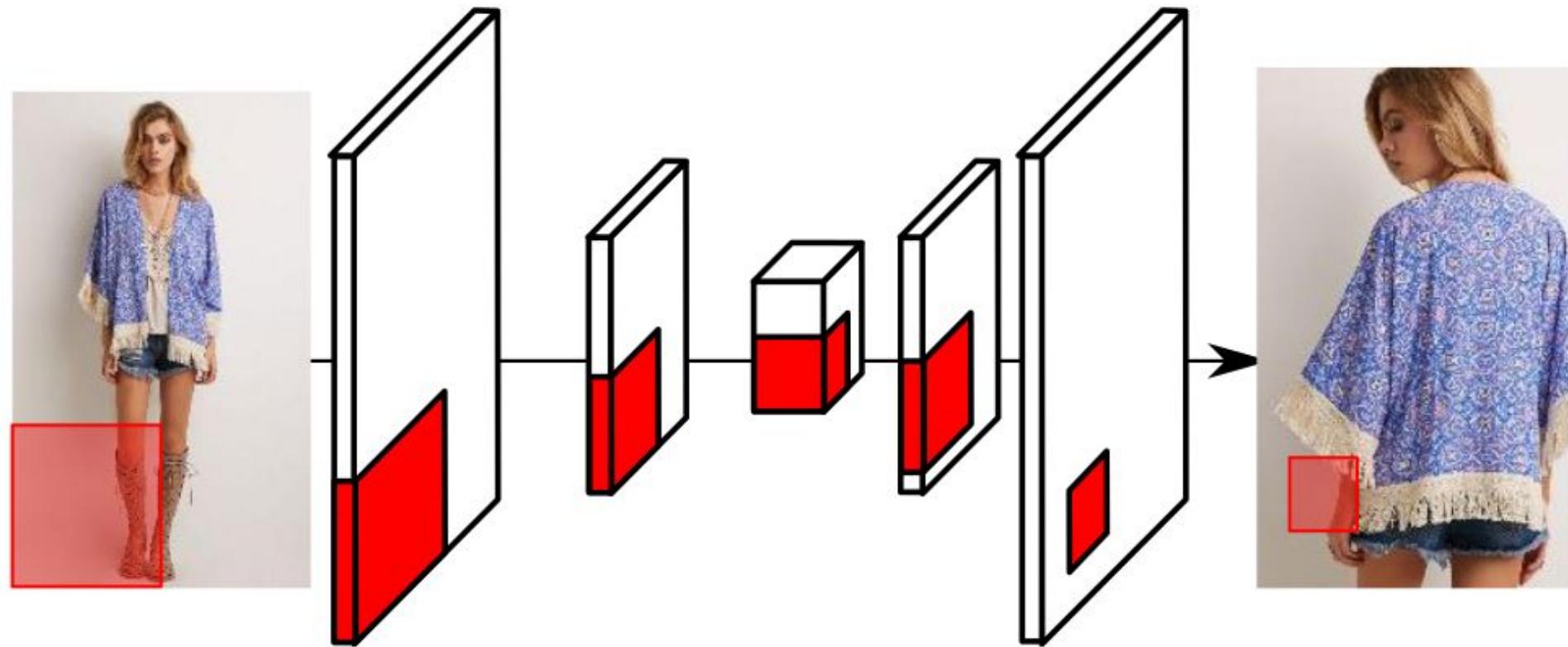
They give more weight to the foreground:

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

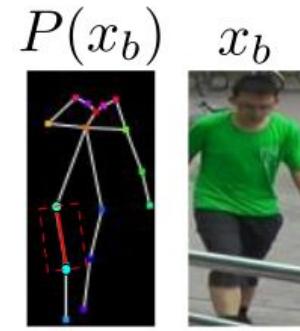
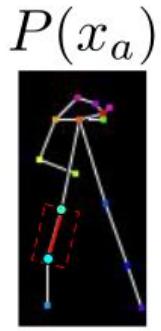
# Pose guided



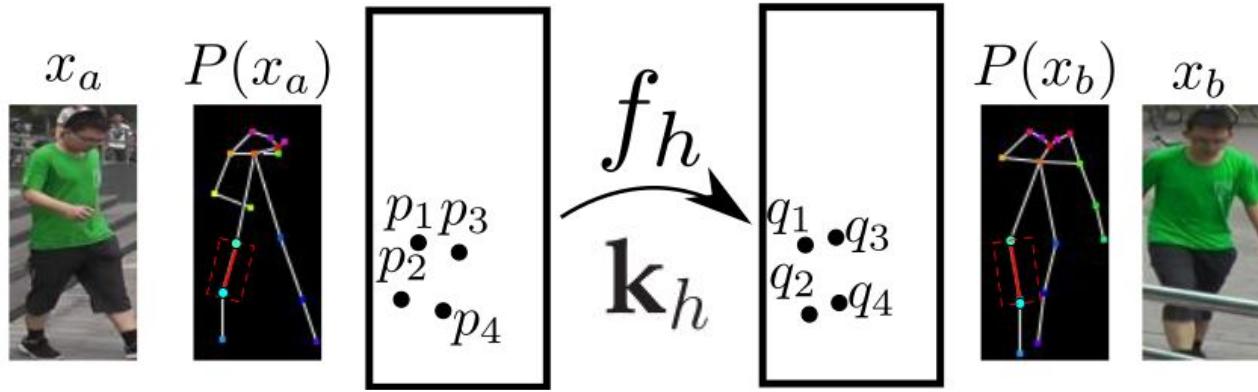
# Pix2Pix for pose guided: problem



# Pose guided: deformable GAN

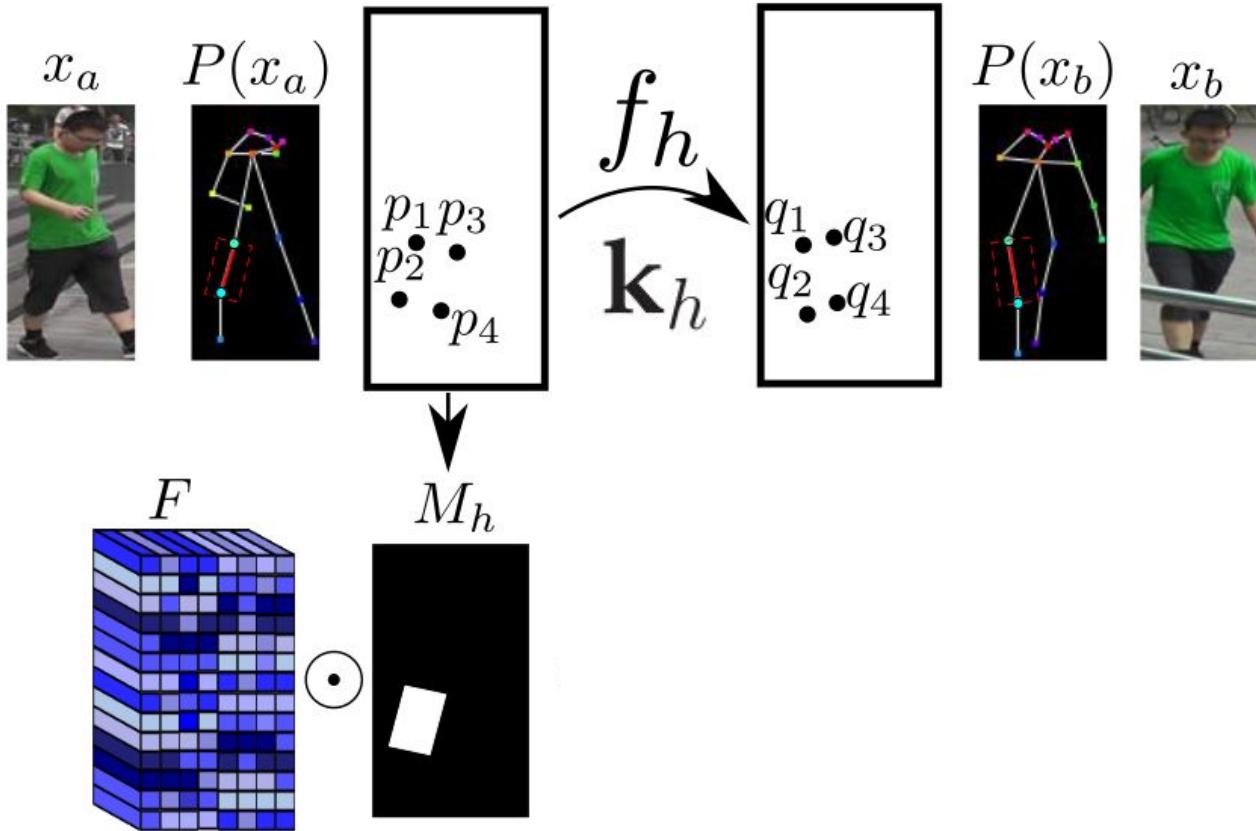


# Pose guided: deformable GAN

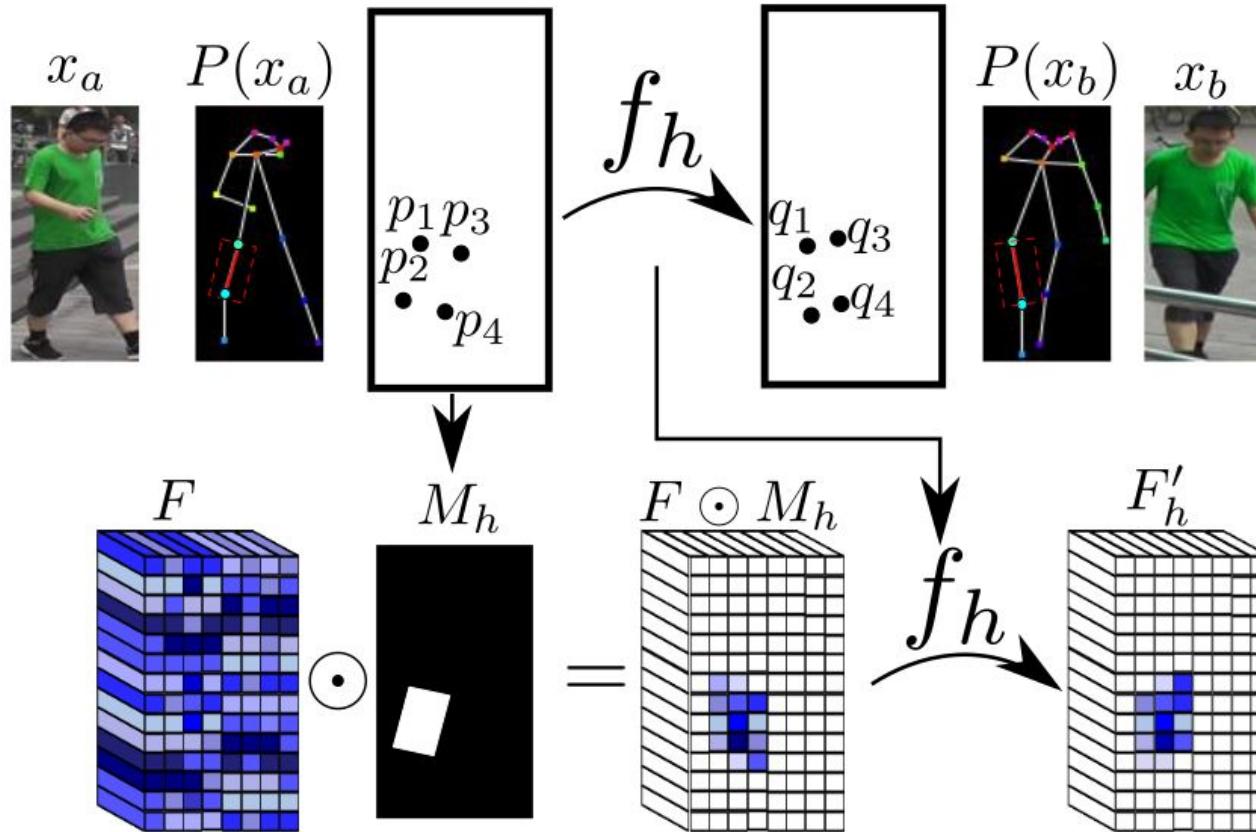


$$\min_{\mathbf{k}_h} \sum_{\mathbf{p}_j \in R_h^a, \mathbf{q}_j \in R_h^b} \|\mathbf{q}_j - f_h(\mathbf{p}_j; \mathbf{k}_h)\|_2^2.$$

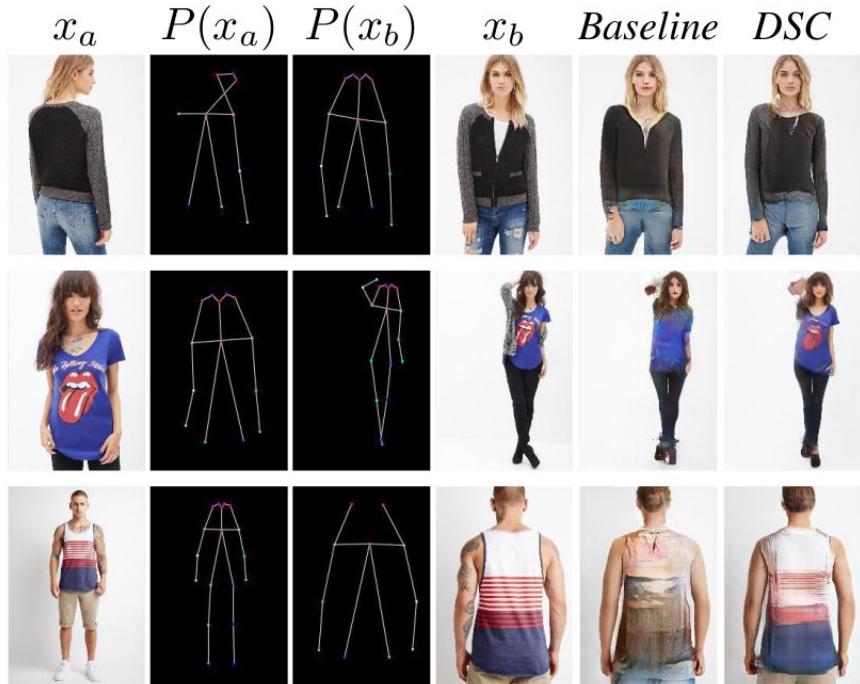
# Pose guided: deformable GAN



# Pose guided: deformable GAN



# Pose guided: deformable GAN



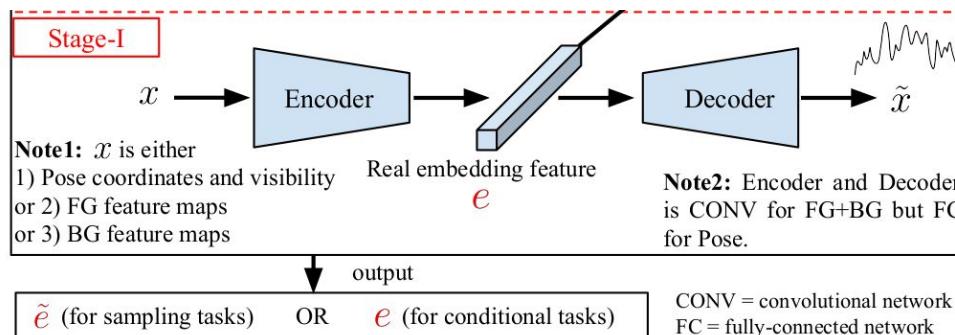
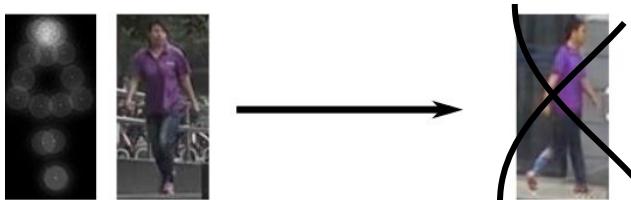
# Pose guided: reconstruction vs GAN

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

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



# Pose guided: Disentangled Generation



# Pose guided: Disentangled Generation

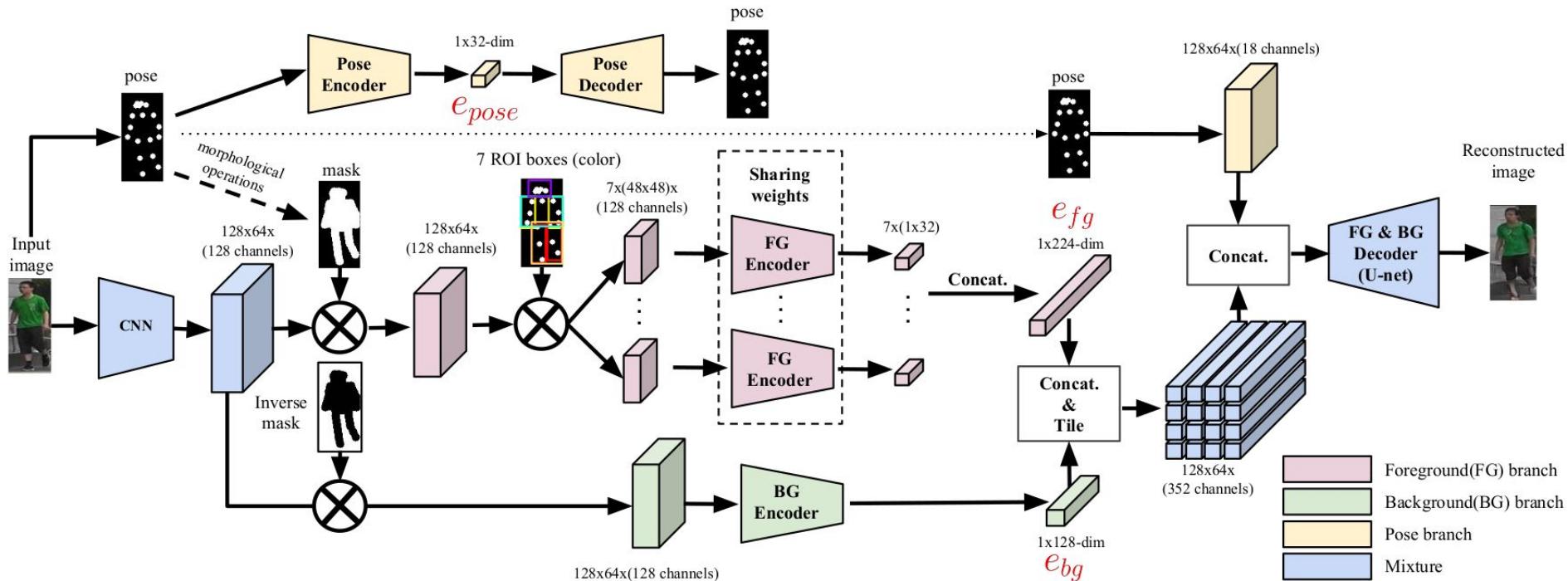


Figure 3: Stage-I: disentangled image reconstruction. This framework is composed of three branches: foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features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

# Pose guided: Disentangled Generation

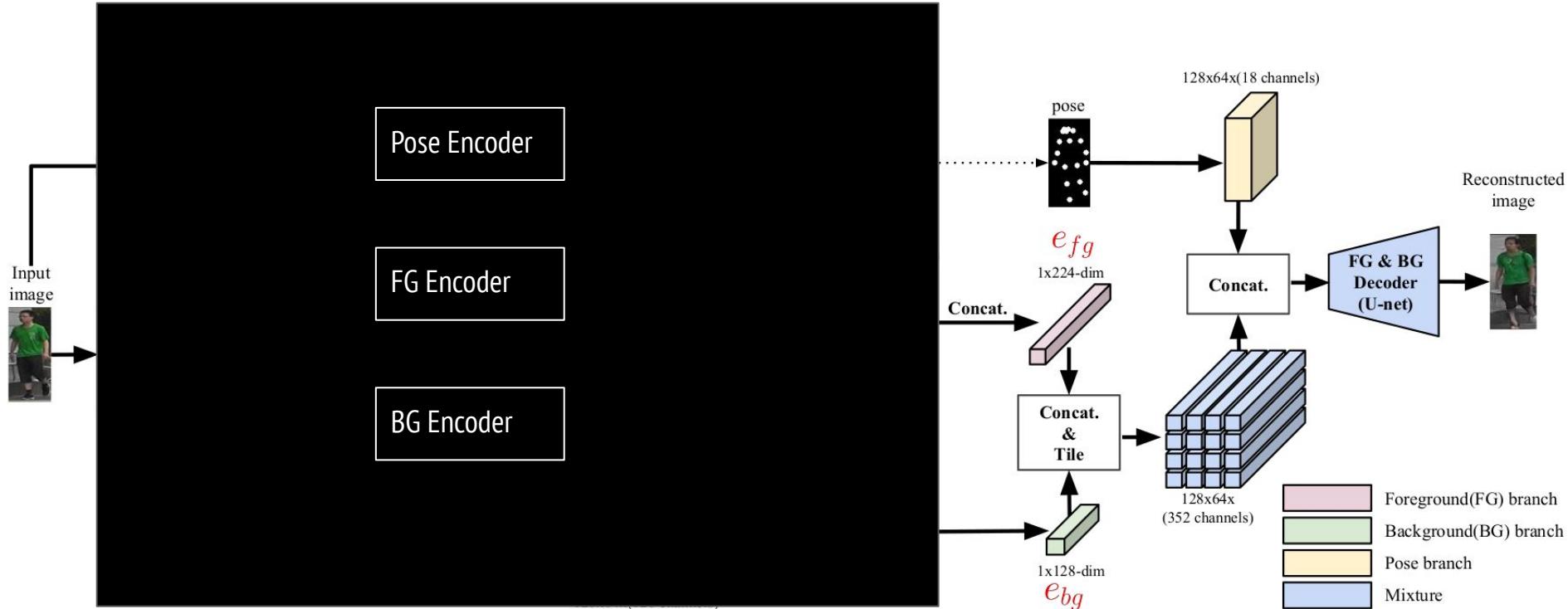
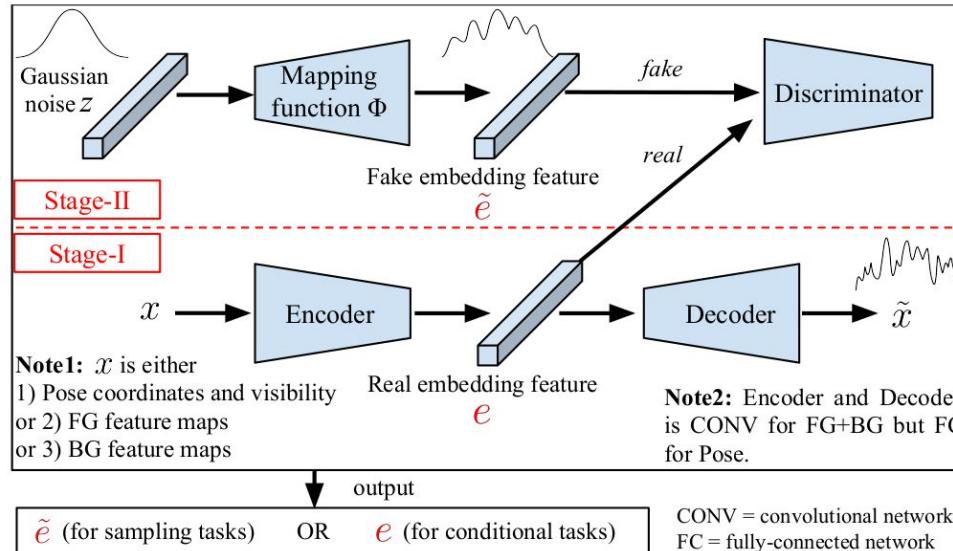


Figure 3: Stage-I: disentangled image reconstruction. This framework is composed of three branches: foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features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

# Pose guided: Disentangled Generation



# Pose guided: Disentangled Generation

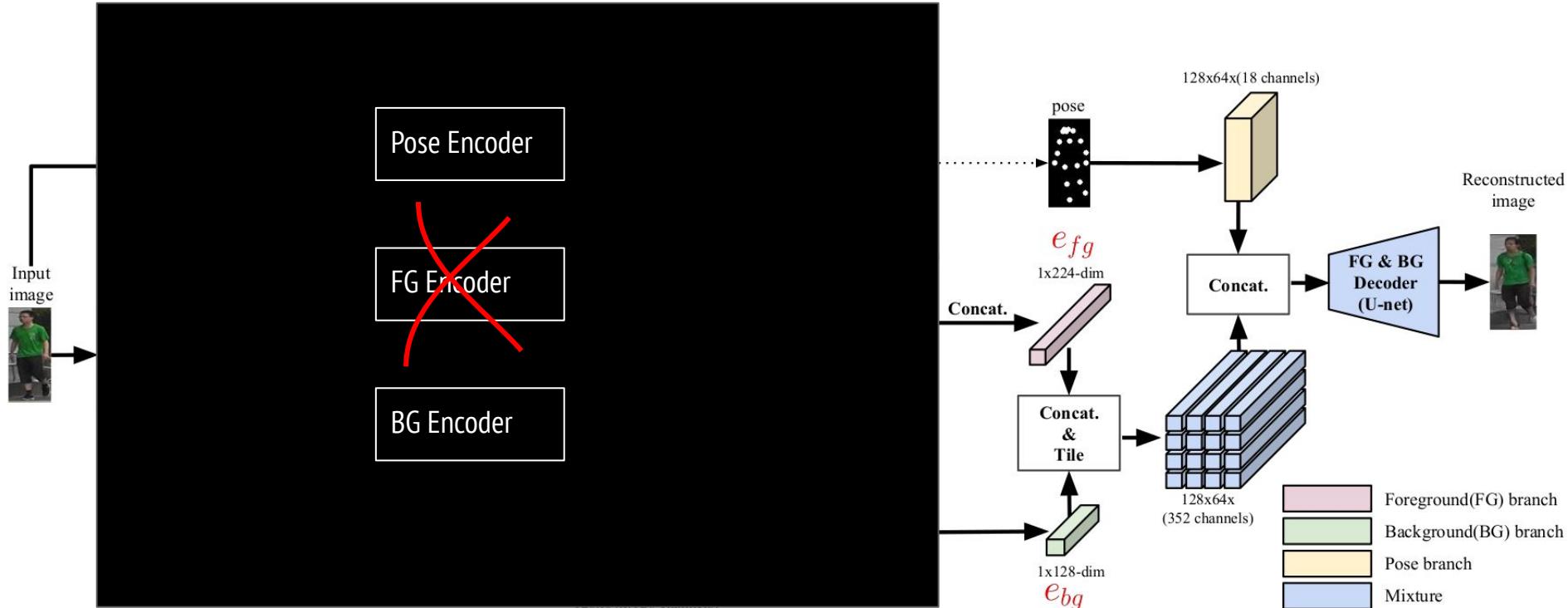


Figure 3: Stage-I: disentangled image reconstruction. This framework is composed of three branches: foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features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

# Pose guided: Disentangled Generation

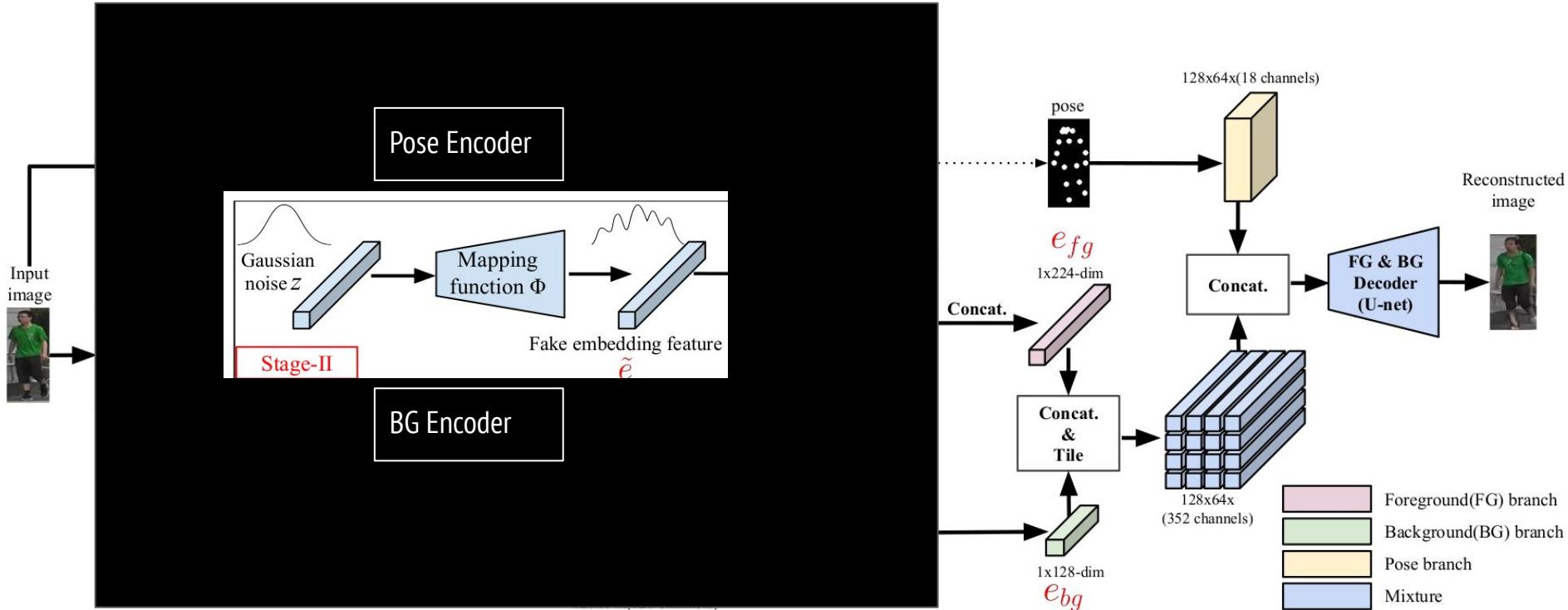


Figure 3: Stage-I: disentangled image reconstruction. This framework is composed of three branches: foreground, background and pose. Note that we use a fully-connected auto-encoder network to reconstruct the pose (incl. keypoint coordinates and visibility), so that we can decode the embedded pose features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

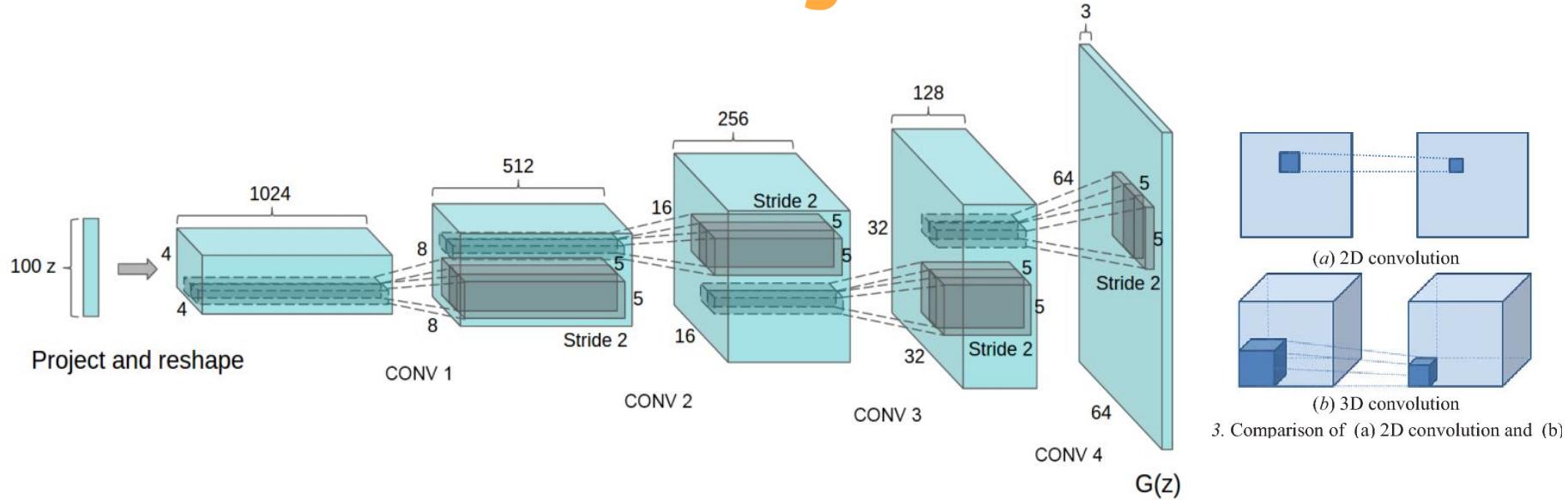


FG, BG and Pose sampling

# Video generation

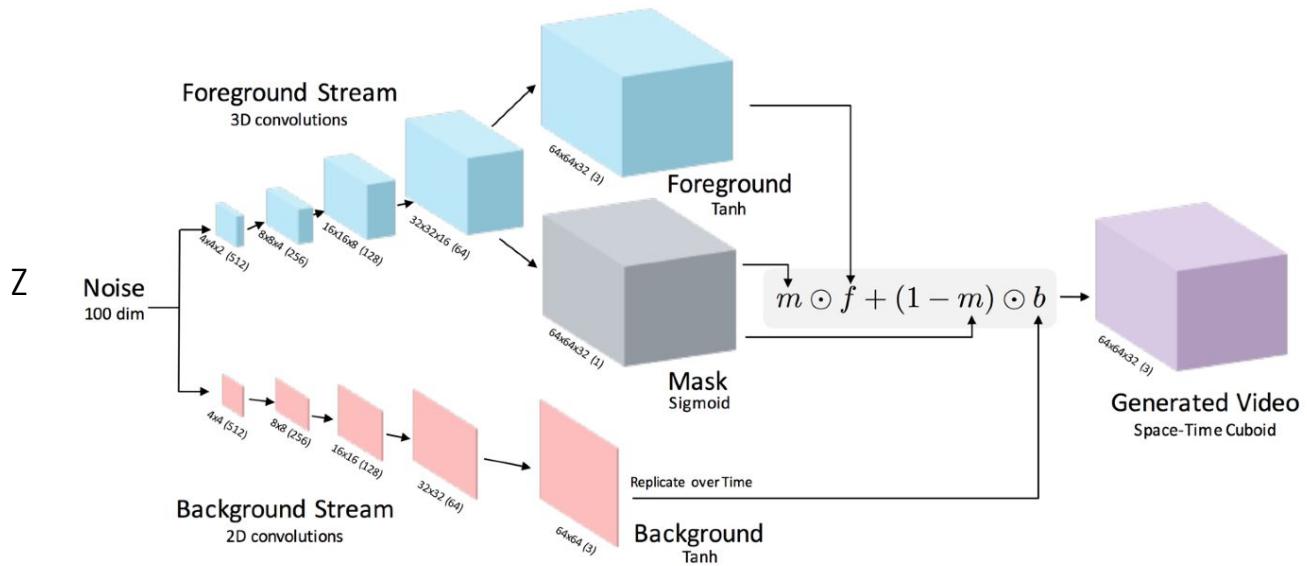


# From DCGAN to video generation



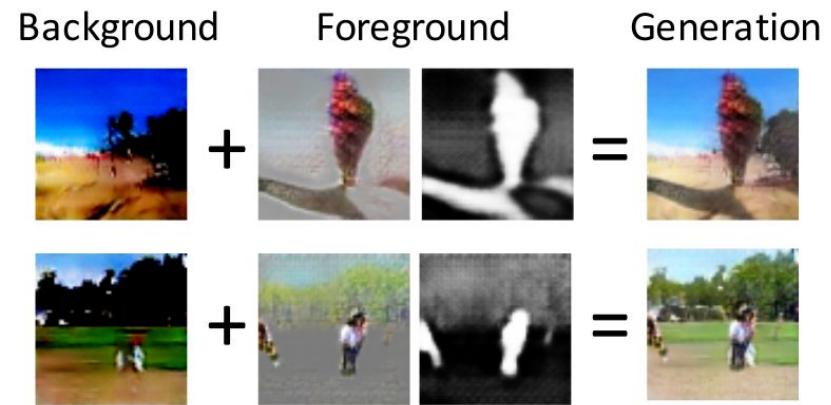
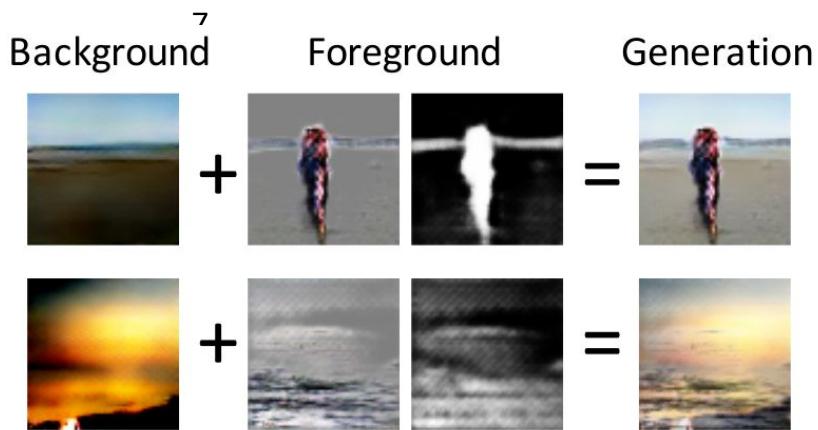
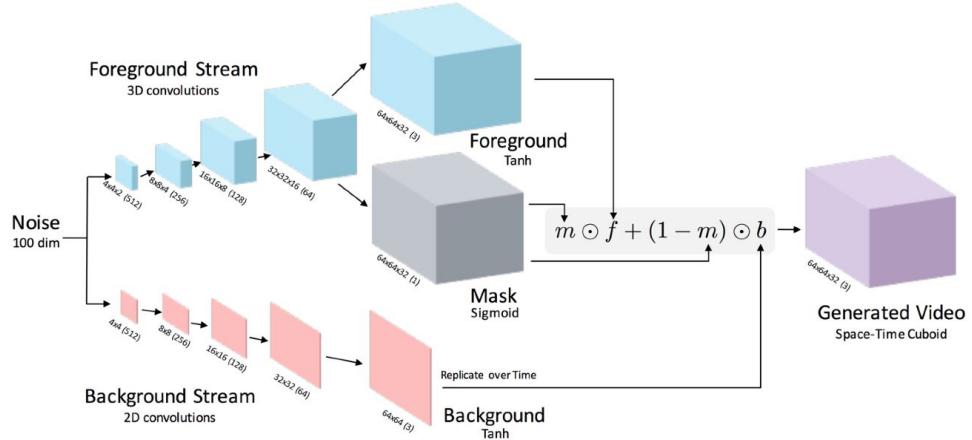
A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, ICLR 2016  
Figure from Lung nodule detection based on 3D convolutional neural networks Lei Fan et al.

# Video Generation



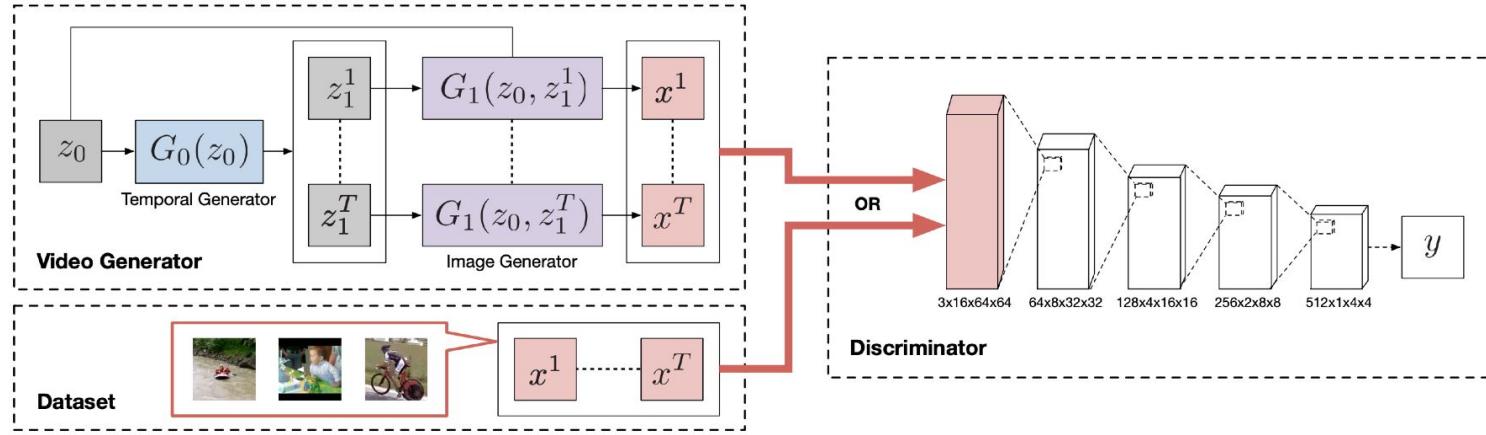
Saito, Matsumoto, and Saito. "Temporal generative adversarial nets with singular value clipping." ICCV'2017.

# Video Generation



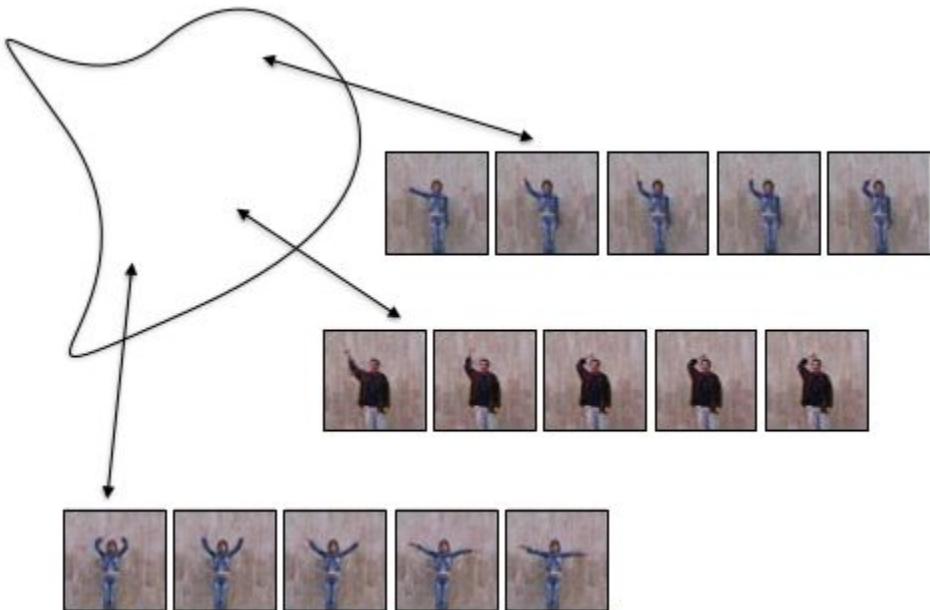
Saito, Matsumoto, and Saito. "Temporal generative adversarial nets with singular value clipping." ICCV'2017.

# Video Generation



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

# Video Generation



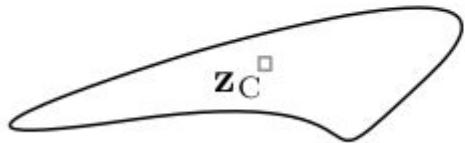
## Limitations:

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

# Video Generation: MoCoGAN

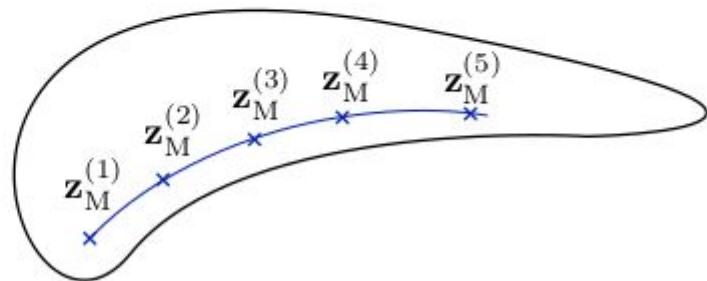
Sampled content

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



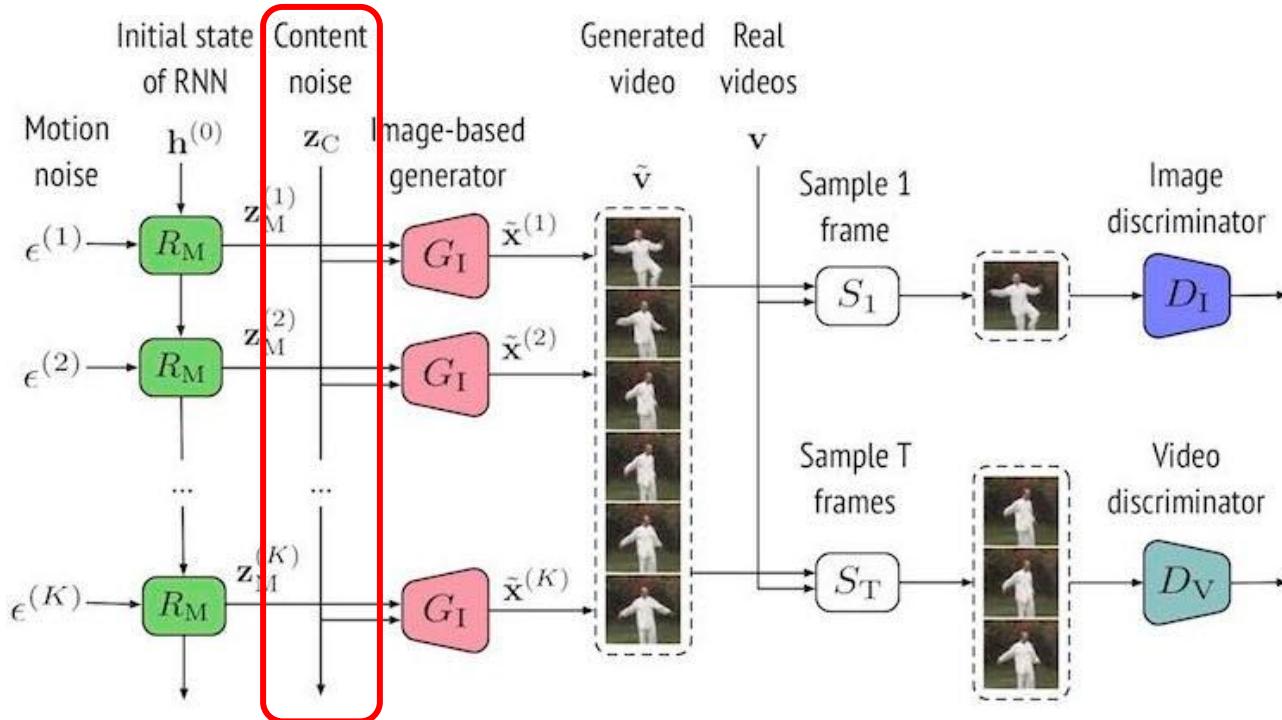
Motion trajectory

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



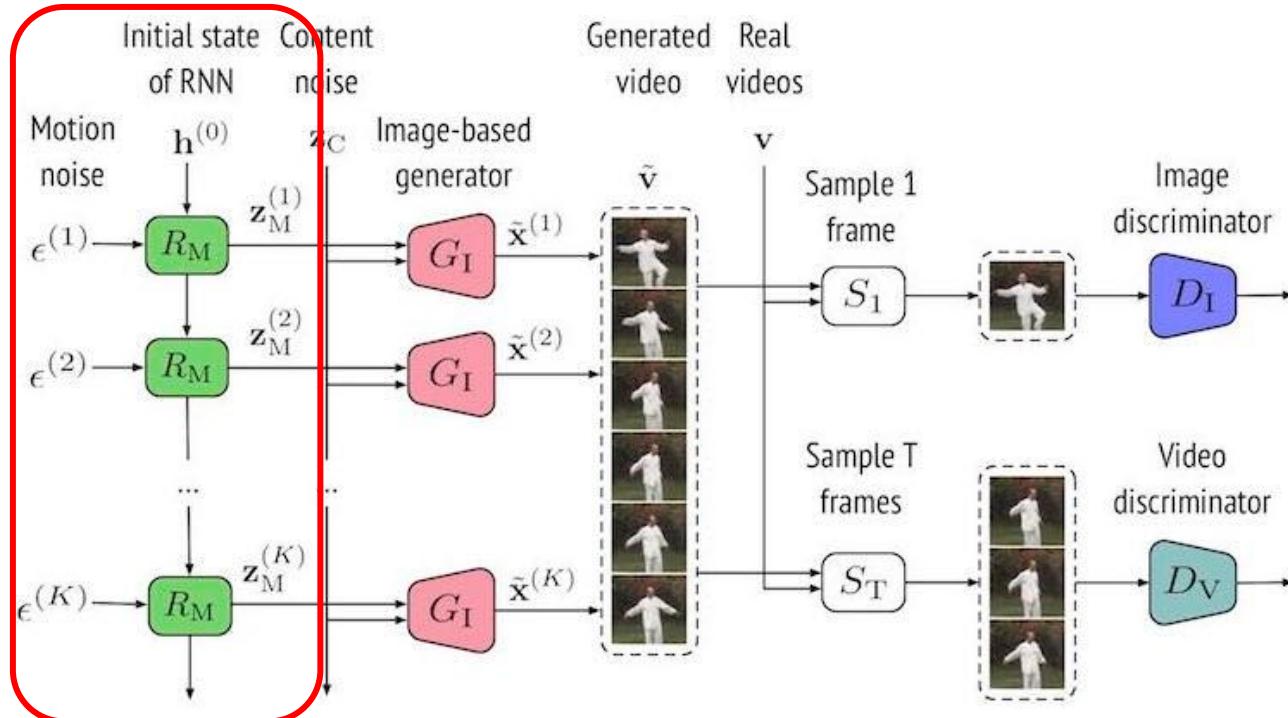
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

# Video Generation: MoCoGAN



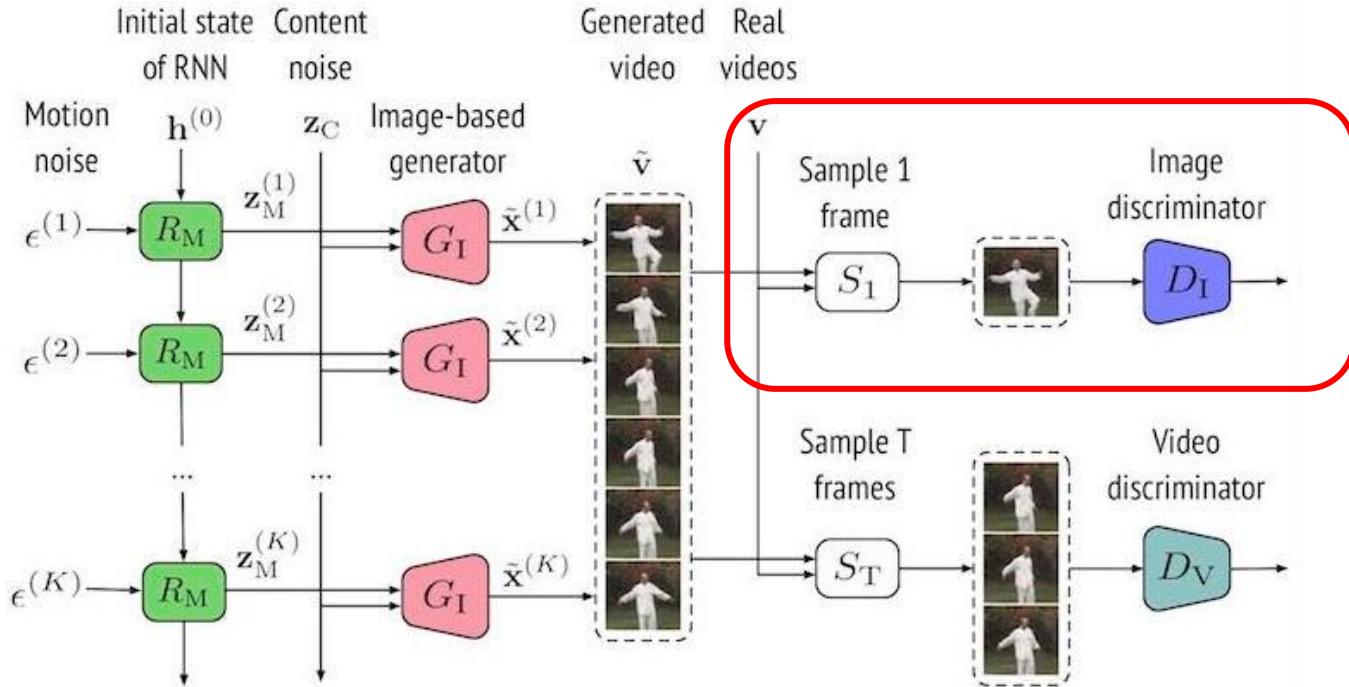
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

# Video Generation: MoCoGAN



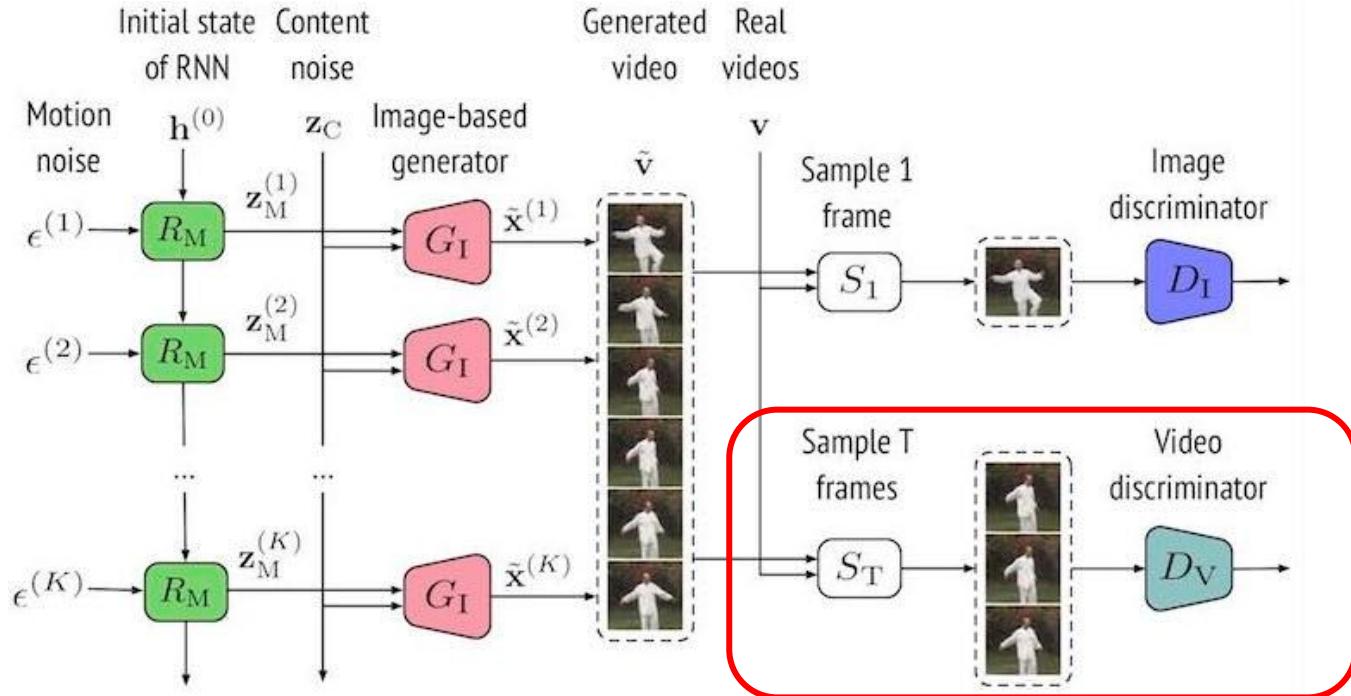
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

# 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



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

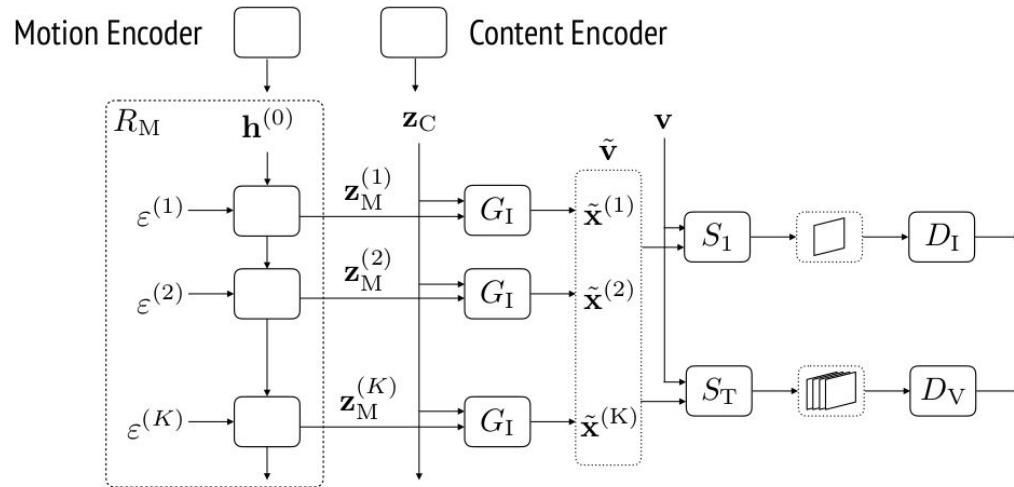
Video prediction:



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

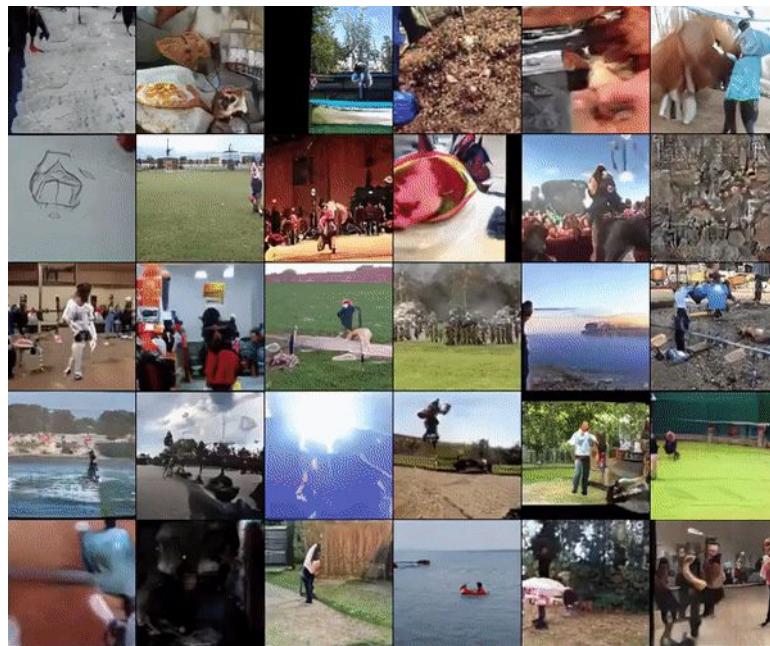
# Video Generation: MoCoGAN

Video prediction:



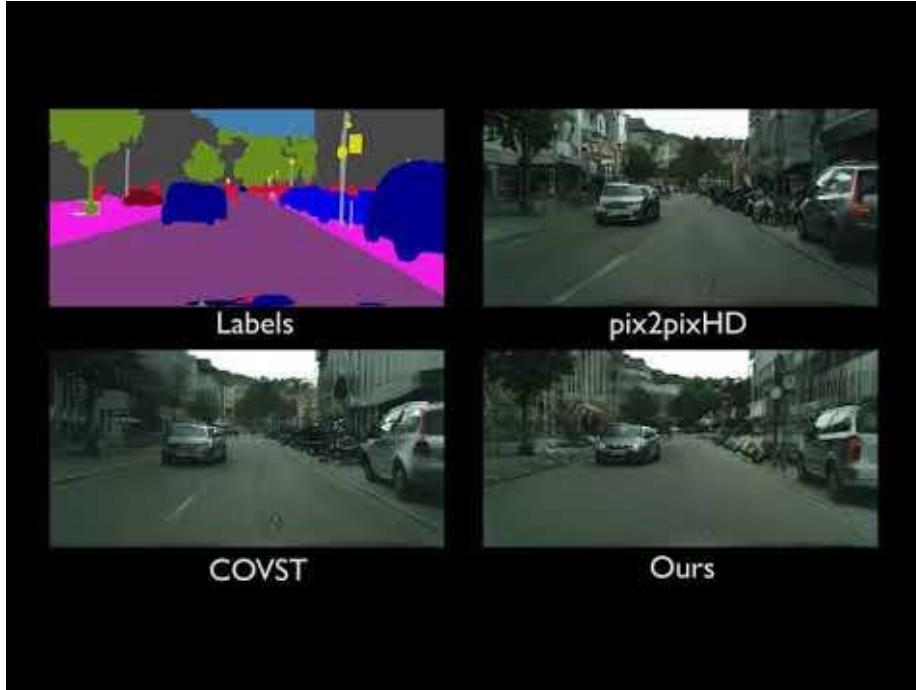
MoCoGAN: Decomposing Motion and Content for Video Generation Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, Jan Kautz CVPR'2018.

# Video Generation: DVGAN



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

# Video-to-video translation



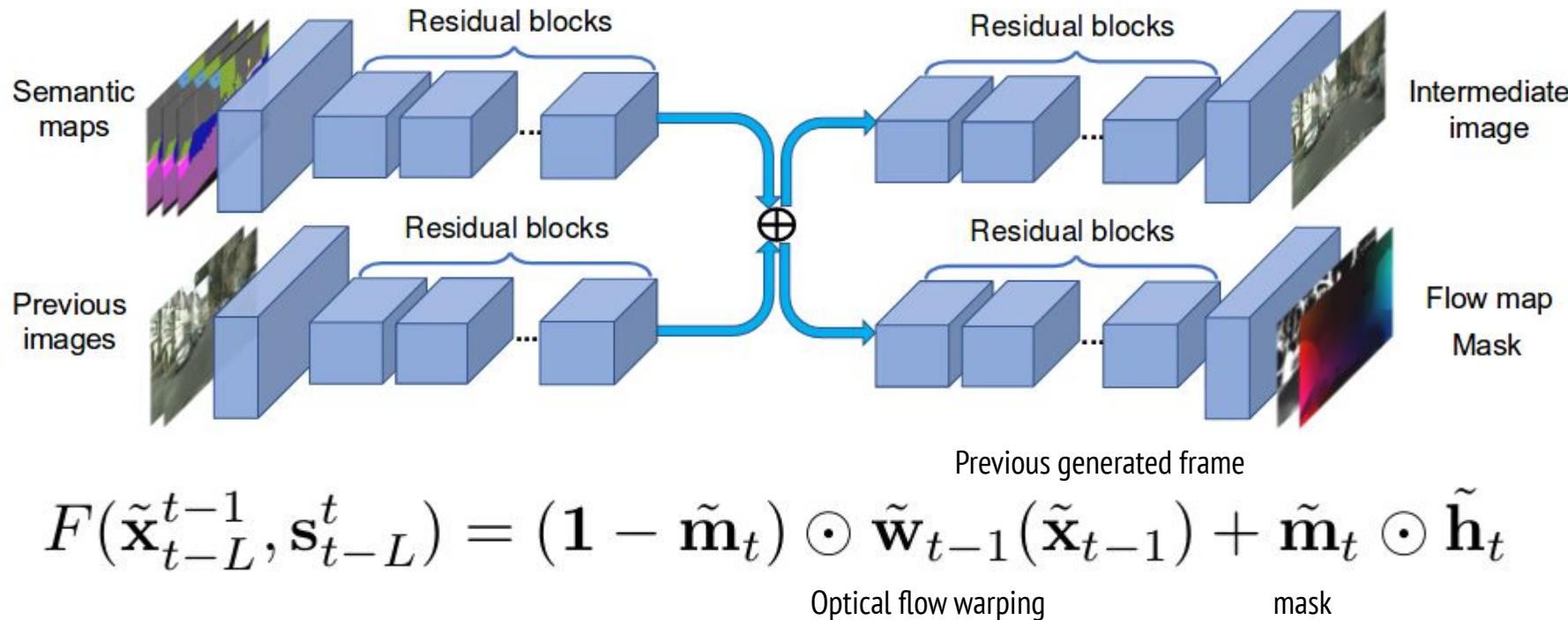
Wang, Ting-Chun, et al. "Video-to-video synthesis." NIPS'2018

# Video-to-video translation



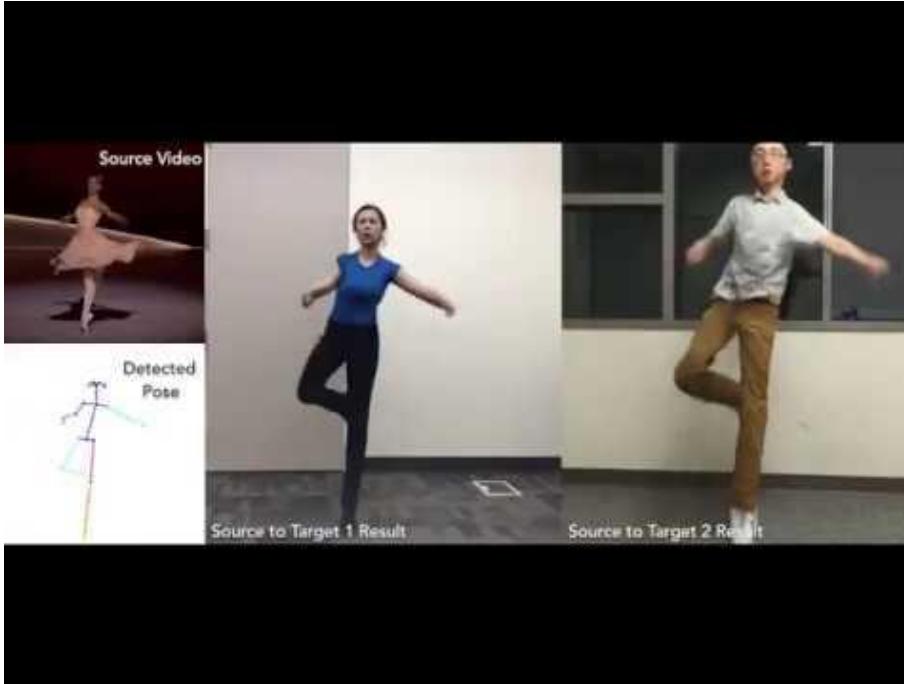
Wang, Ting-Chun, et al. "Video-to-video synthesis." NIPS'2018

# Video-to-video translation



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

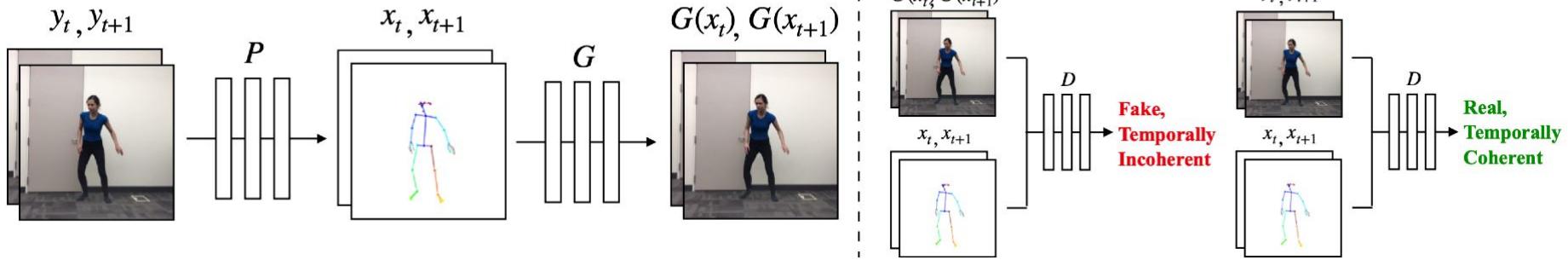
# Animating Single Subject



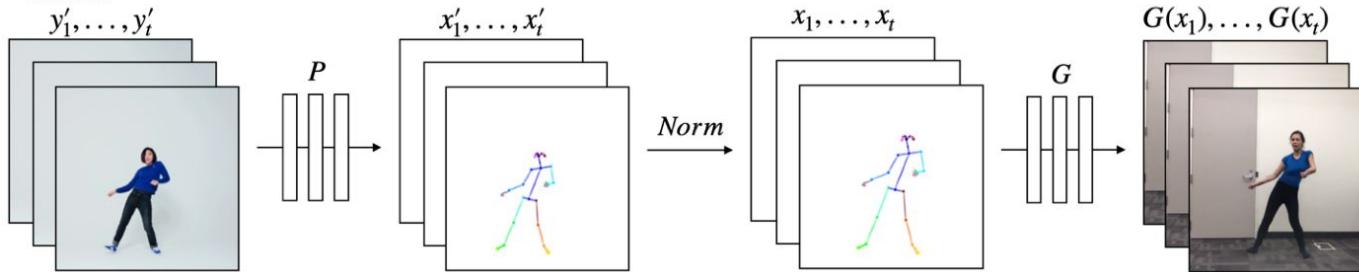
Chan, Caroline, et al. "Everybody dance now." ICCV'2019.

# Animating Single Subject

Training

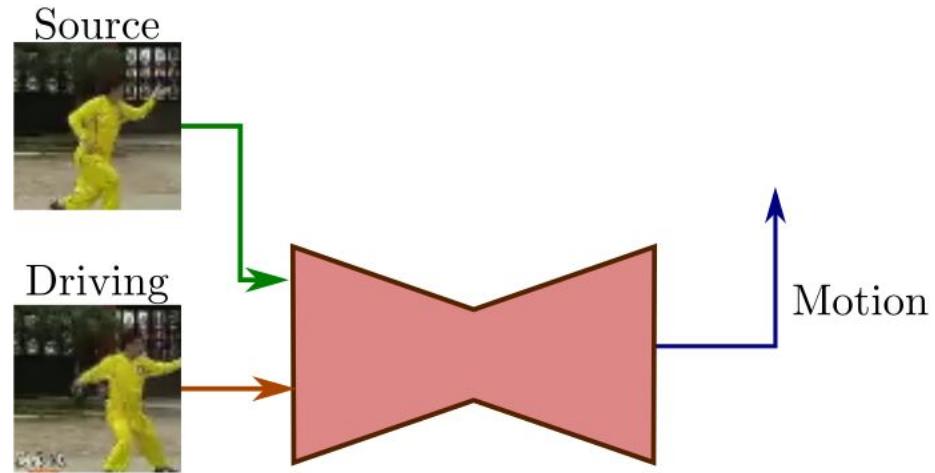


Transfer

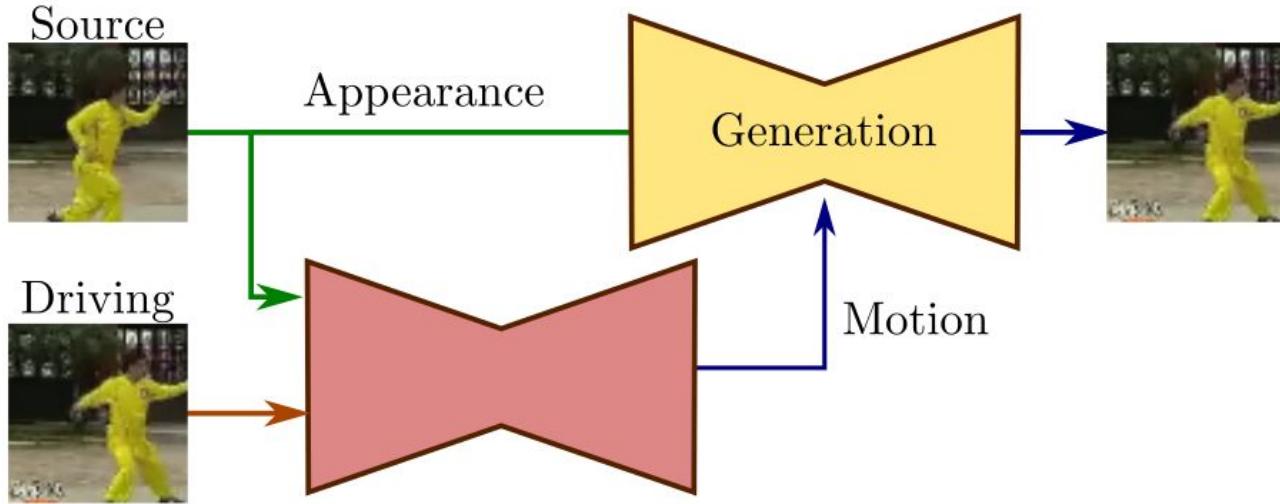


Chan, Caroline, et al. "Everybody dance now." ICCV'2019.

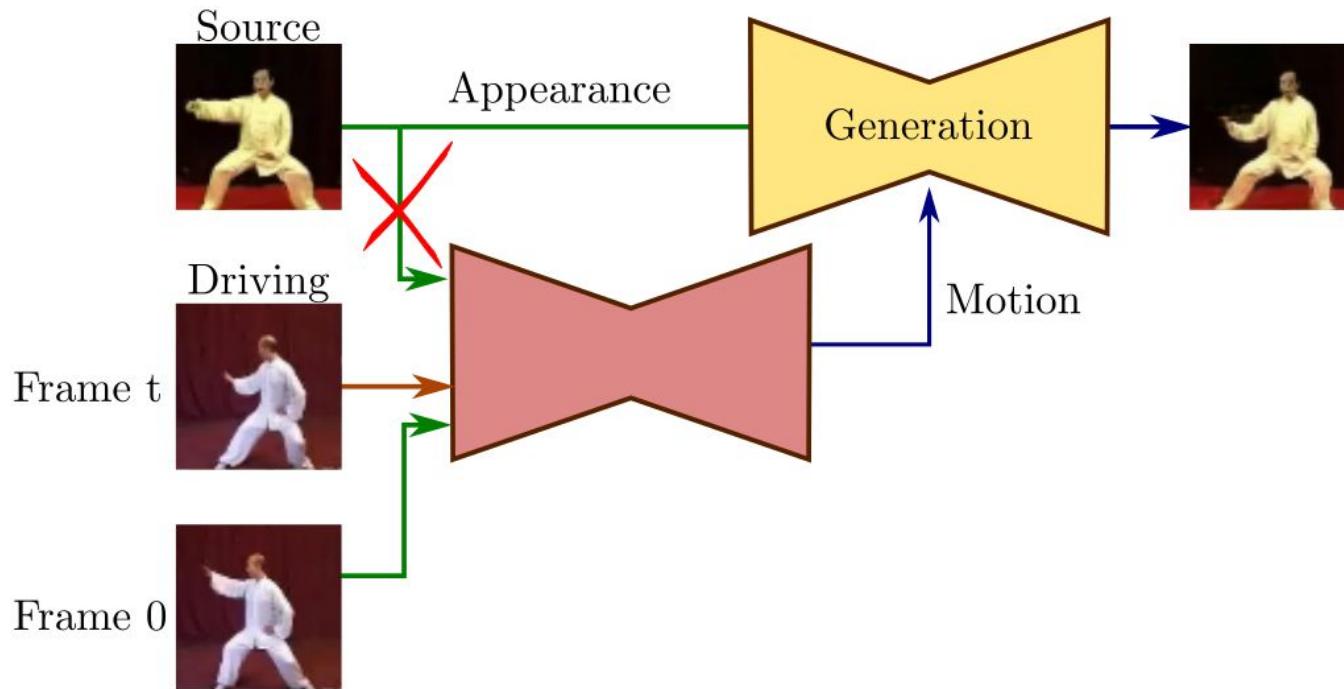
# Retargeting



# Retargeting



# Retargeting



# Retargeting



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