

# Lecture 7: Feed-forward generative models

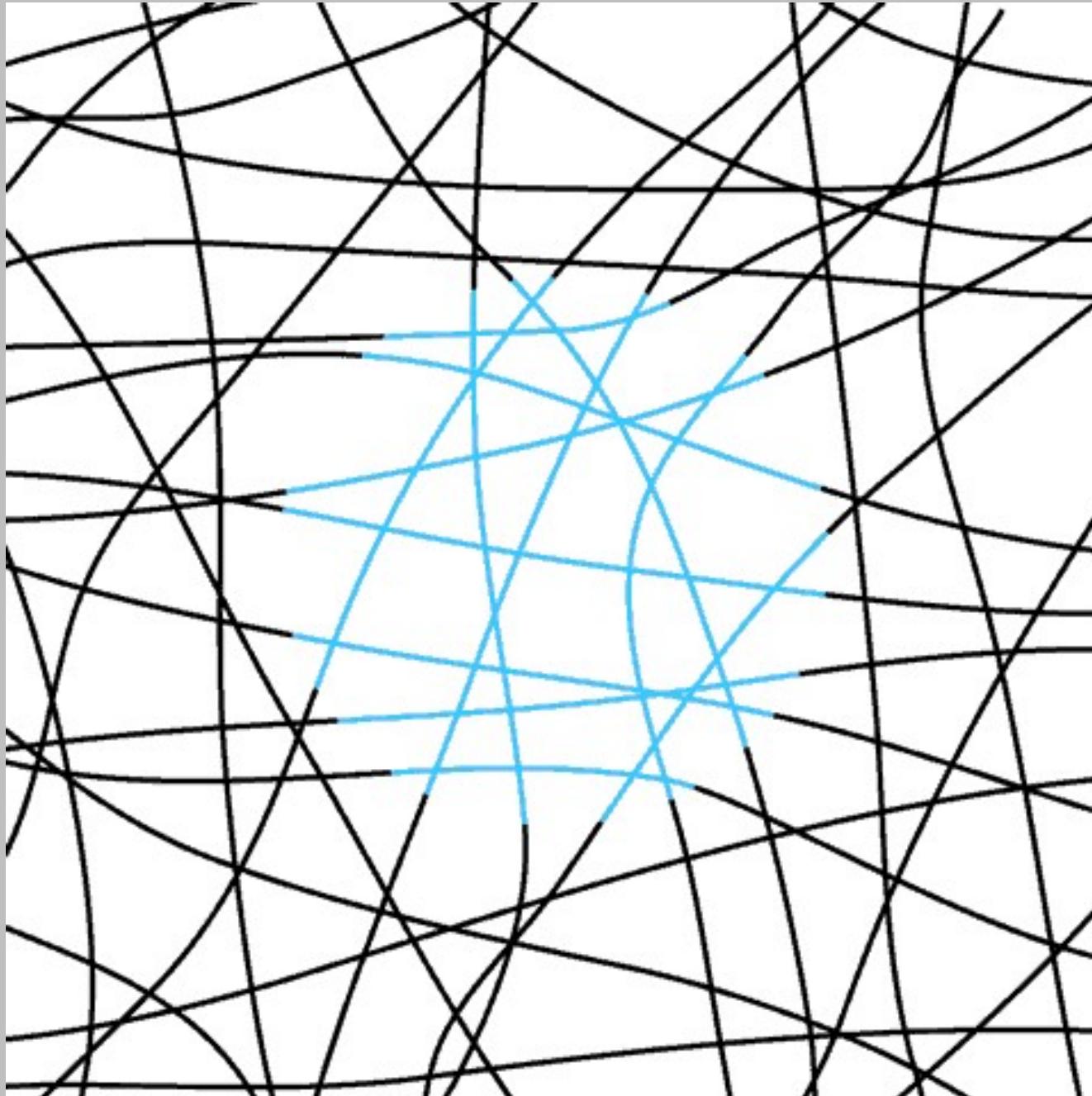
# Recap

- Previous lectures: image classification and related tasks (*networks that take images and output high-level semantic information*)
- Next two lectures: networks that *generate* images
- Reminder: we have seen network-generated images before (*pre-image method*). More of that today.

# Why do we care about image synthesis

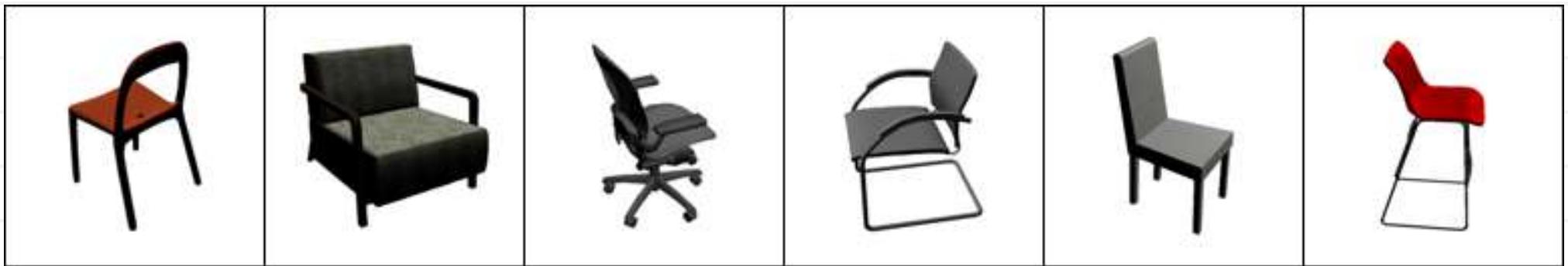
1. Fun
2. A harder problem, a frontier in AI
3. Computer graphics, design, image editing
4. To improve computer vision (“analysis-by-synthesis”)

# Analysis by synthesis in your brain



# Feed-forward conditional generation

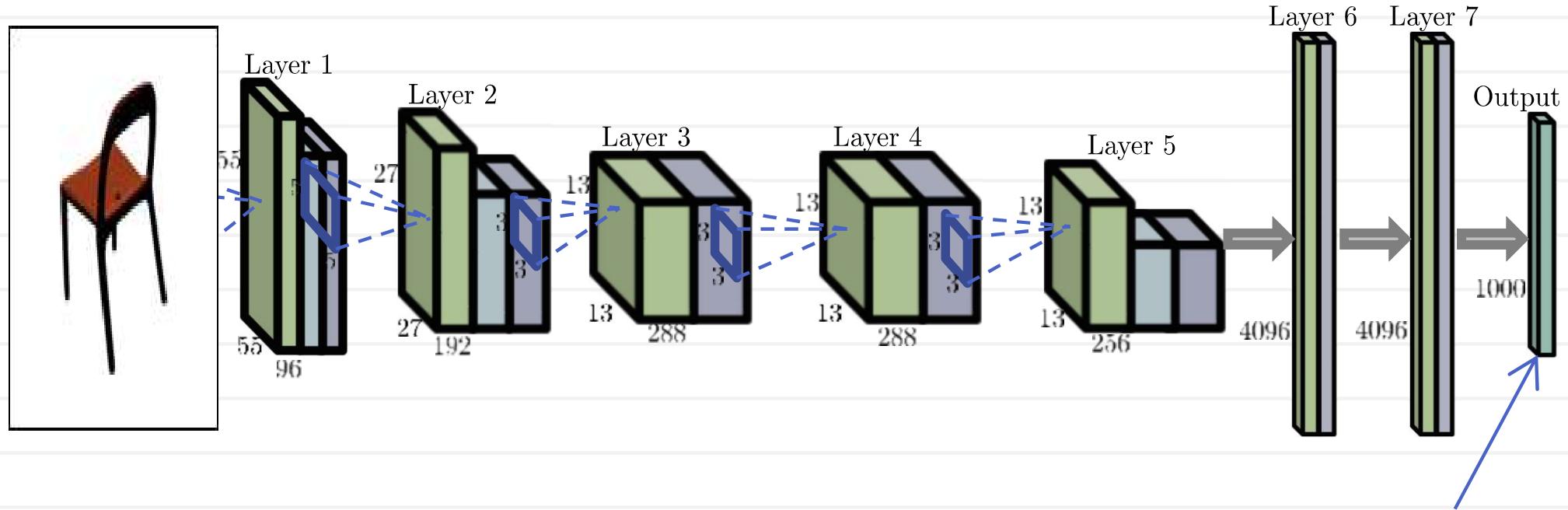
Source images:



Each image has “chair ID”, and viewpoint direction

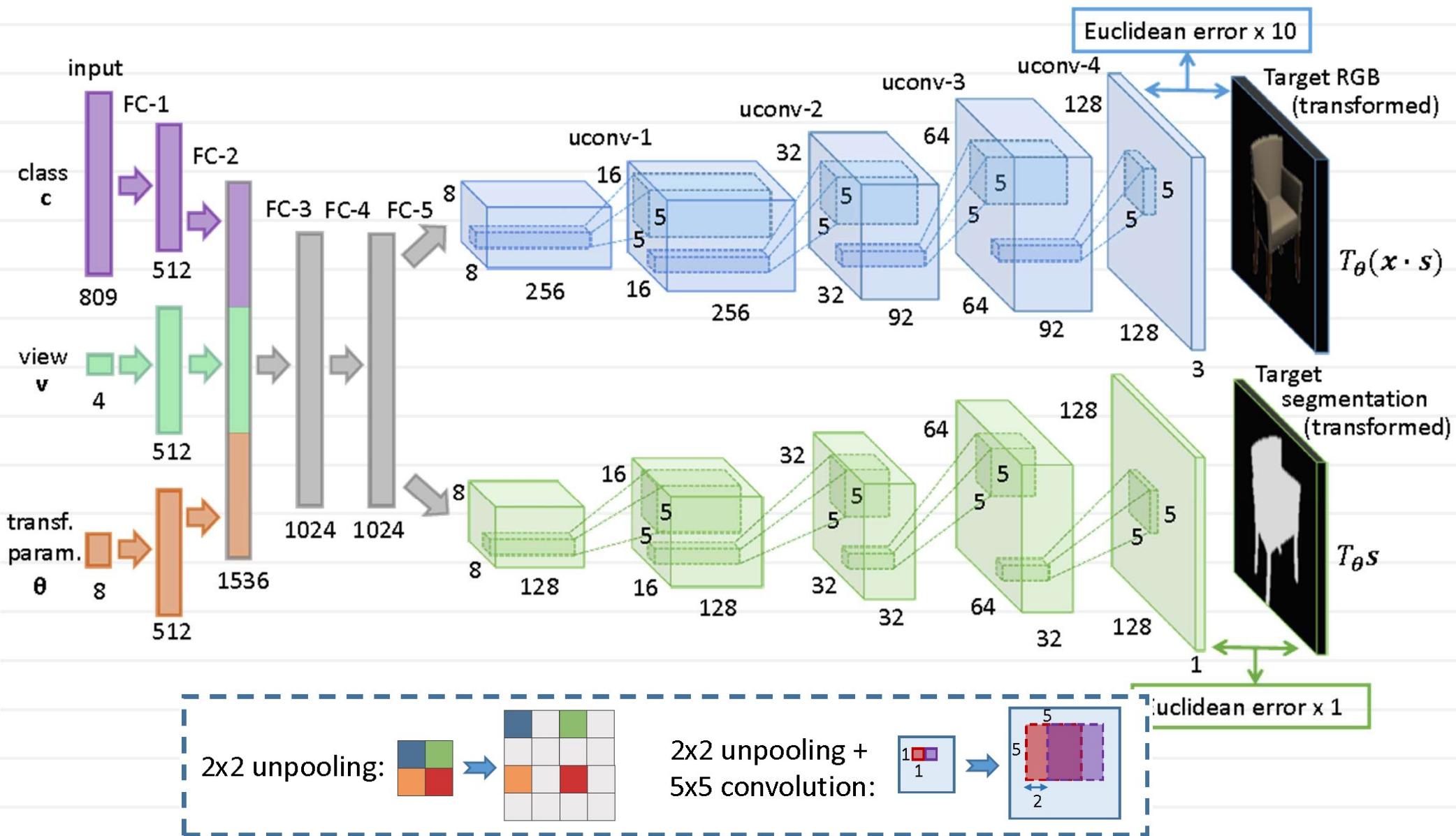
[Dosovitskiy et al. CVPR 2015]

# “Natural” dataflow



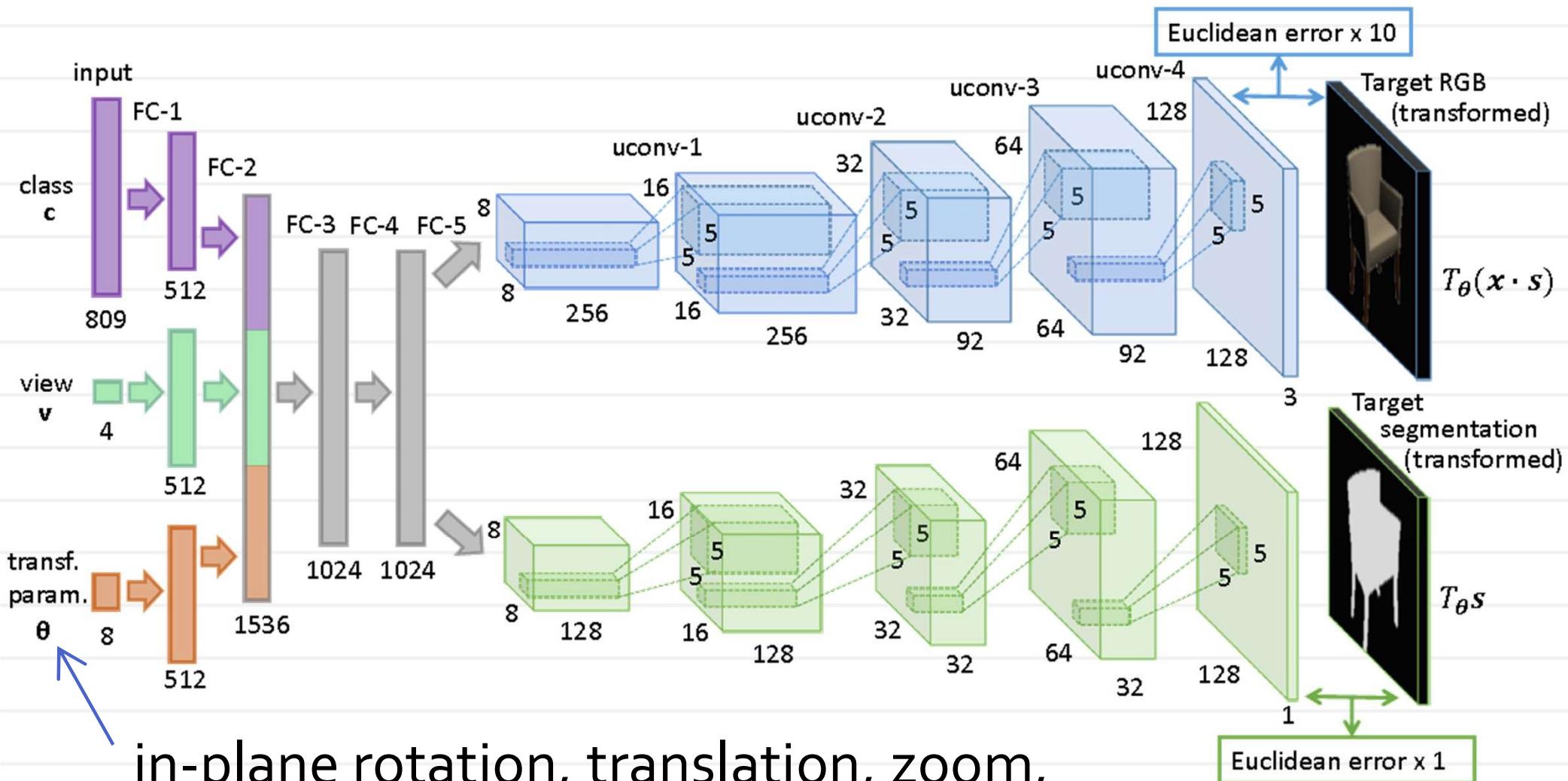
chair type  
view angle  
etc.

# Feed-forward conditional generation



[Dosovitskiy et al. CVPR 2015]

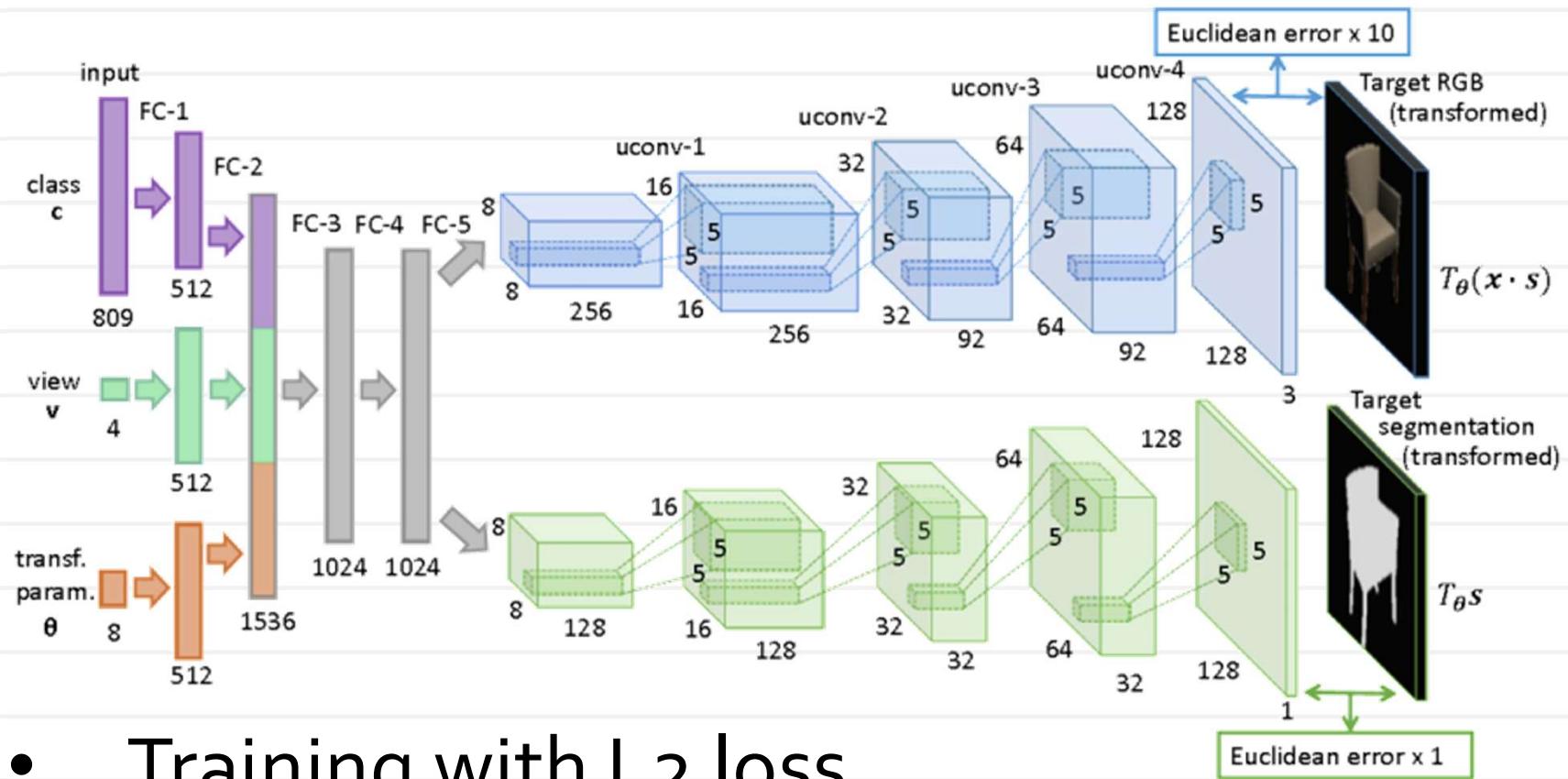
# Feed-forward conditional generation



in-plane rotation, translation, zoom,  
stretching horizontally or vertically,  
changing hue, changing saturation,  
changing brightness

[Dosovitskiy et al. CVPR 2015]

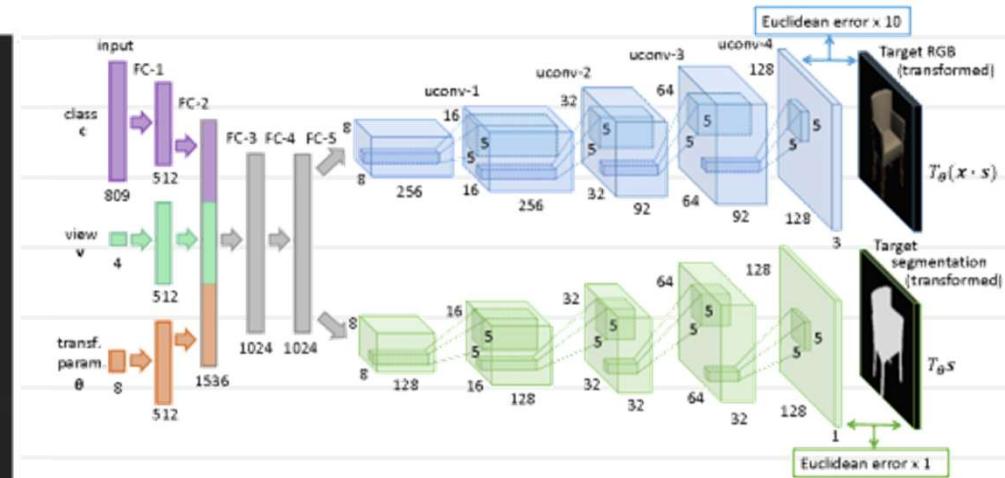
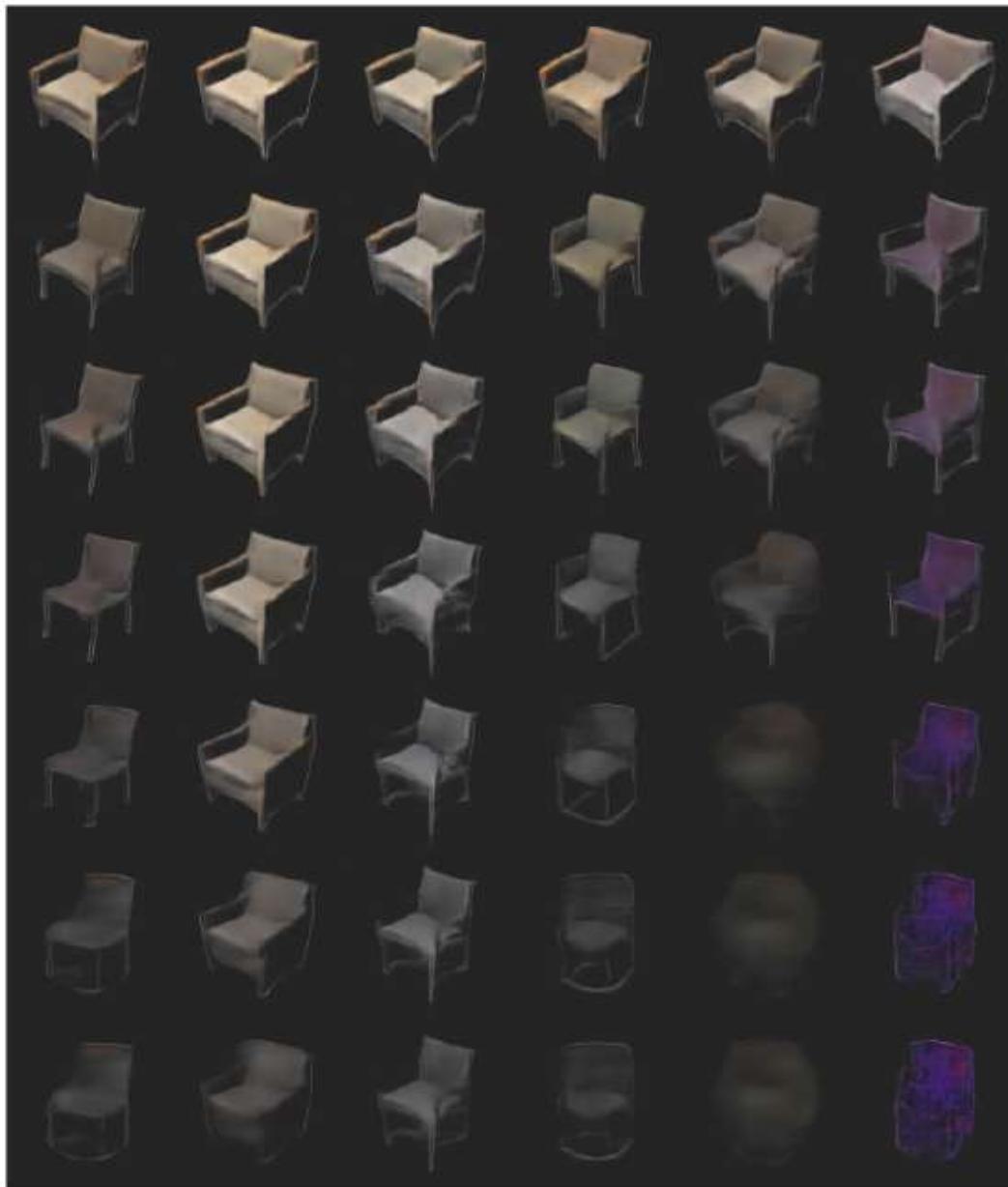
# Feed-forward conditional generation



- Training with L<sub>2</sub> loss
- Data augmentation
- Mask is used to regularize (multi-task learning) and to clip the results

[Dosovitskiy et al. CVPR 2015]

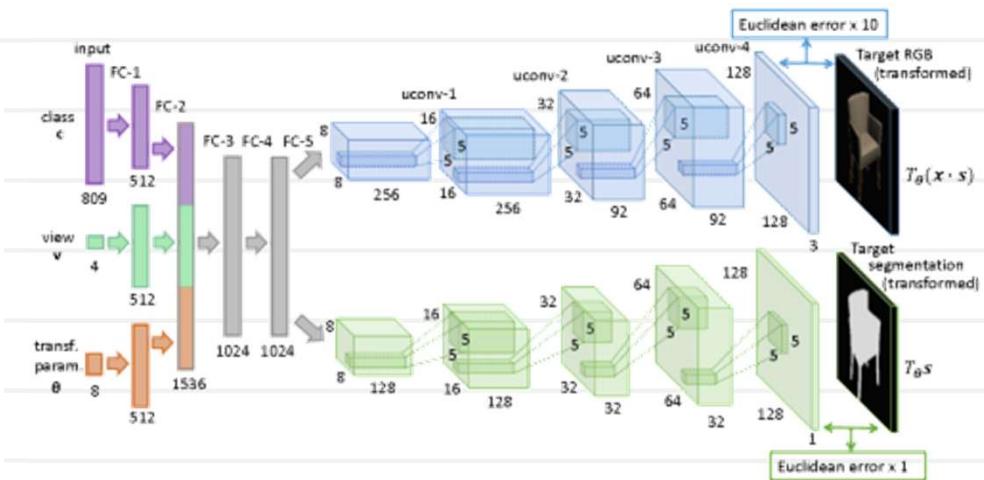
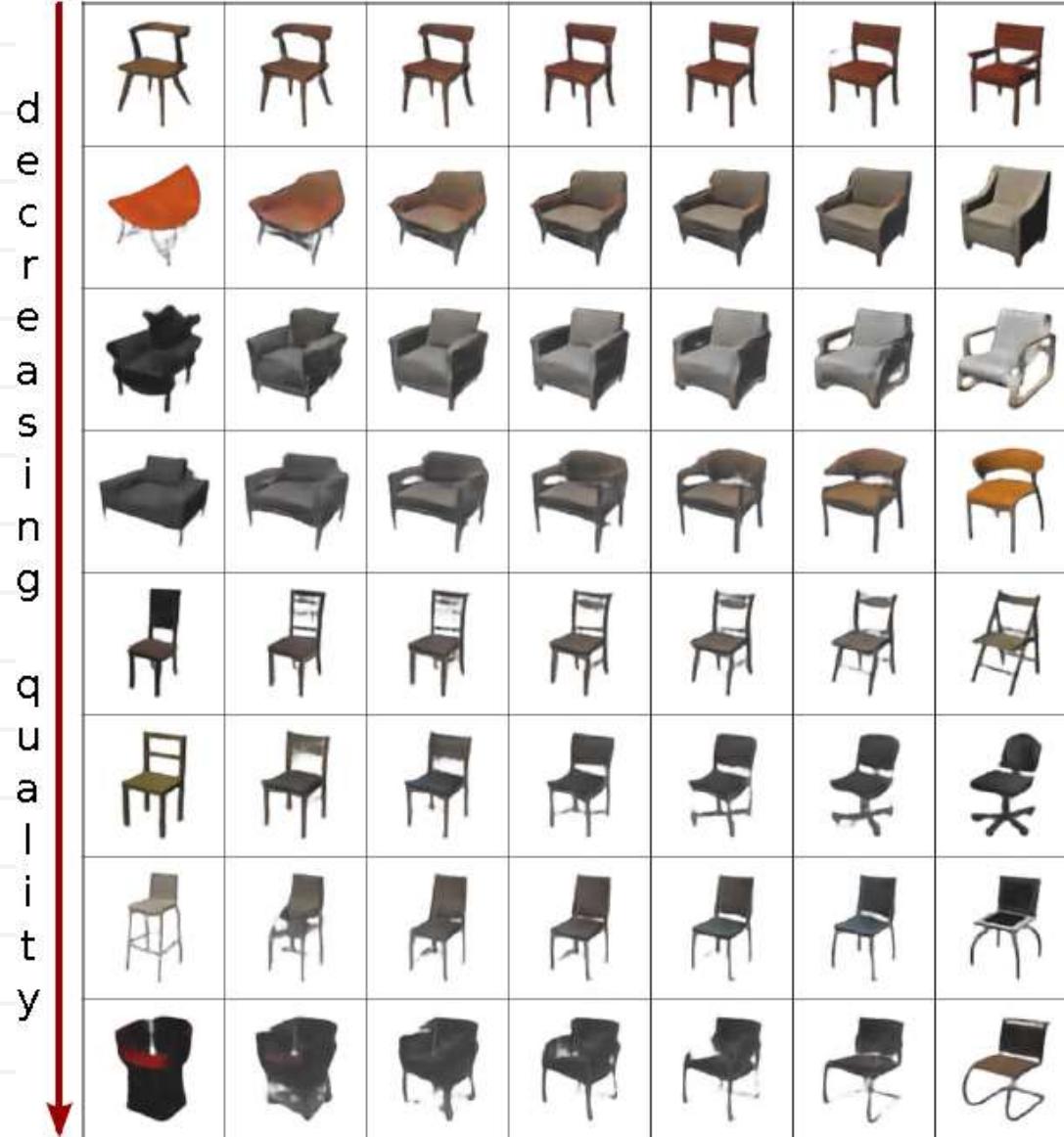
# Feed-forward conditional generation



Increasing a  
single neuron in  
FC1

[Dosovitskiy et al. CVPR 2015]

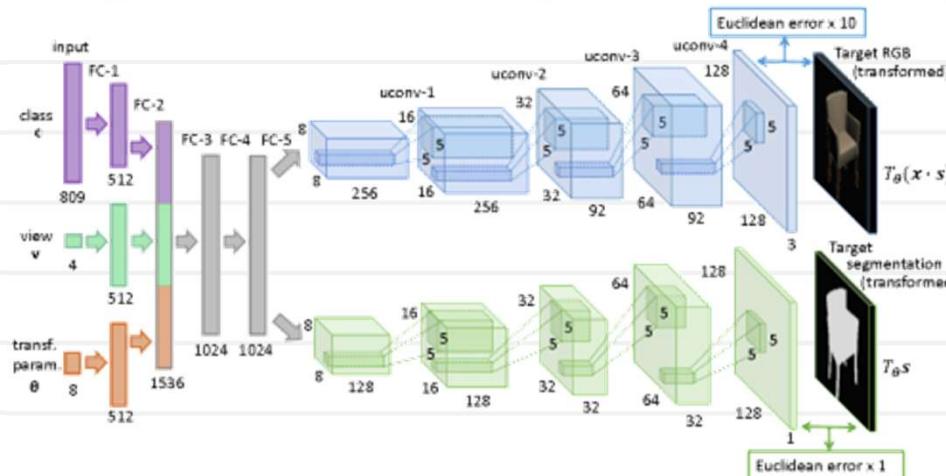
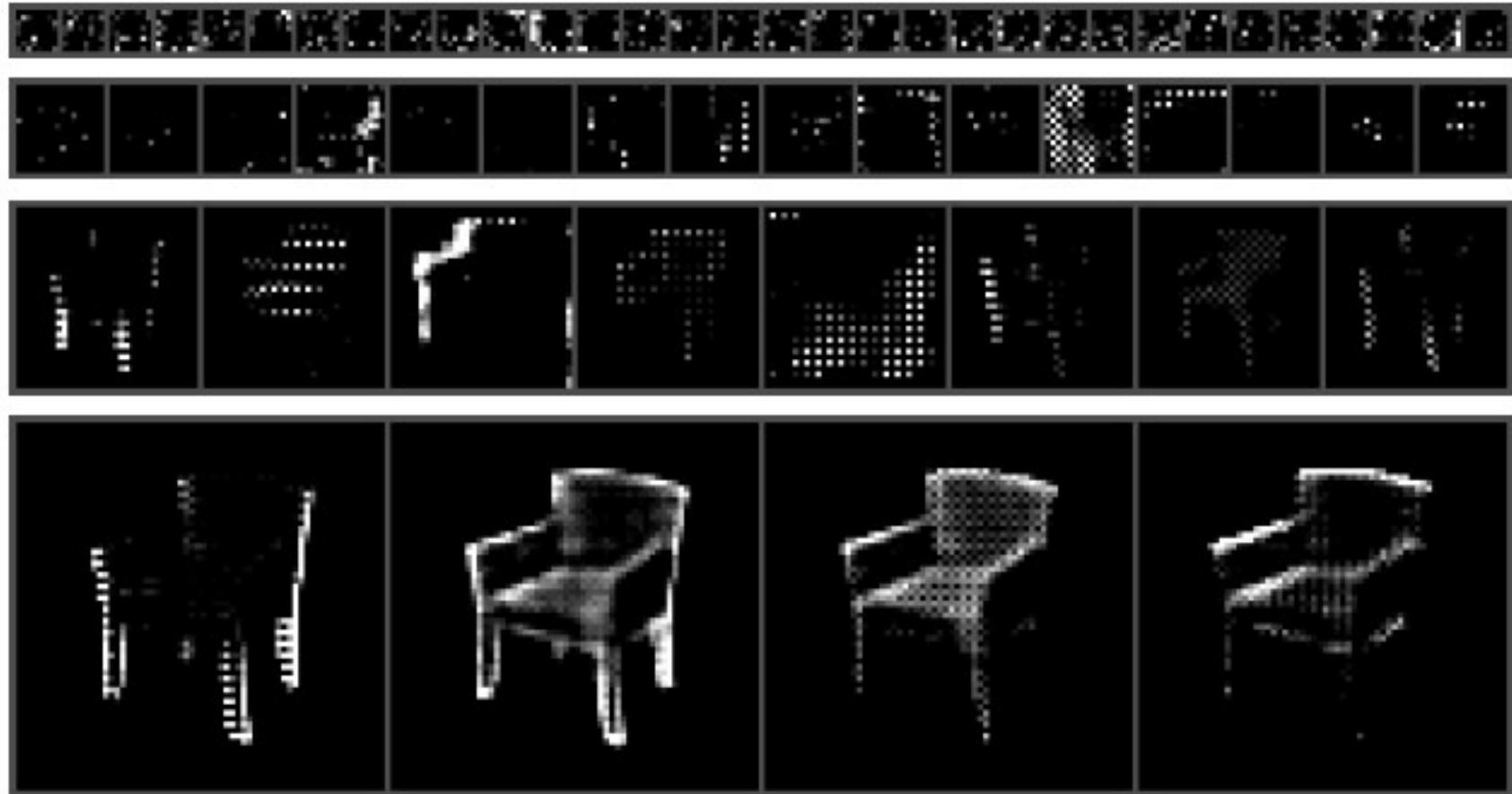
# Feed-forward conditional generation



Morphing  
between chairs

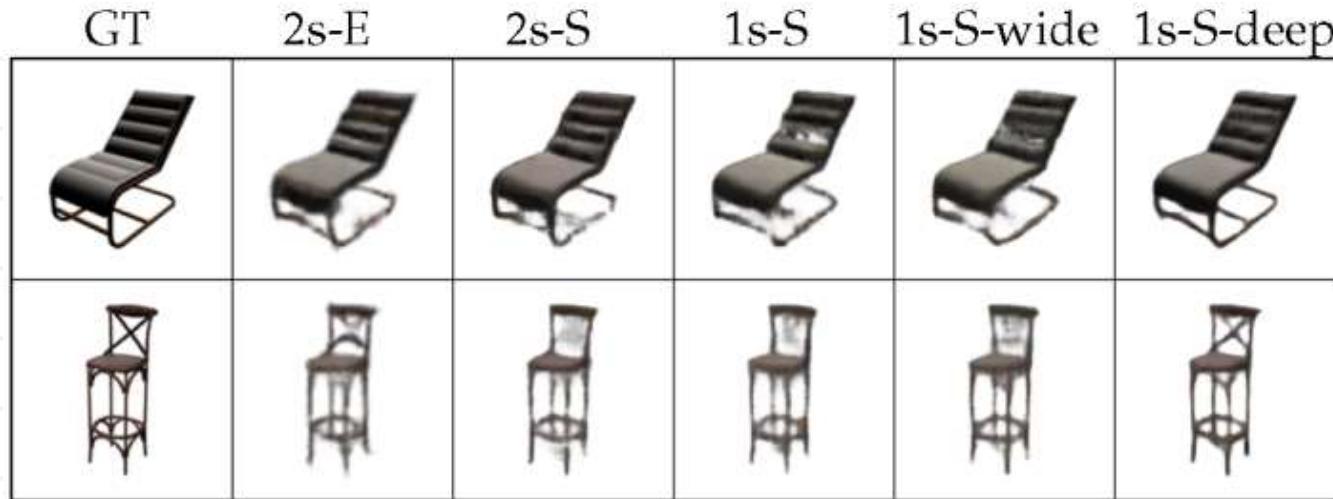
[Dosovitskiy et al. CVPR 2015]

# Feed-forward conditional generation

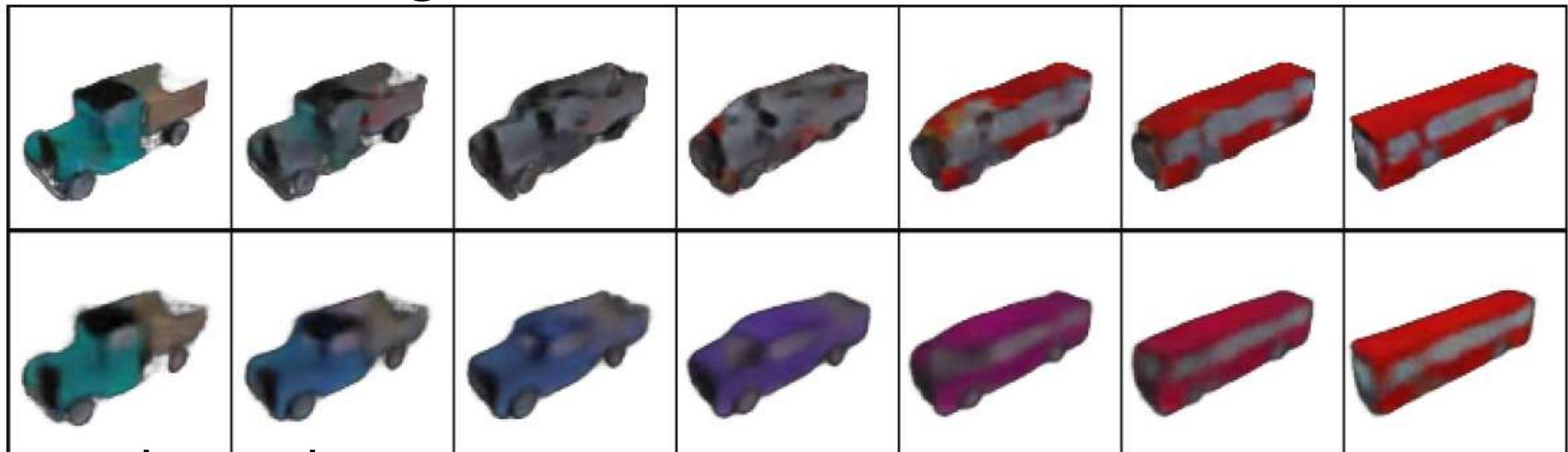


[Dosovitskiy et al. CVPR 2015]

# More results



with data augmentation



without data augmentation

Blurriness appears as a result of the L<sub>2</sub>-loss

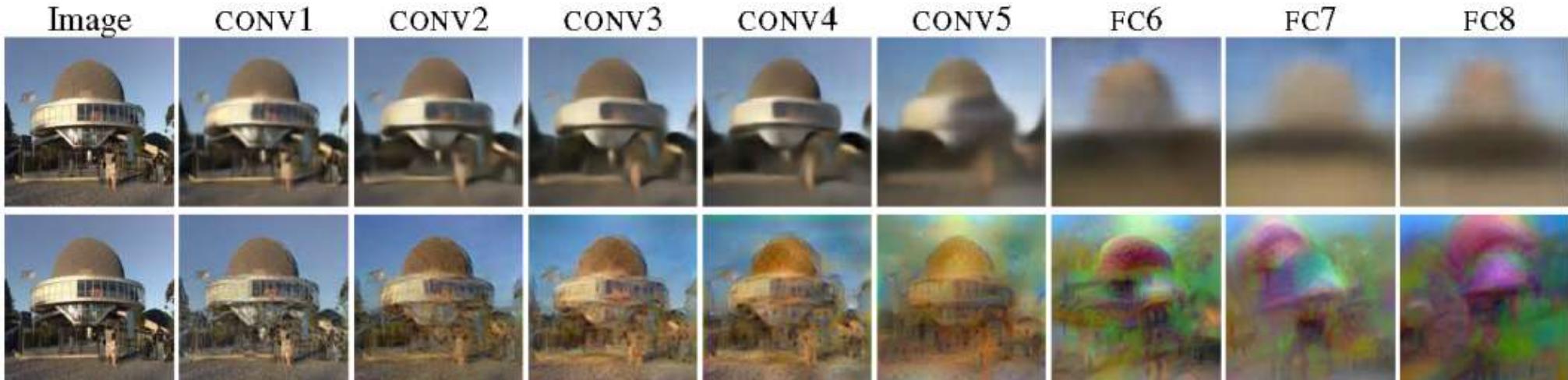
[Dosovitskiy et al. CVPR 2015]

# Inverting CNNs with ConvNets

[Dosovitskiy and Brox, Arxiv15]:

- Train a special generating network to reconstruct the original image from AlexNet activations
- Their architecture for FC: 3 FC + 5 up-conv layers
- Their architecture for Conv layers: several conv + several up-convolutional layers

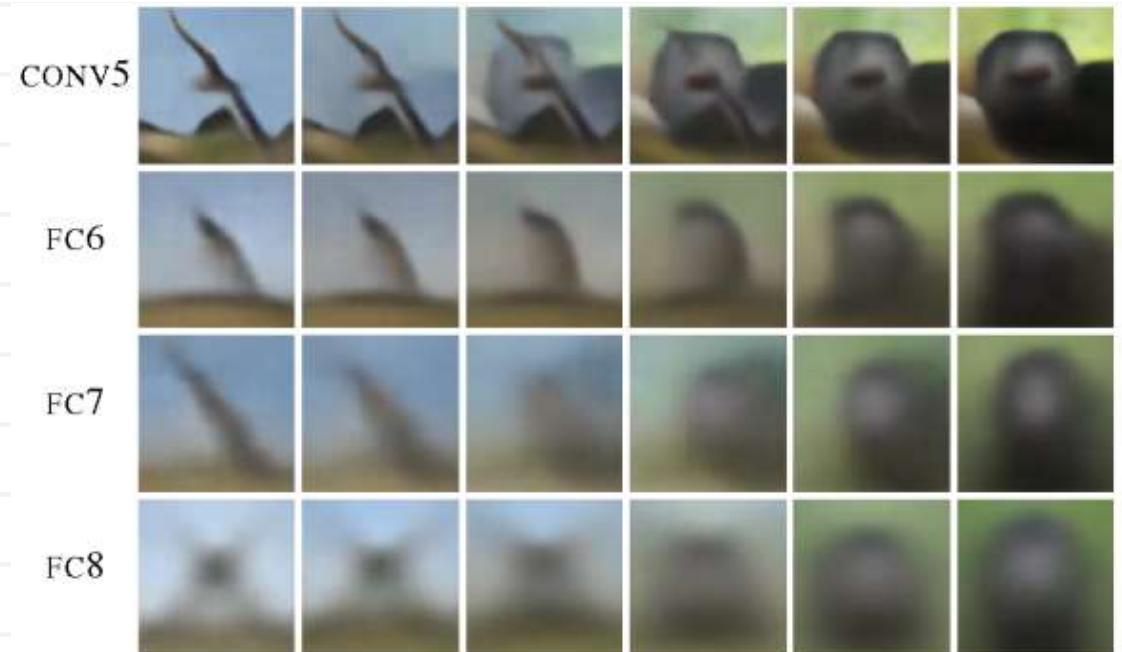
Their result:



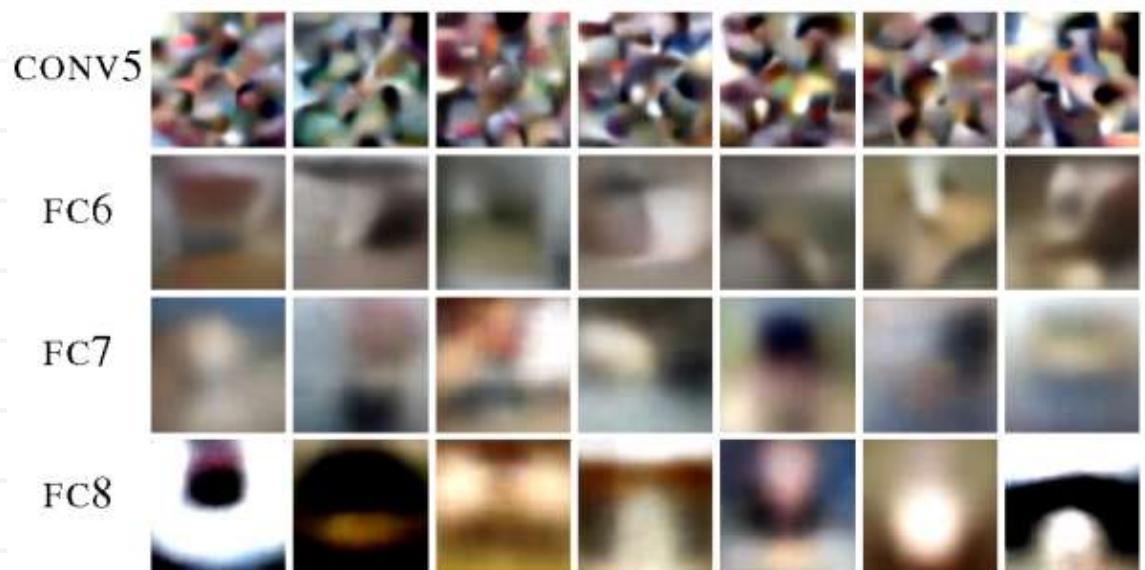
Pre-image method of Mahendran and Vedaldi (Lecture 5)

# Experiments with the inverting architecture

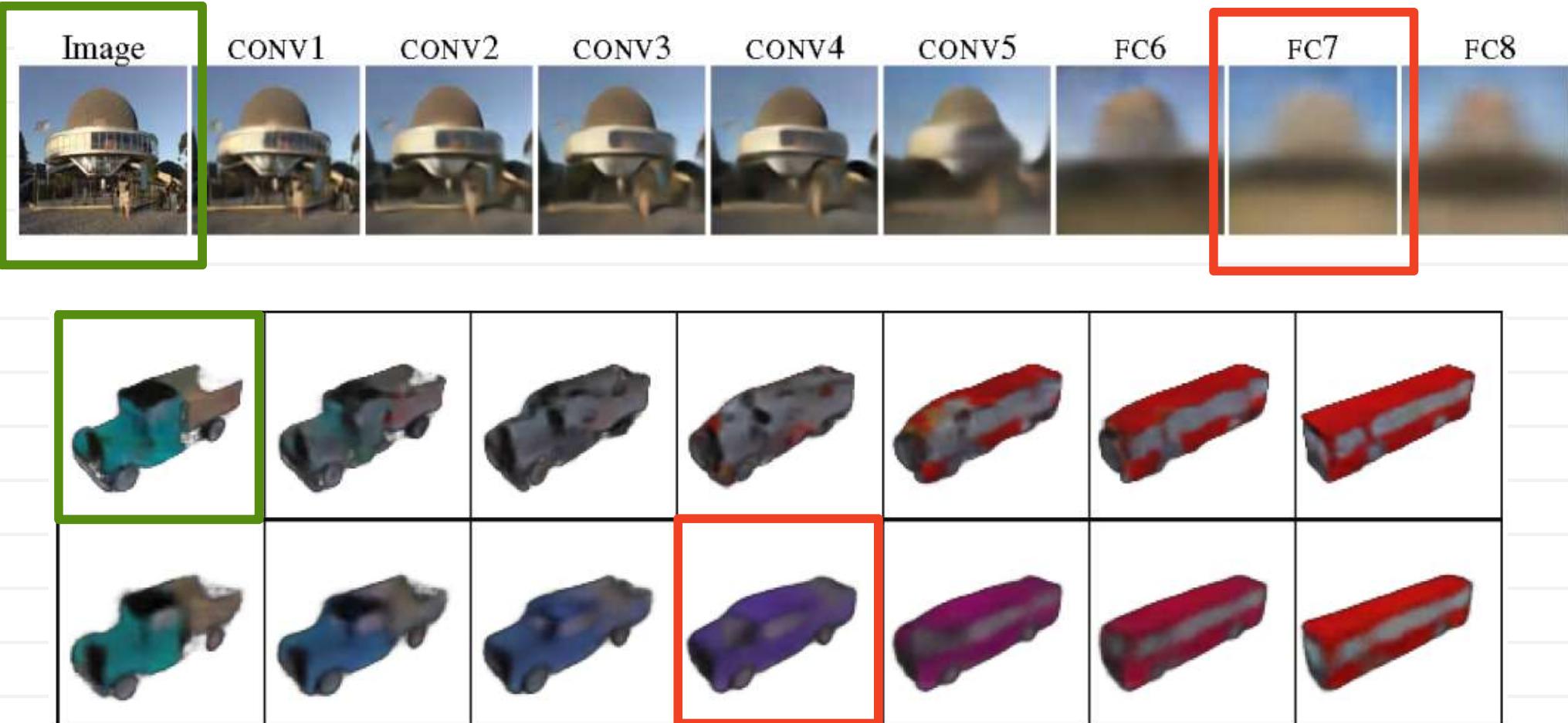
Interpolation  
between  
two images:



Sampling  
from random  
activations:



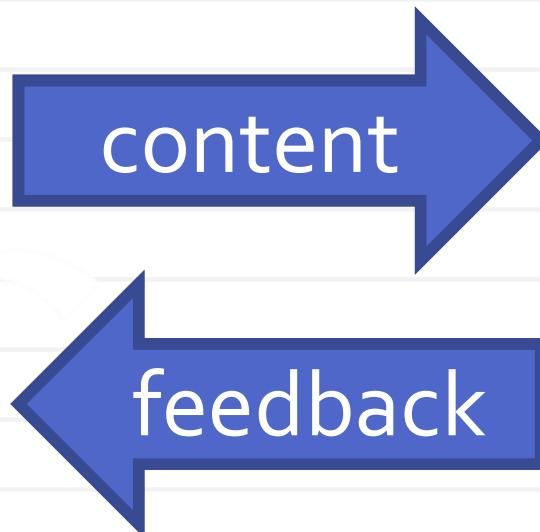
# Let us look at reconstruction deficiencies



- How can we tell between **real** and **synthetic**?
- Can a **computer** tell between them?

# Let us look at reconstruction deficiencies

- How can we tell between **real** and **synthetic**?
- Can a **computer** tell between them?



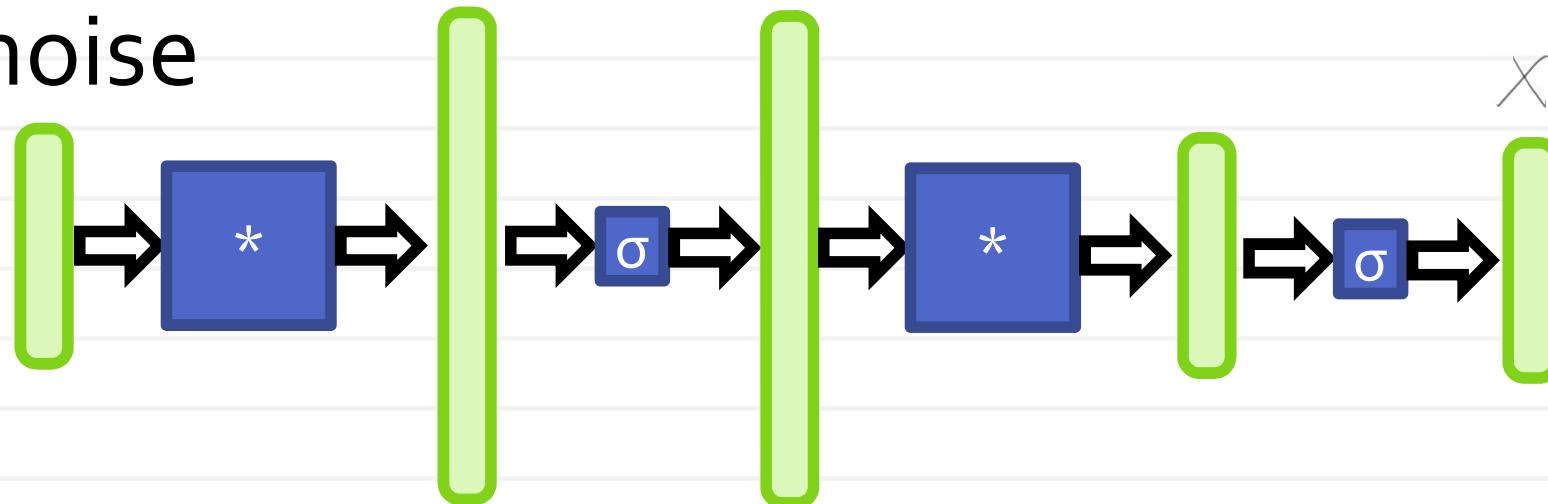
**Generator Network:**  
generates content  
(e.g. images)

**Discriminator**  
**Network:**  
assesses the result  
(e.g. images)

# Adversarial generative networks

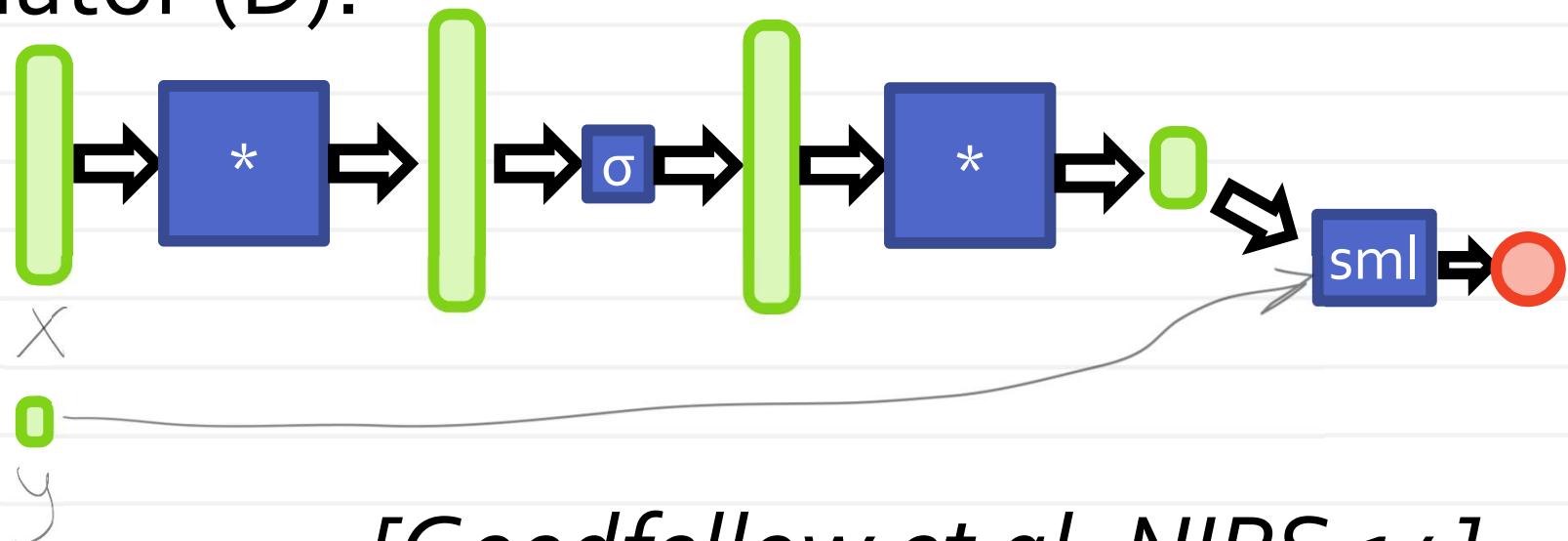
Generator (G):

noise



Discriminator (D):

A 4x4 grid of handwritten digits from 0 to 9, arranged in a 4x4 pattern. Below the grid, there are handwritten labels: "X" on the left, "Y" below the grid, and "sml" to the right of the discriminator's output.



[Goodfellow et al. NIPS 14]

# Adversarial learning

for number of training iterations do

  for  $k$  steps do

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

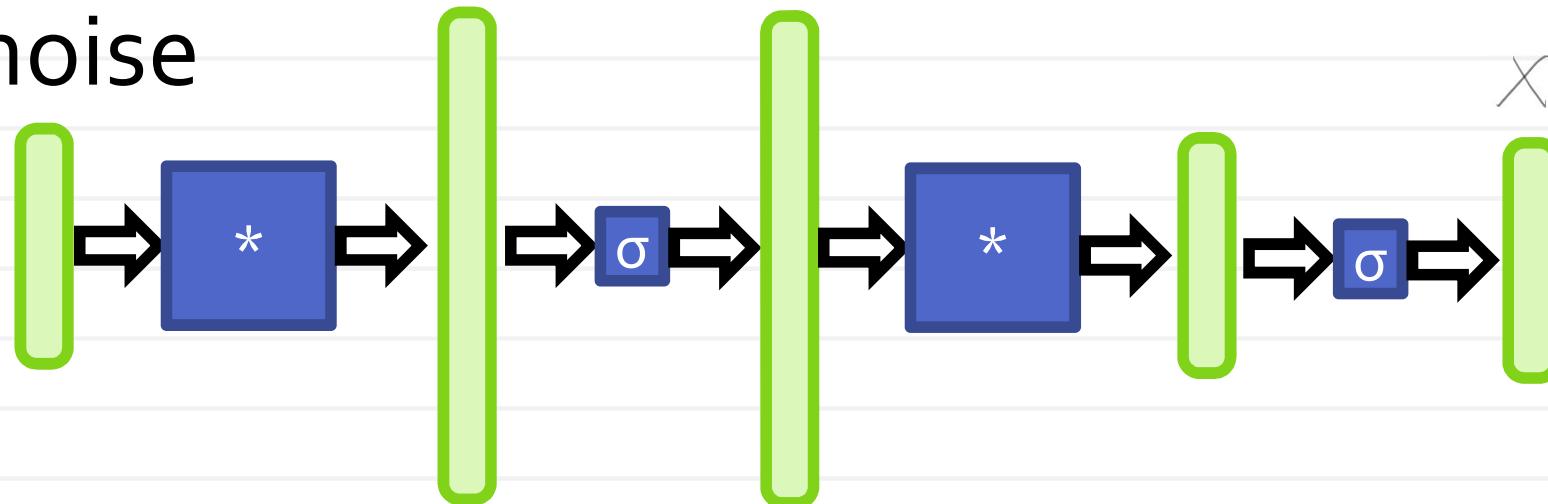
end for

[Goodfellow et al. NIPS 14]

# Adversarial generative networks

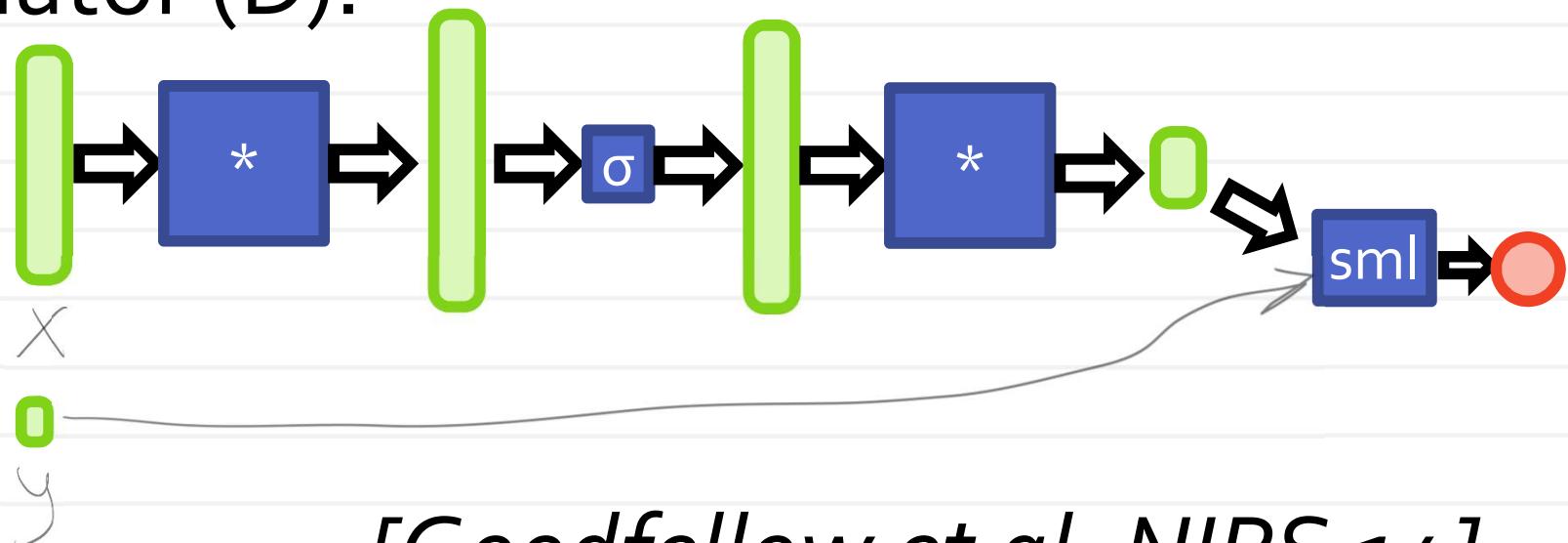
Generator (G):

noise



Discriminator (D):

5 0 4 1 9 2 1 3  
4 4 6 0 4 5 6 7  
2 0 2 7 1 8 6 4  
2 3 5 9 1 7 6 2  
8 6 3 7 5 8 0 9  
8 7 6 0 9 7 5 7  
2 3 9 4 9 2 1 6  
5 6 7 9 9 3 7 0

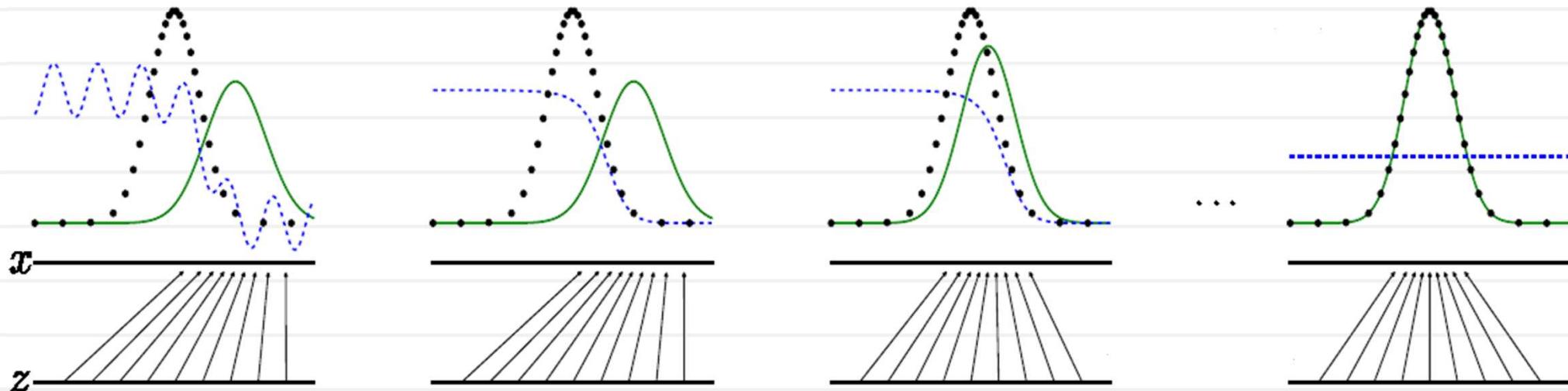


[Goodfellow et al. NIPS 14]

# Adversarial learning

Looking for a *saddle point*:

$$\min_G \max_D V(D, G) = E_{x \sim P_D(x)} [\log D(x)] + \\ + E_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$



[Goodfellow et al. NIPS 14]

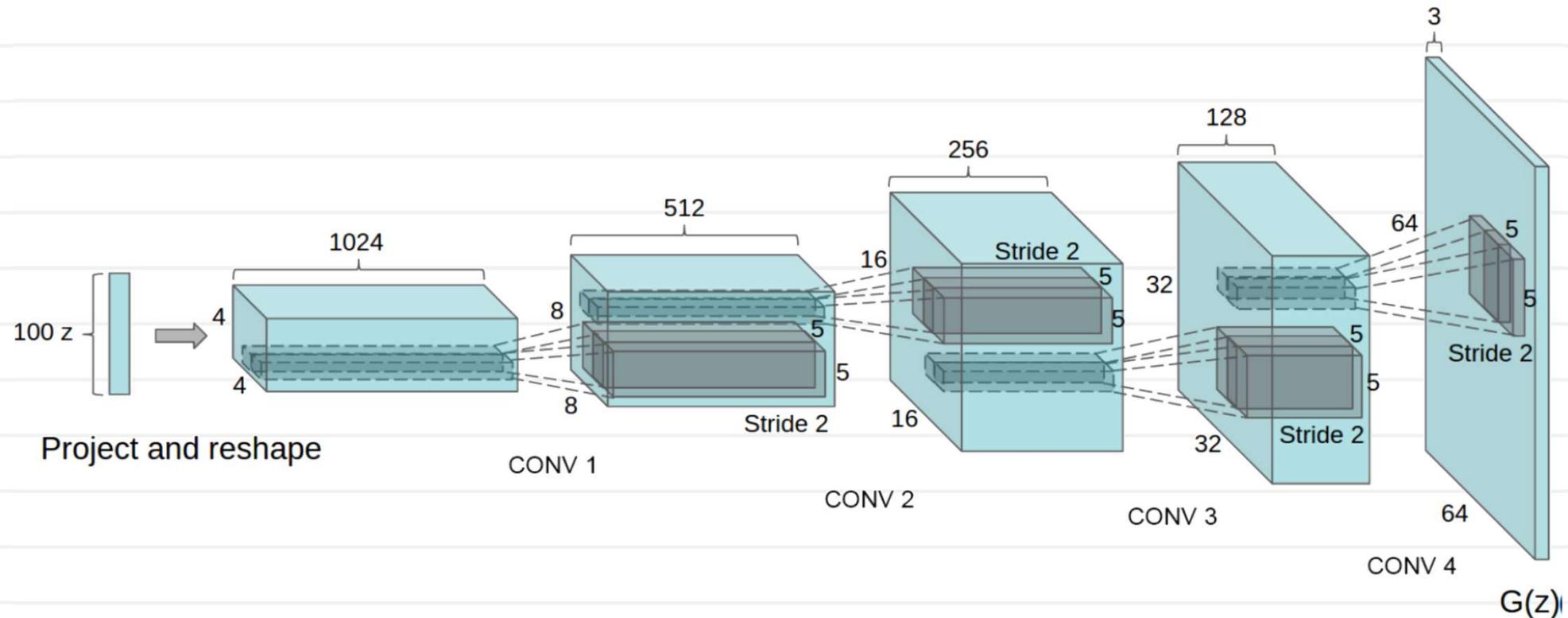
# Modified GANs

$$\min_G \max_D V(D, G) = E_{x \sim P_D(x)} [\log D(x)] + \\ + E_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$
$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

- Not a zero-sum game anymore
- Works better in the initial stage, when discriminator is much better than generator

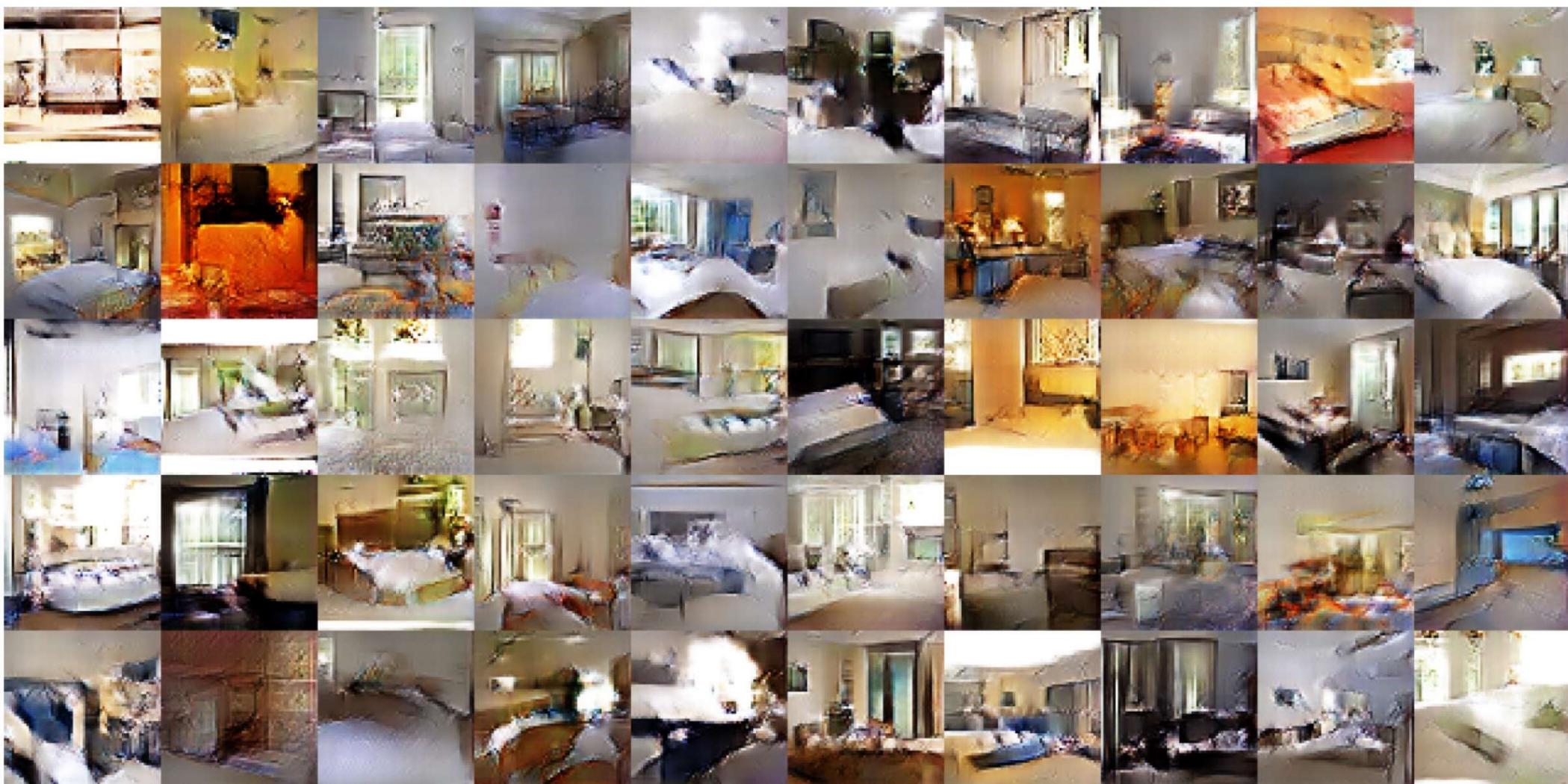
# DCGAN



- Discriminator uses Leaky ReLU + tanH (last layer)

[Radford et al. 2015]

# DCGAN



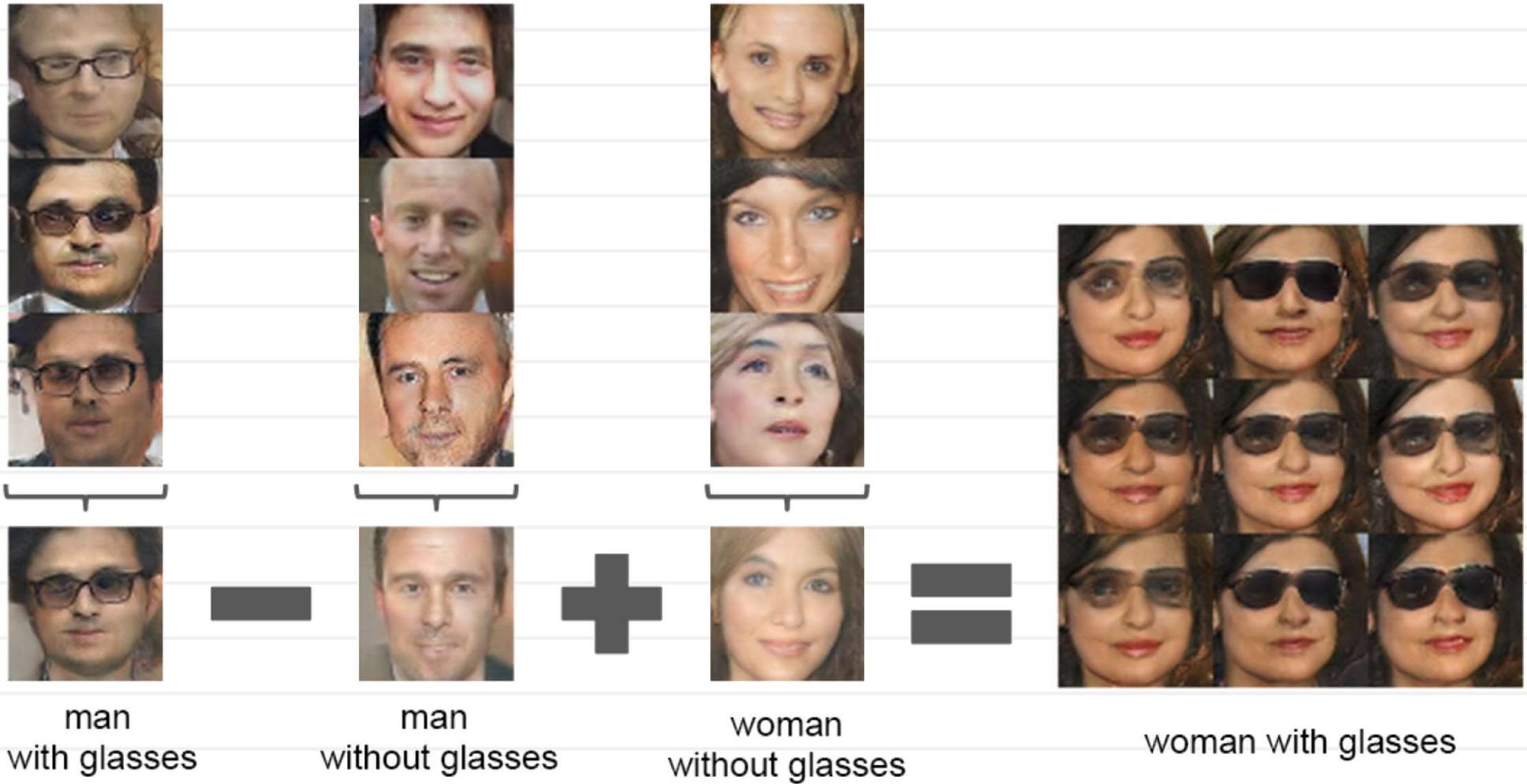
[Radford et al. 2015]

# DCGAN



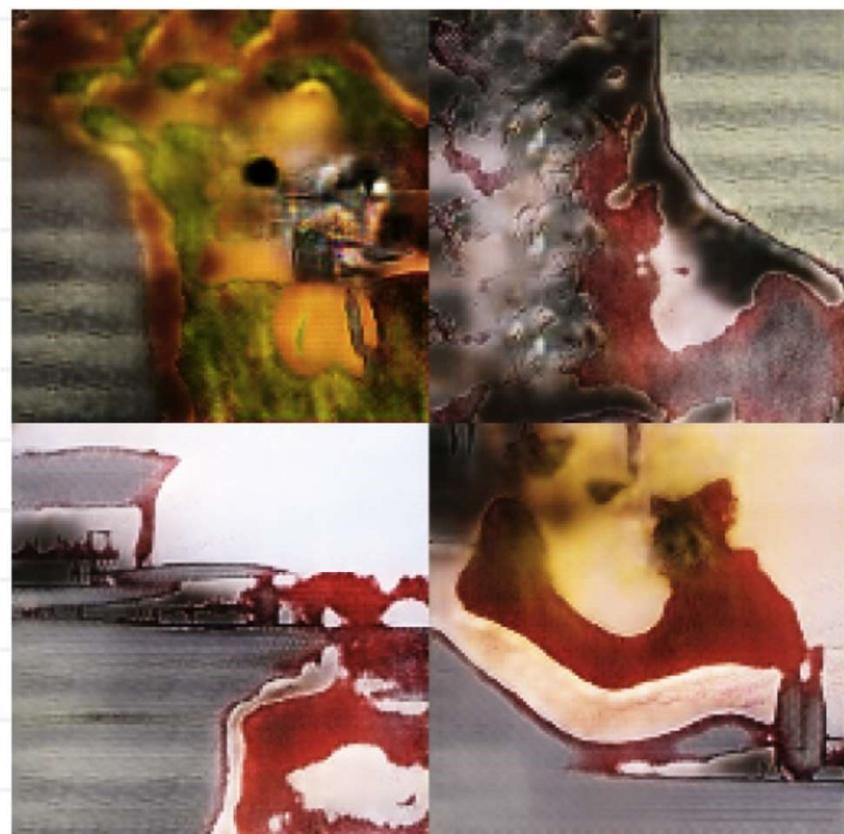
[Radford et al. 2015]

# DCGAN



[Radford et al. 2015]

# SOTA with GANs



DCGAN



[Salimans 16]

# Optimal discriminator

$$\min_G \max_D V(D, G) = E_{x \sim p_D(x)} [\log D(x)] +$$

$$+ E_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

$p_{model}(x)$

vs

$p_{data}(x)$

Optimizing discriminator gives:

$$D(x) = \frac{p_{data}(x)}{p_{data}(x) + p_{model}(x)}$$

Optimal discriminator works as *ratio estimator*!

# Optimal generator

$$\min_G \max_D V(D, G) = E_{x \sim P_D(x)} [\log D(x)] + \\ + E_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

What is the optimal generator for the fixed discriminator?

Answer:

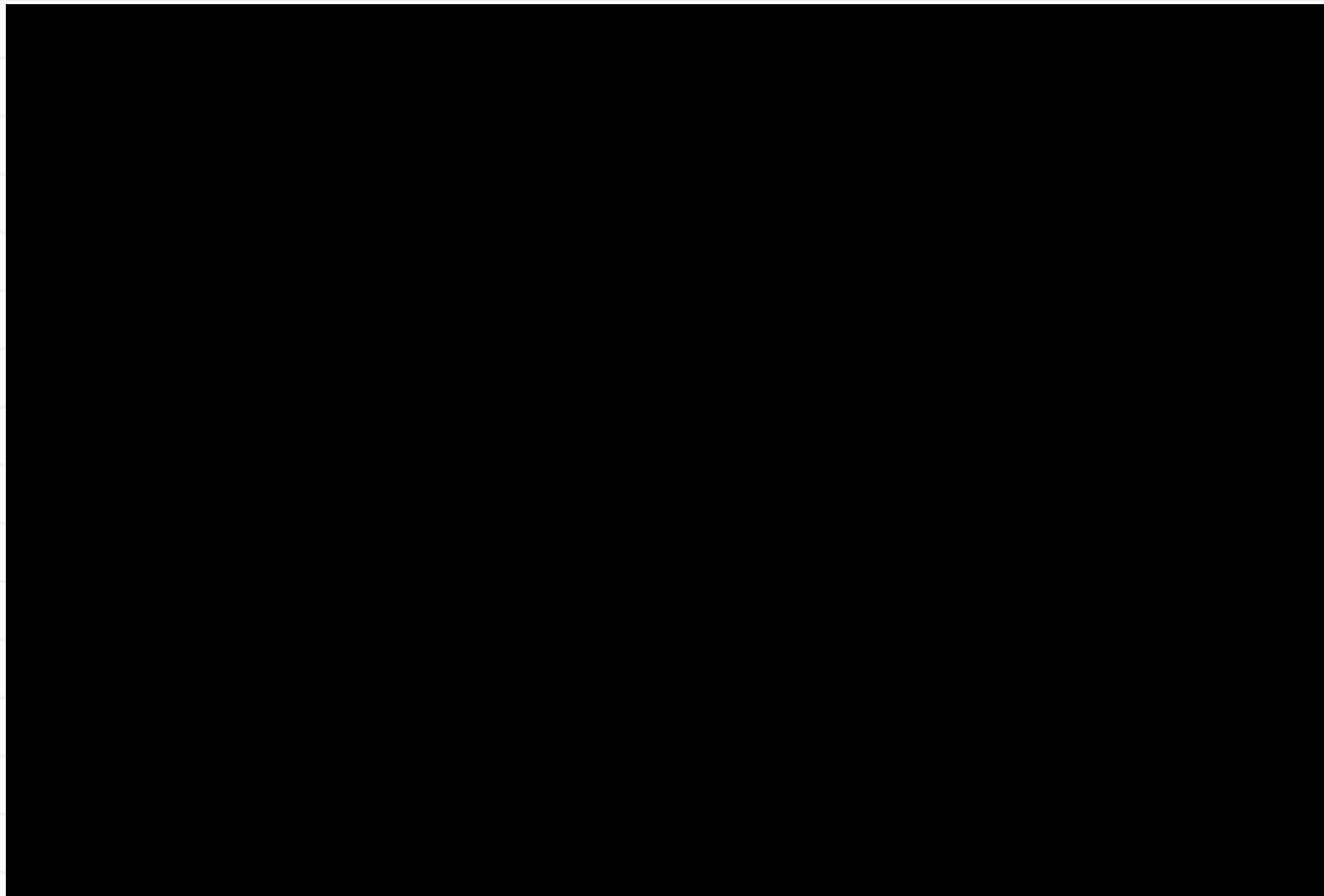
$$G(z) = \arg \max_x D(x)$$

Main source of instability (“mode collapse”).

# Problem with GAN convergence

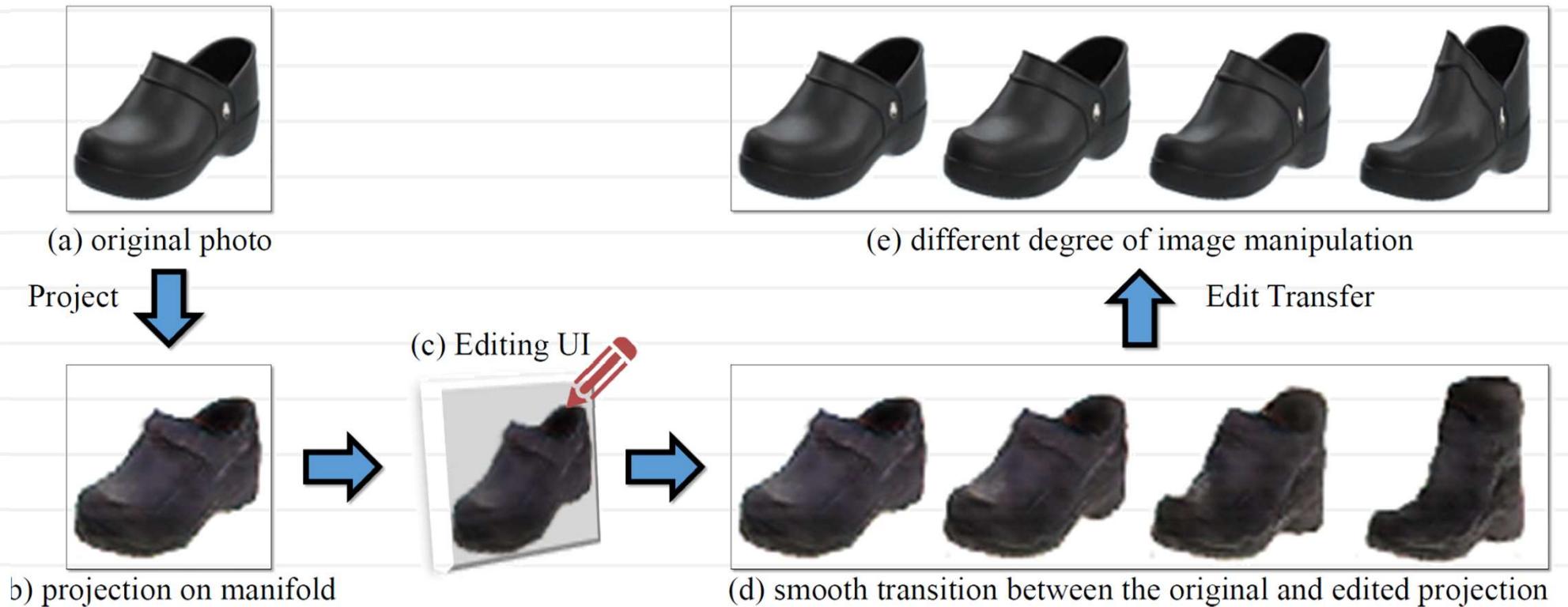
- Gradient updates are not guaranteed to converge to a saddle point
- Local generator collapse (parts of distribution collapse to a single point)
- Low diversity solution still look good for a naked eye
- A lot of recent research is about preventing GAN collapse (e.g. Unrolled GAN [Metz et al. 2016])

# Image editing on image manifold



[Zhu et al. 2016]

# Image editing on image manifold



Projecting  $x^R$  on manifold: 
$$x^* = \arg \min_{x \in \tilde{\mathbb{M}}} \mathcal{L}(x, x^R)$$

Combination of L<sub>2</sub>-distance and VGG-activation-based distance

[Zhu et al. 2016]

# Image editing on image manifold

Projecting  $x^R$  on manifold:

$$x^* = \arg \min_{x \in \tilde{\mathbb{M}}} \mathcal{L}(x, x^R)$$

Method of projection:

- Optimizing in the latent space  $z$  (hard, slow)
- Predicting  $z$  by a feedforward network:

$$\theta_P^* = \arg \min_{\theta_P} \sum_n \mathcal{L}(G(P(x_n^R; \theta_P)), x_n^R)$$

- Both: predict initialization and then optimize

[Zhu et al. 2016]

# Image editing on image manifold

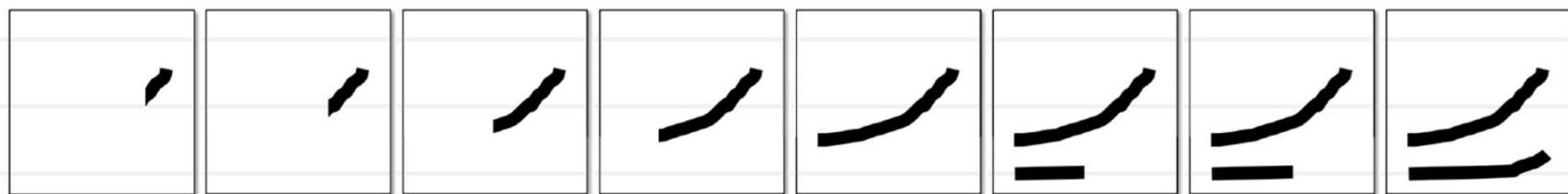


[Zhu et al. 2016]

# Image editing on image manifold

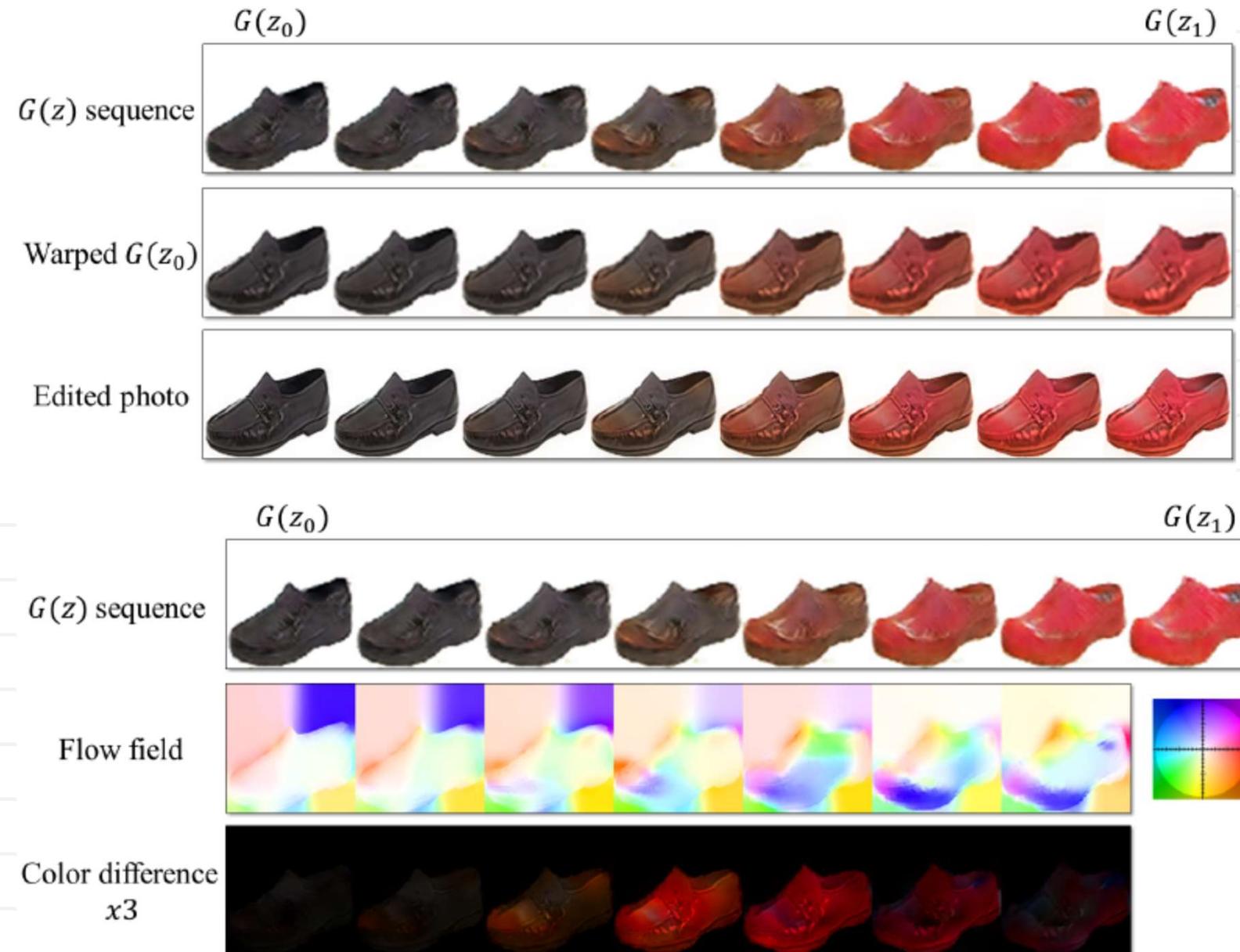
Editing process:

$$z^* = \arg \min_{z \in \mathbb{Z}} \left\{ \underbrace{\sum_g \|f_g(G(z)) - v_g\|^2}_{\text{data term}} + \underbrace{\lambda_s \cdot \|z - z_0\|^2}_{\text{manifold smoothness}} + E_D \right\}$$



[Zhu et al. 2016]

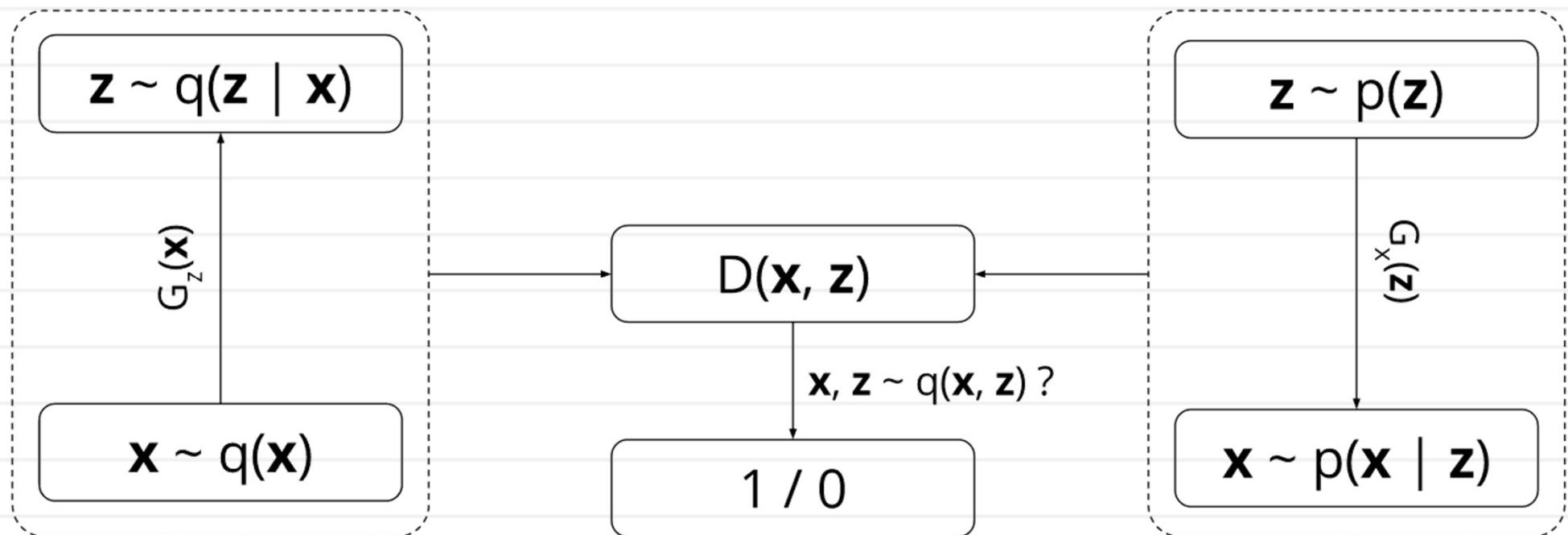
# Transferring edits



[Zhu et al. 2016]

# Adversarially-learned inference

Idea: train  $z \rightarrow x$  and  $x \rightarrow z$  mappings simultaneously



[Dumolin et al .2016]

# Adversarially-learned inference



Samples  $z \rightarrow x$ :



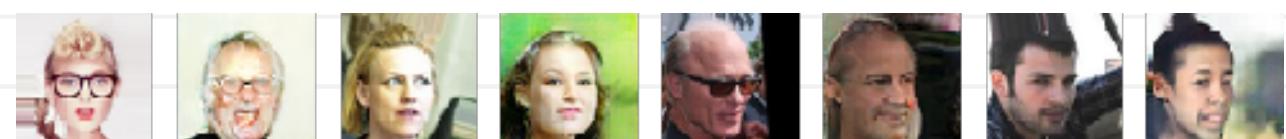
[Dumolin et al .2016]

# Adversarially-learned inference



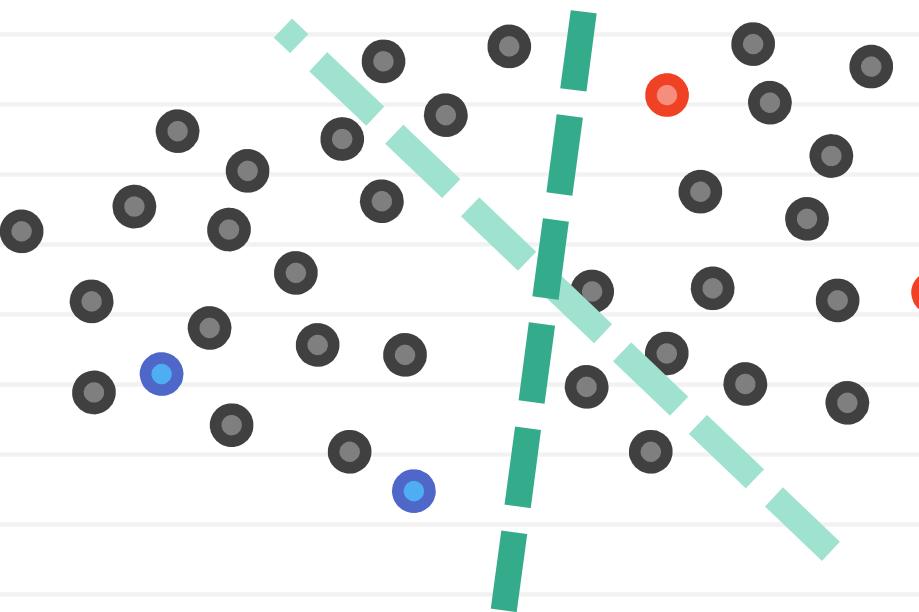
Reconstructions

$X \rightarrow Z \rightarrow X$ :

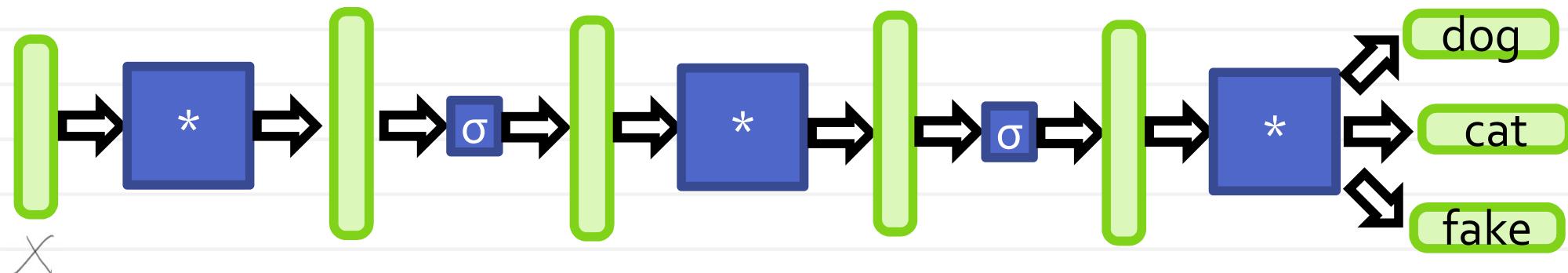


[Dumolin et al .2016]

# Semi-supervised learning



New discriminator:



[Salimans 2016]

# Semi-supervised learning

$$L = -\mathbb{E}_{\mathbf{x}, y \sim p_{\text{data}}(\mathbf{x}, y)} [\log p_{\text{model}}(y|\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim G} [\log p_{\text{model}}(y = K+1|\mathbf{x})]$$

=  $L_{\text{supervised}} + L_{\text{unsupervised}}$ , where

$$L_{\text{supervised}} = -\mathbb{E}_{\mathbf{x}, y \sim p_{\text{data}}(\mathbf{x}, y)} \log p_{\text{model}}(y|\mathbf{x}, y < K+1)$$

$$L_{\text{unsupervised}} = -\{\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} \log[1 - p_{\text{model}}(y = K+1|\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G} \log[p_{\text{model}}(y = K+1|\mathbf{x})]\}$$

## Perturbed MNIST results:

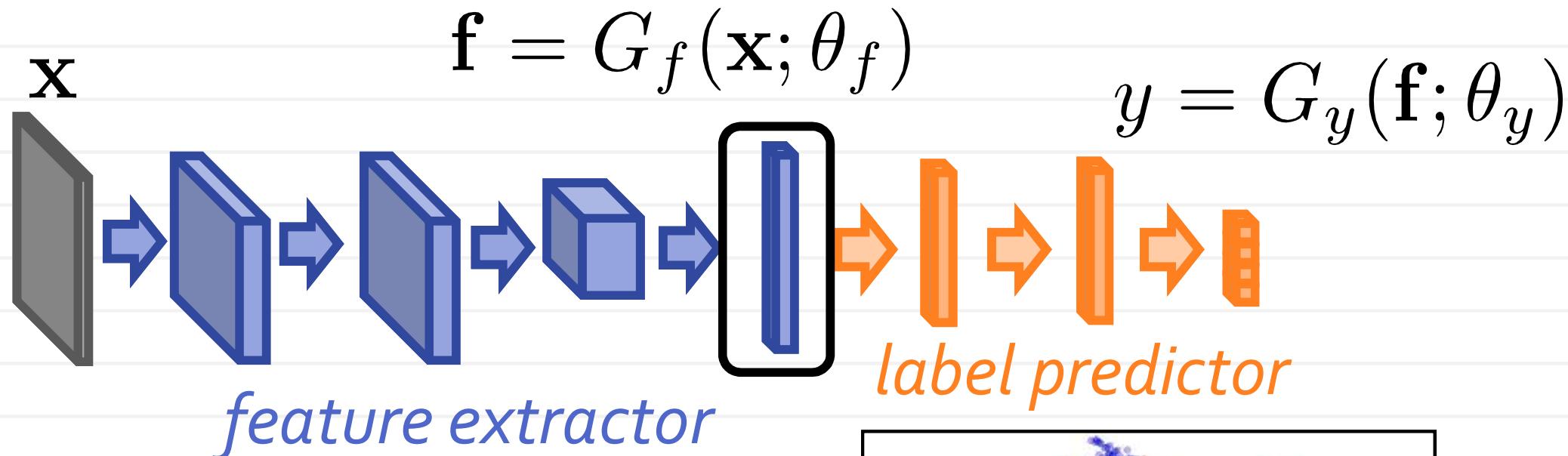
Model	Number of incorrectly predicted test examples for a given number of labeled samples			
	20	50	100	200
DGN [21]			$333 \pm 14$	
Virtual Adversarial [22]			$212$	
CatGAN [14]			$191 \pm 10$	
Skip Deep Generative Model [23]			$132 \pm 7$	
Ladder network [24]			$106 \pm 37$	
Auxiliary Deep Generative Model [23]			$96 \pm 2$	
Our model	$1677 \pm 452$	$221 \pm 136$	$93 \pm 6.5$	$90 \pm 4.2$
Ensemble of 10 of our models	$1134 \pm 445$	$142 \pm 96$	$86 \pm 5.6$	$81 \pm 4.3$

[Salimans 2016]

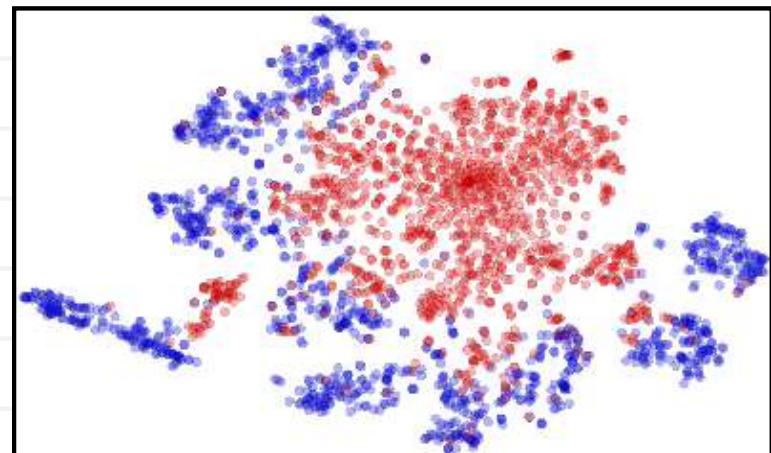
# Example: Internet images -> Webcam sensor



# Domain shift in a deep architecture



When trained on source  
only, feature  
distributions do not  
match:

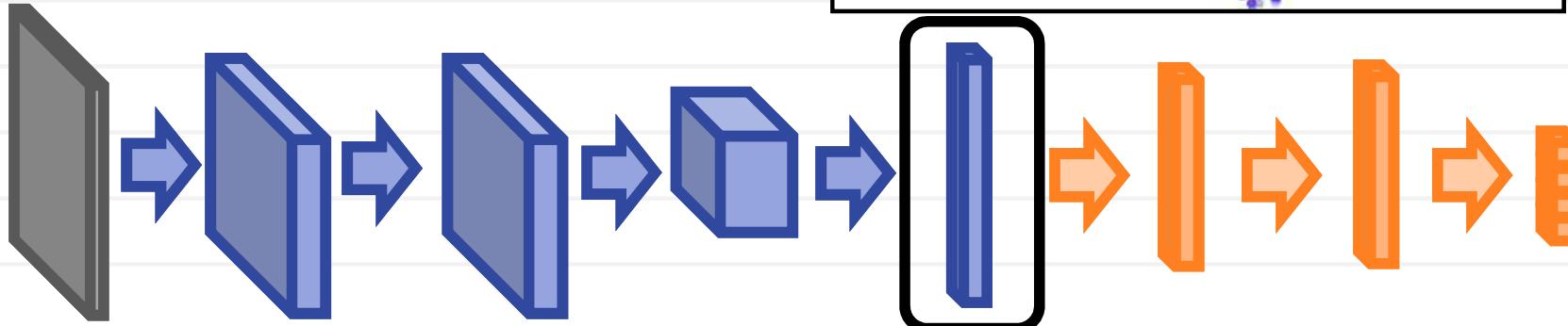
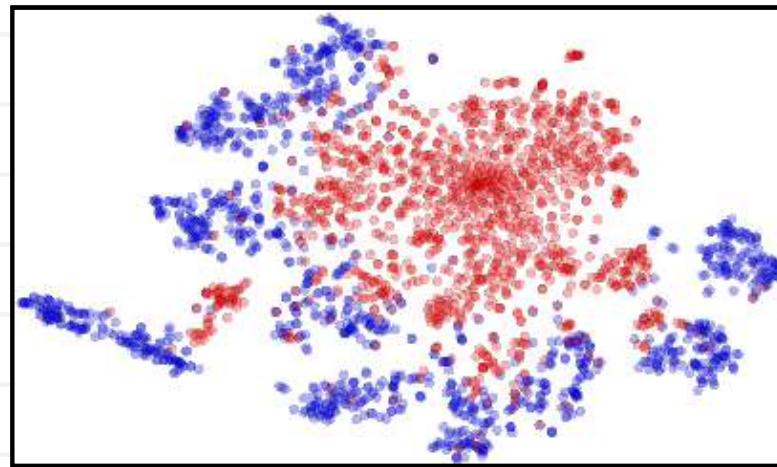


$$S(\mathbf{f}) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim S(\mathbf{x})\}$$

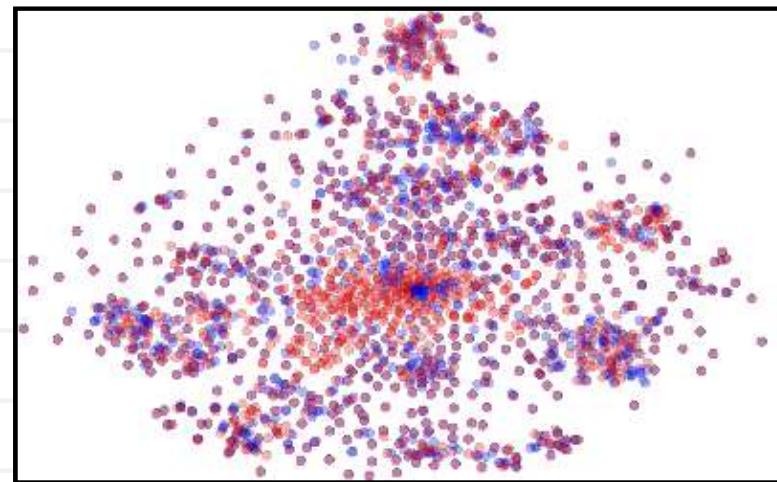
$$T(\mathbf{f}) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim T(\mathbf{x})\}$$

# Idea 1: domain-invariant features wanted

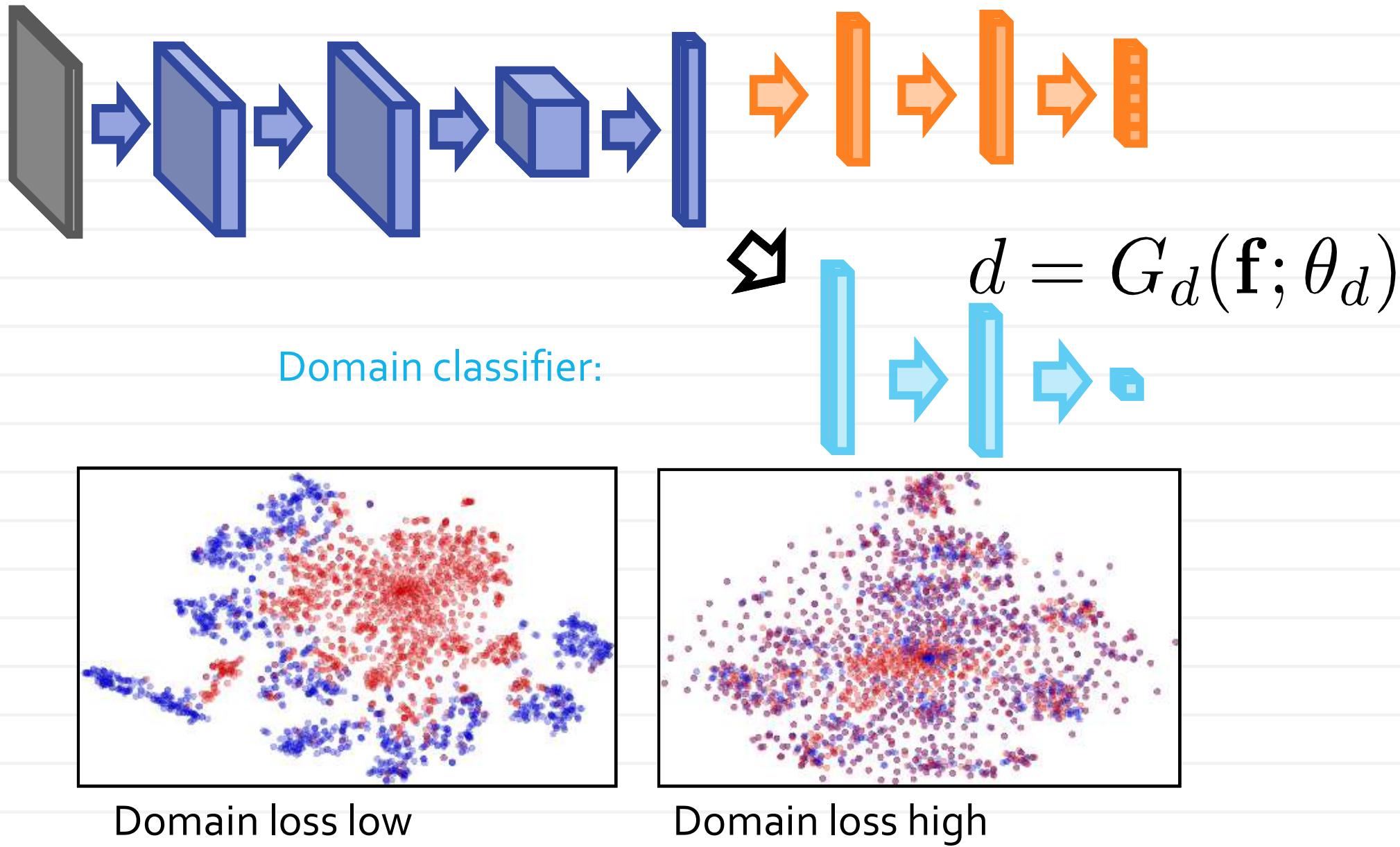
Feature distribution without adaptation:



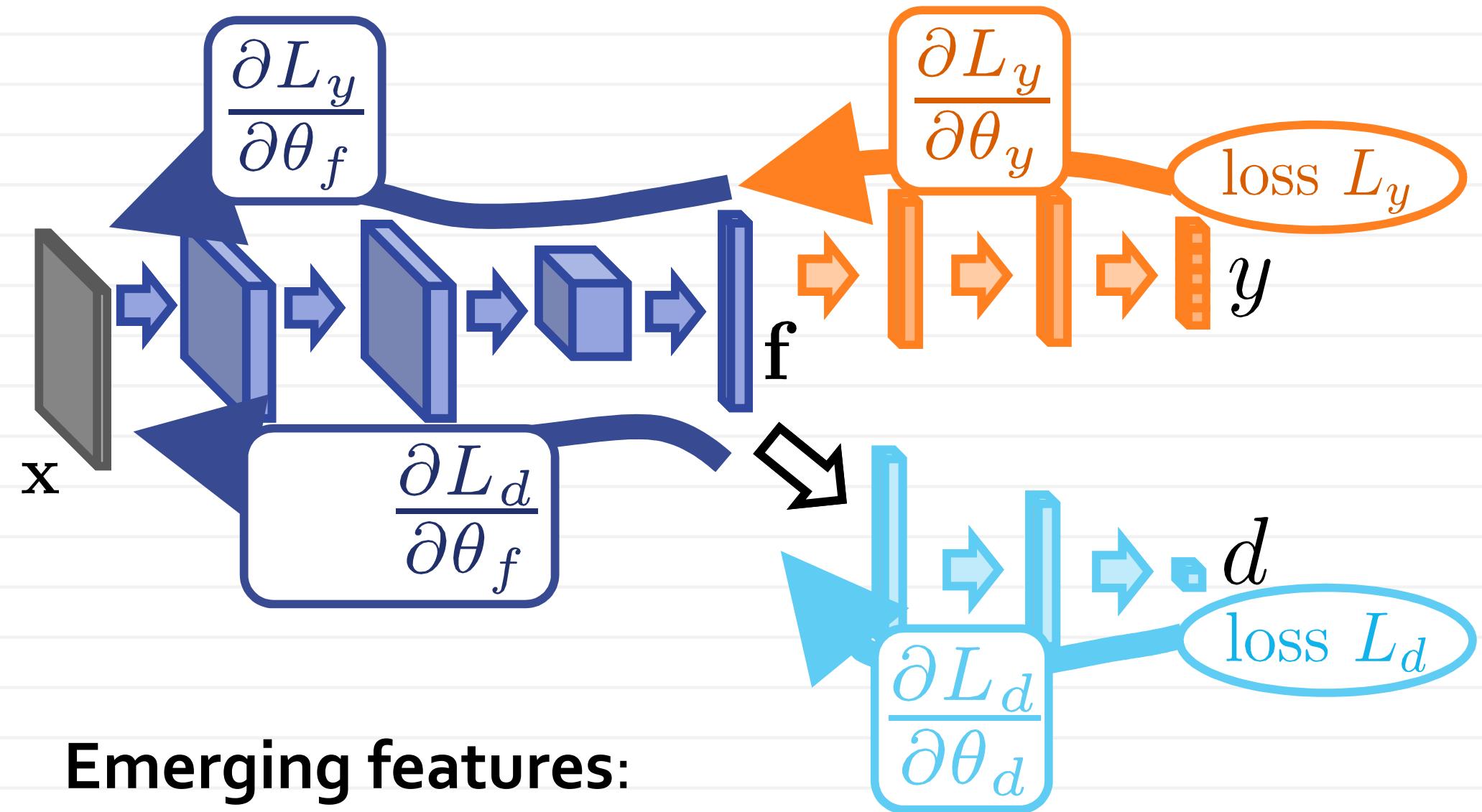
Our goal (after adaptation):



# Idea 2: measuring domain shift



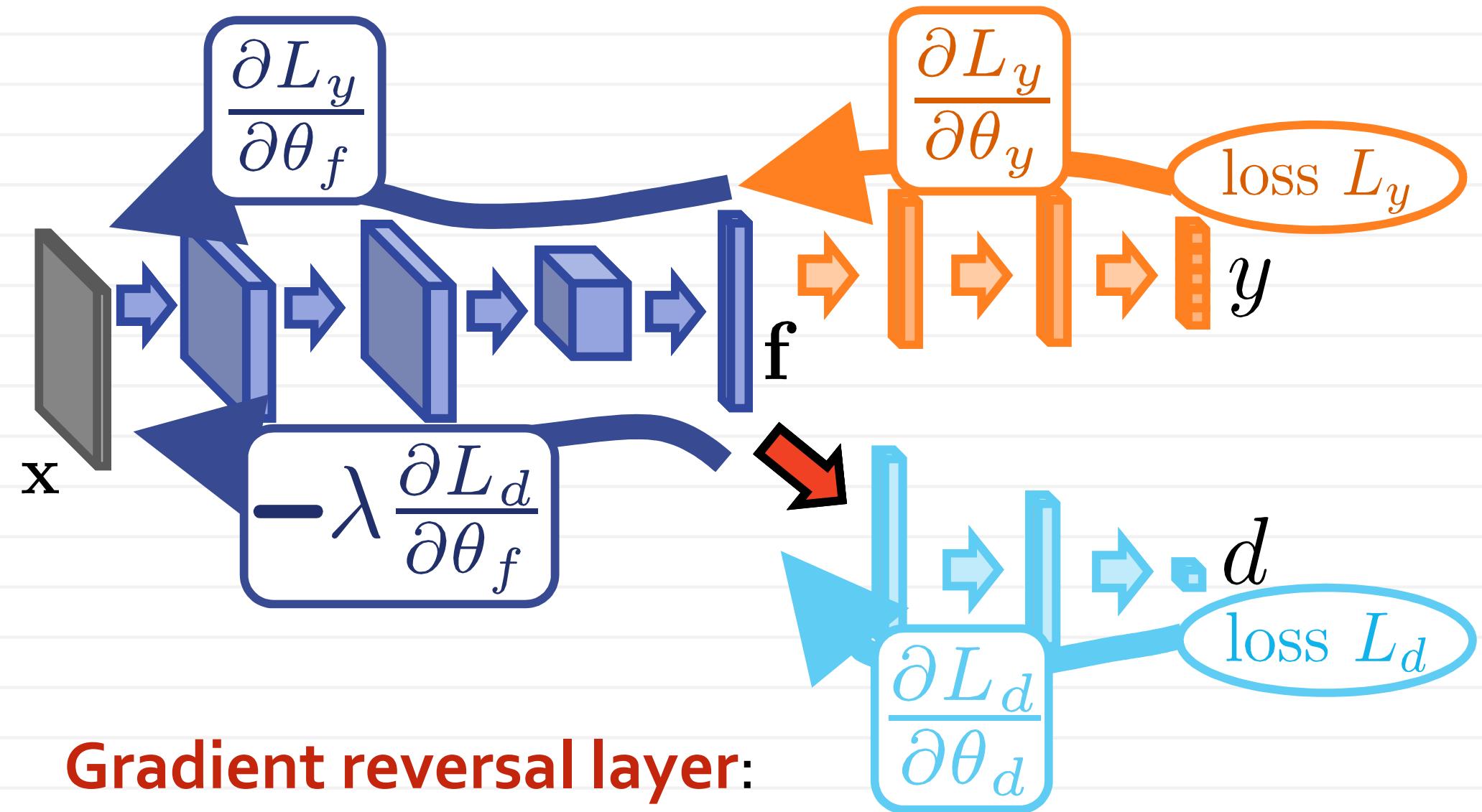
# Idea 3: minimizing domain shift



## Emerging features:

- Discriminative (good for predicting  $y$ )
- Domain-discriminative (good for predicting  $d$ )

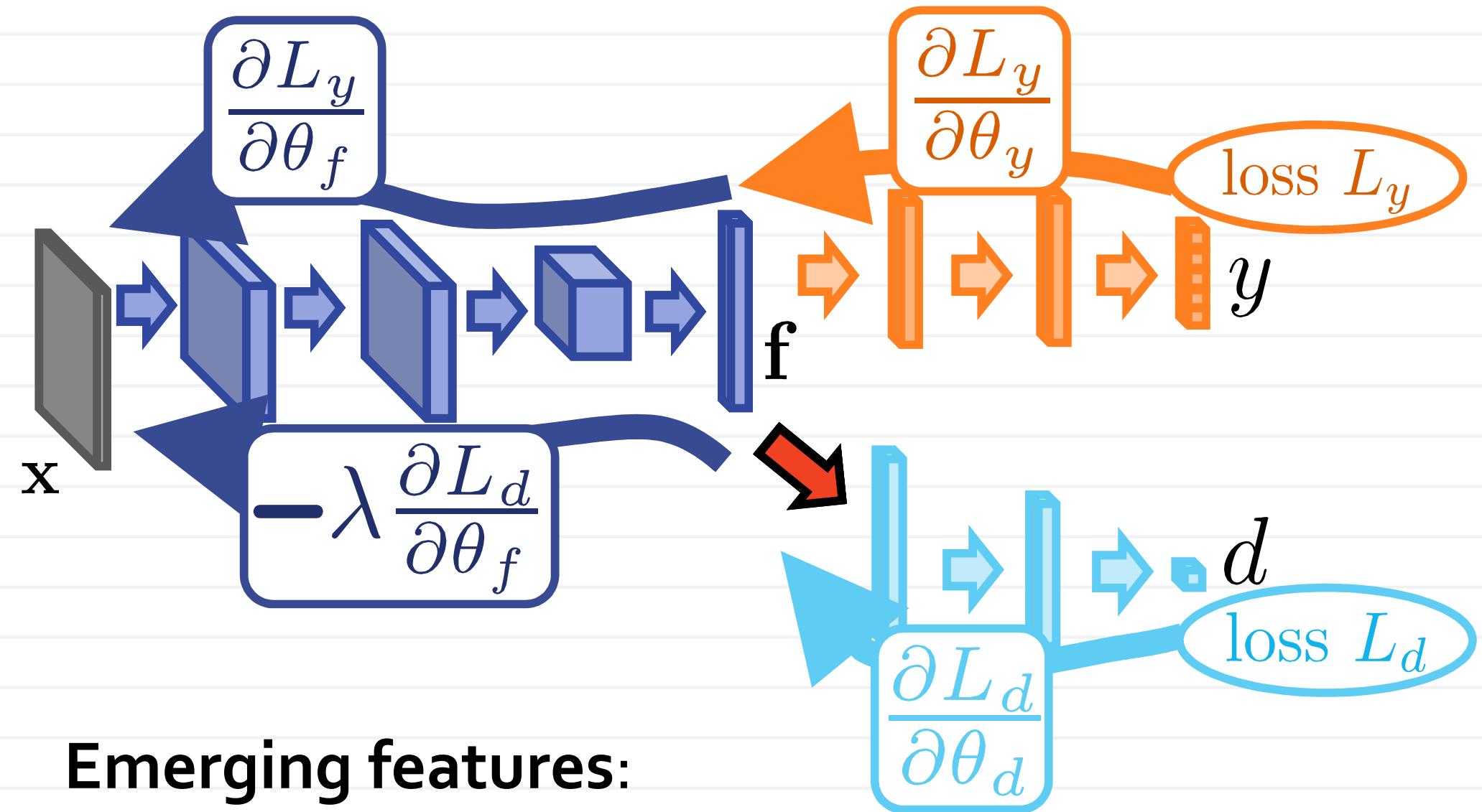
# Idea 3: minimizing domain shift



## Gradient reversal layer:

- Copies data without change at forwardprop
- Multiplies the gradient by  $-\lambda$  at backprop

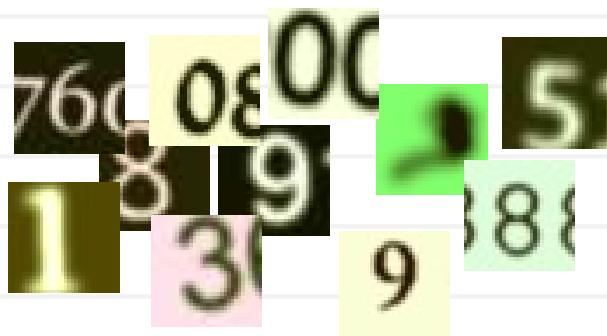
# Idea 3: minimizing domain shift



## Emerging features:

- Discriminative (good for predicting  $y$ )
- Domain-invariant (not good for predicting  $d$ )

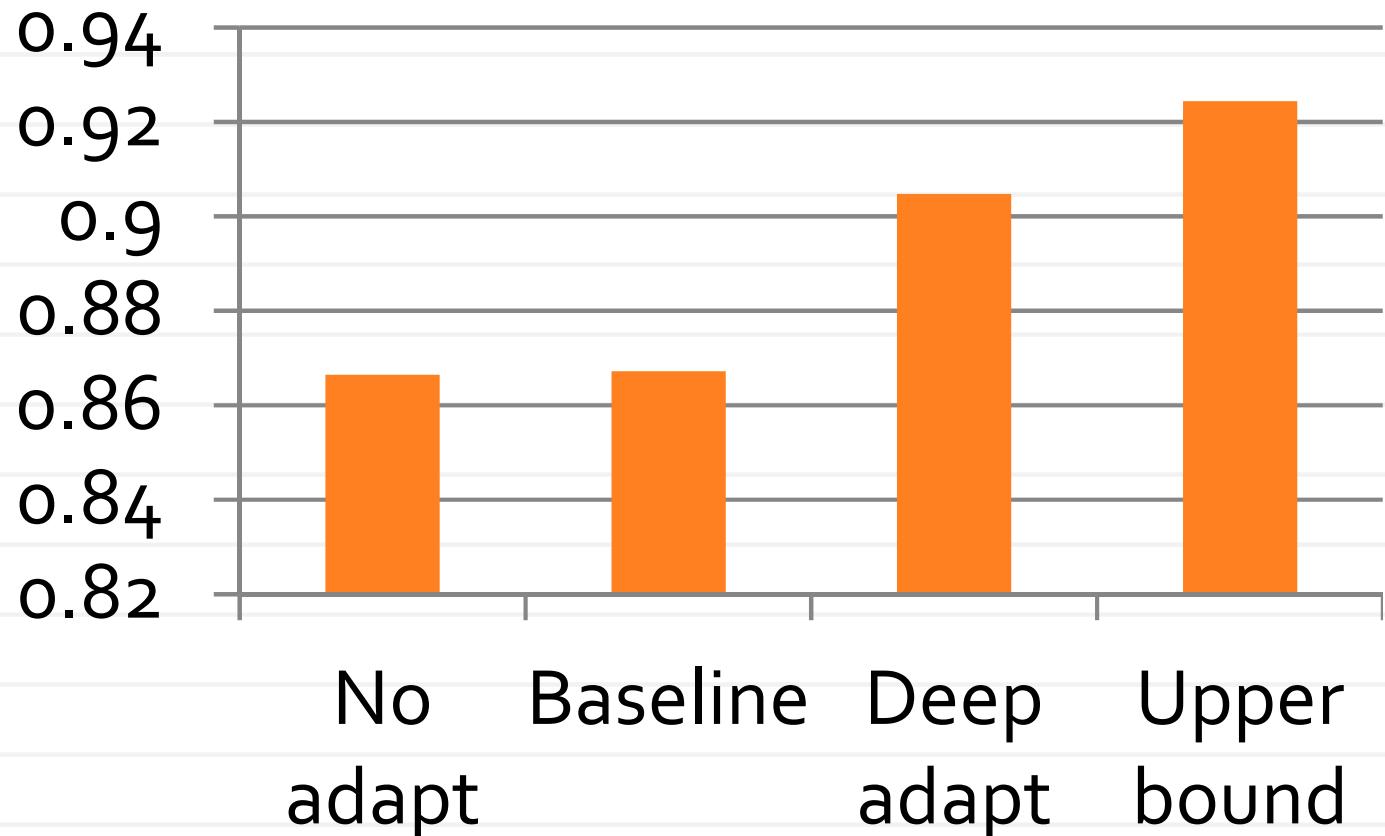
# Example: from synthetic to real



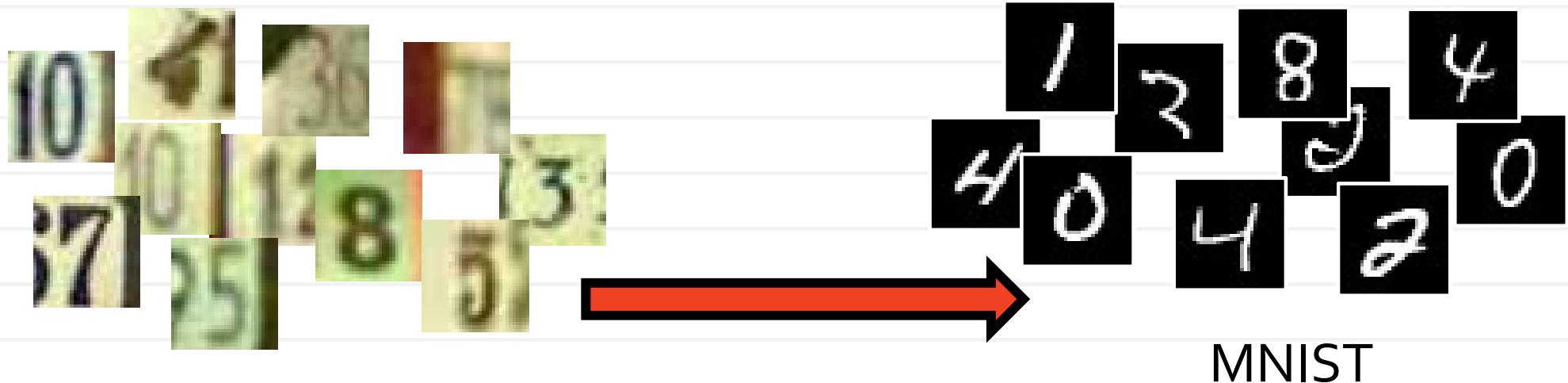
“Windows digits”



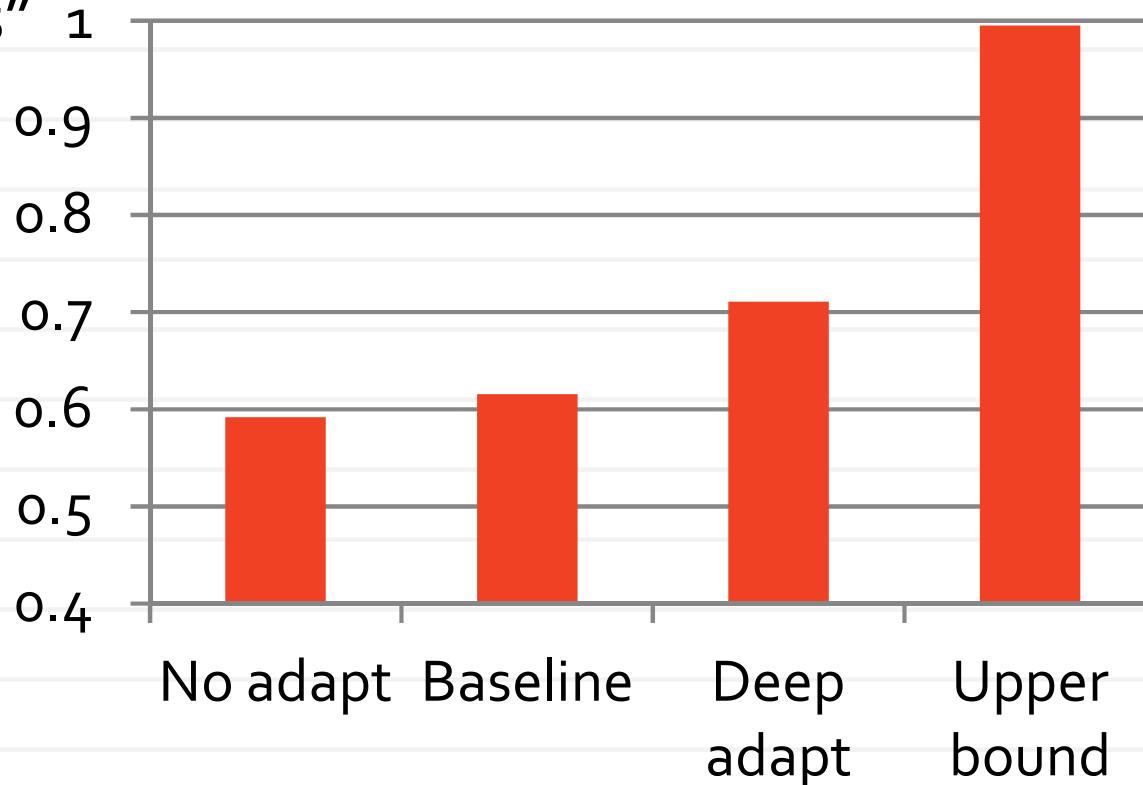
“House numbers”



# Example: large gap

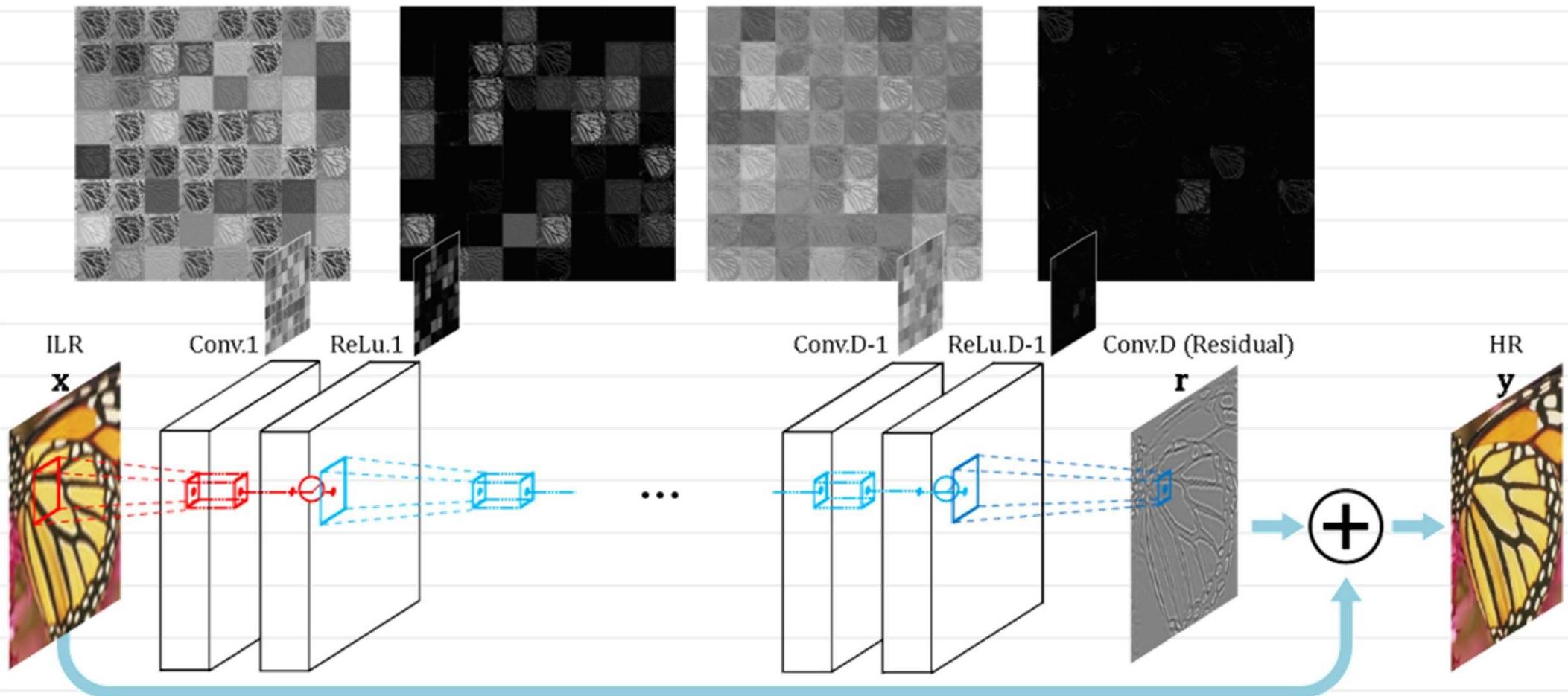


"House numbers"



Reverse  
direction  
does not  
work ☹

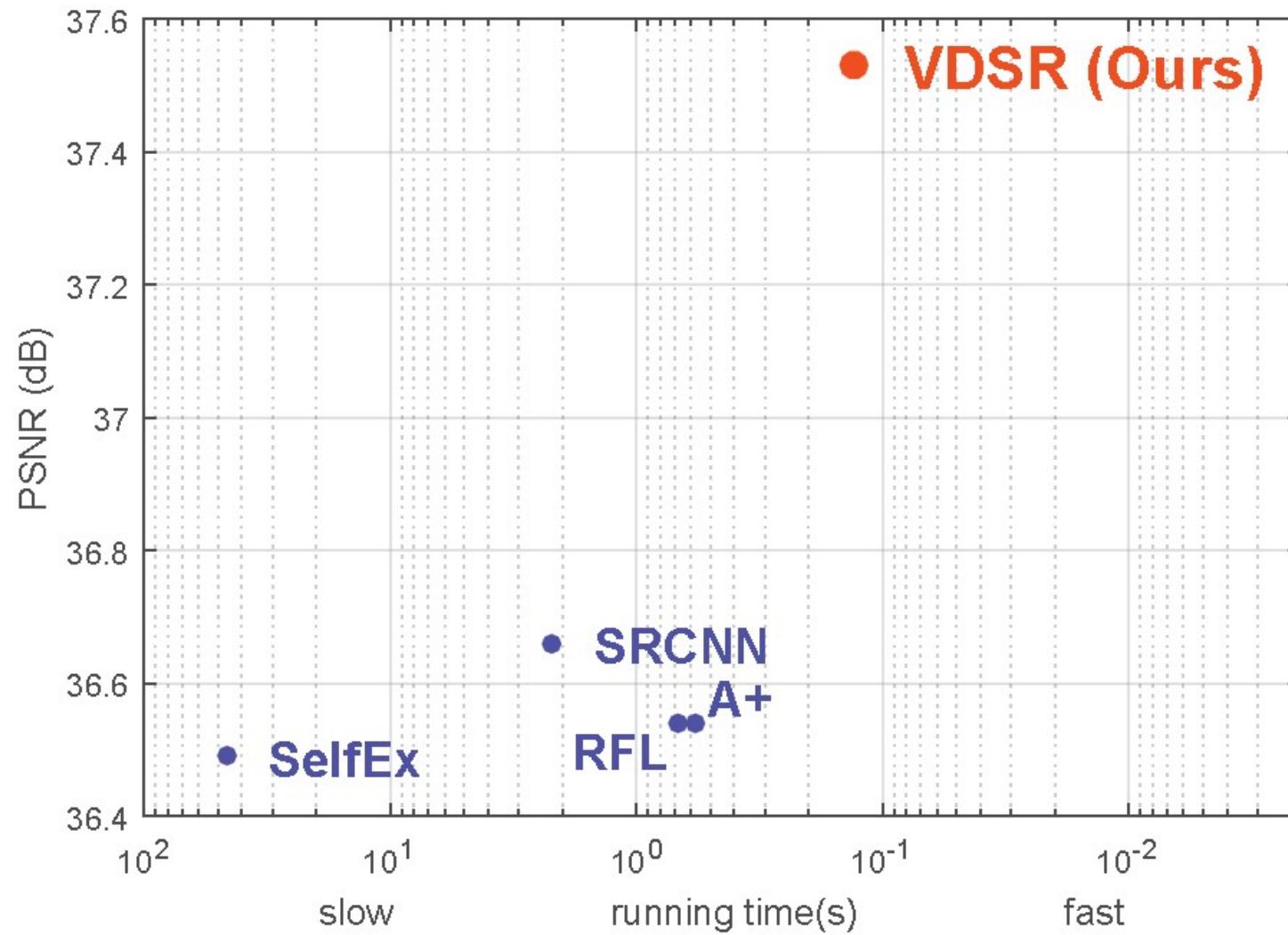
# Image superresolution with ConvNets



- 19 layers with  $3 \times 3$  filters (VGG inspired)

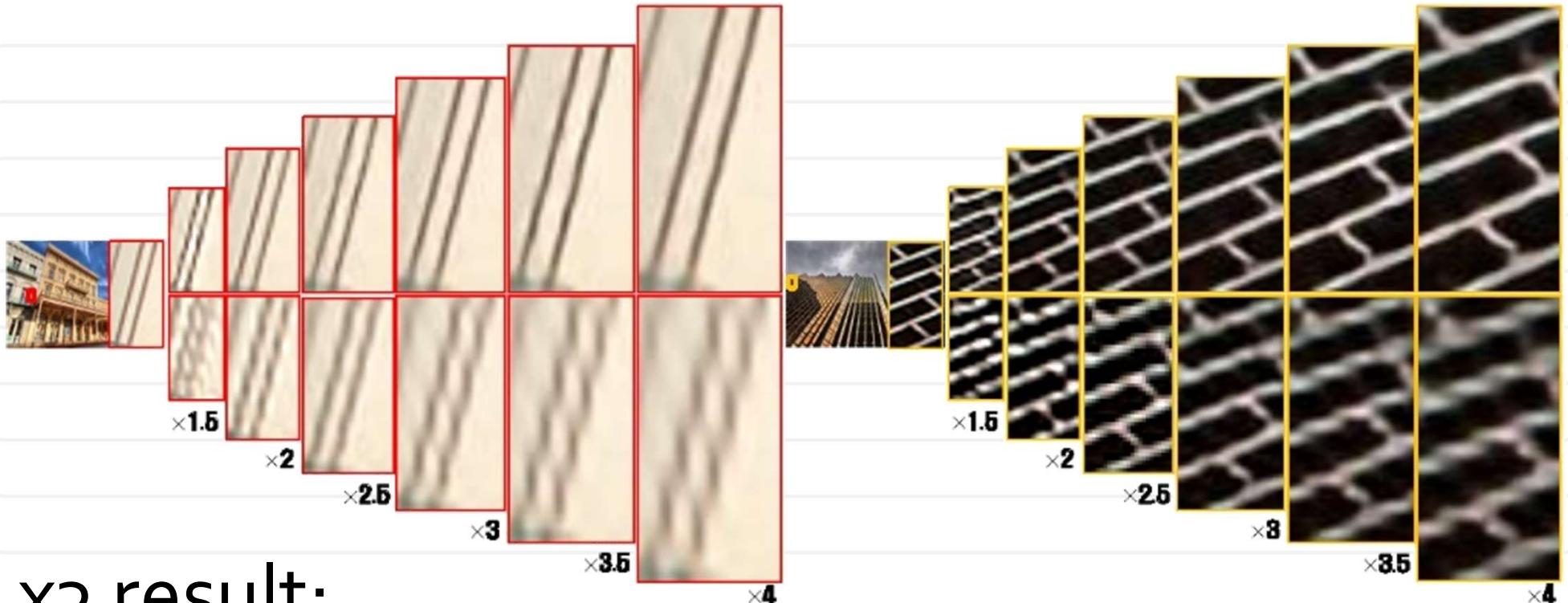
[Kim et al. CVPR16]

# Image superresolution with ConvNets

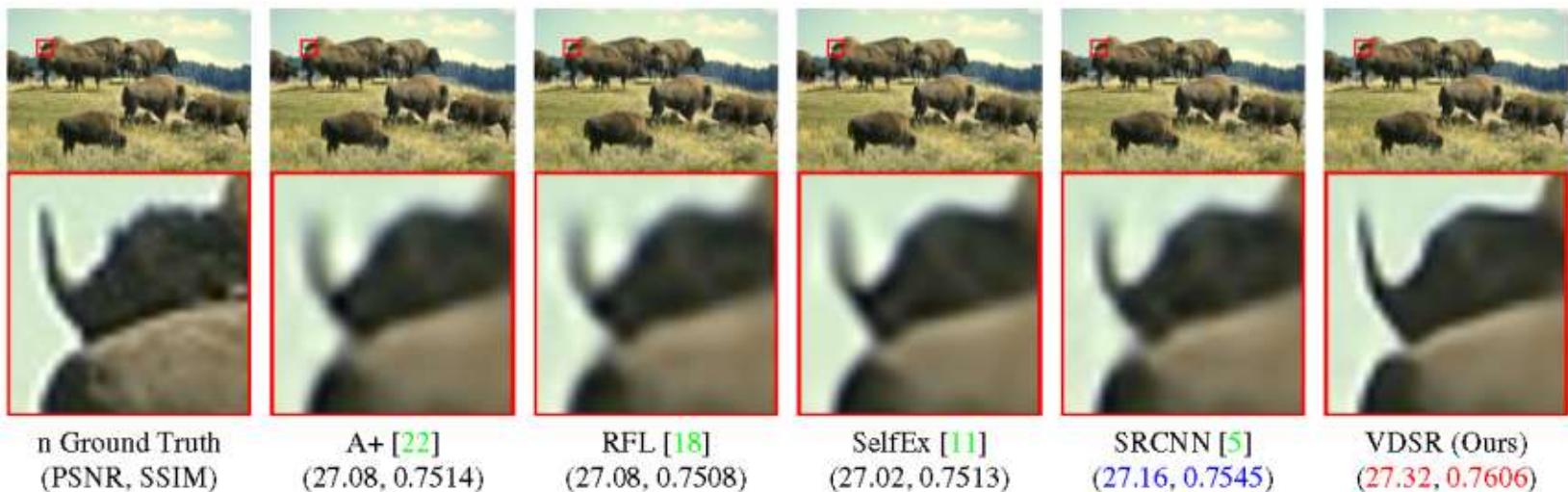


[Kim et al. CVPR16]

# Image superresolution with ConvNets



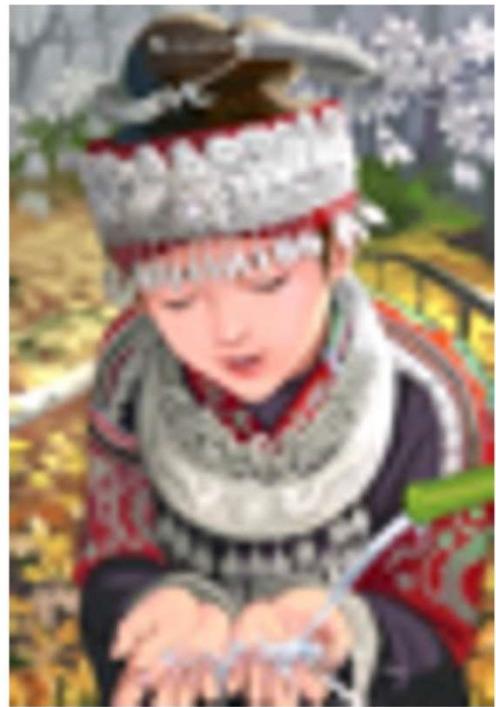
x3 result:



[Kim et al. CVPR16]

# Superresolution with GAN

bicubic  
(21.59dB/0.6423)



SRResNet  
(23.53dB/0.7832)



SRGAN  
(21.15dB/0.6868)



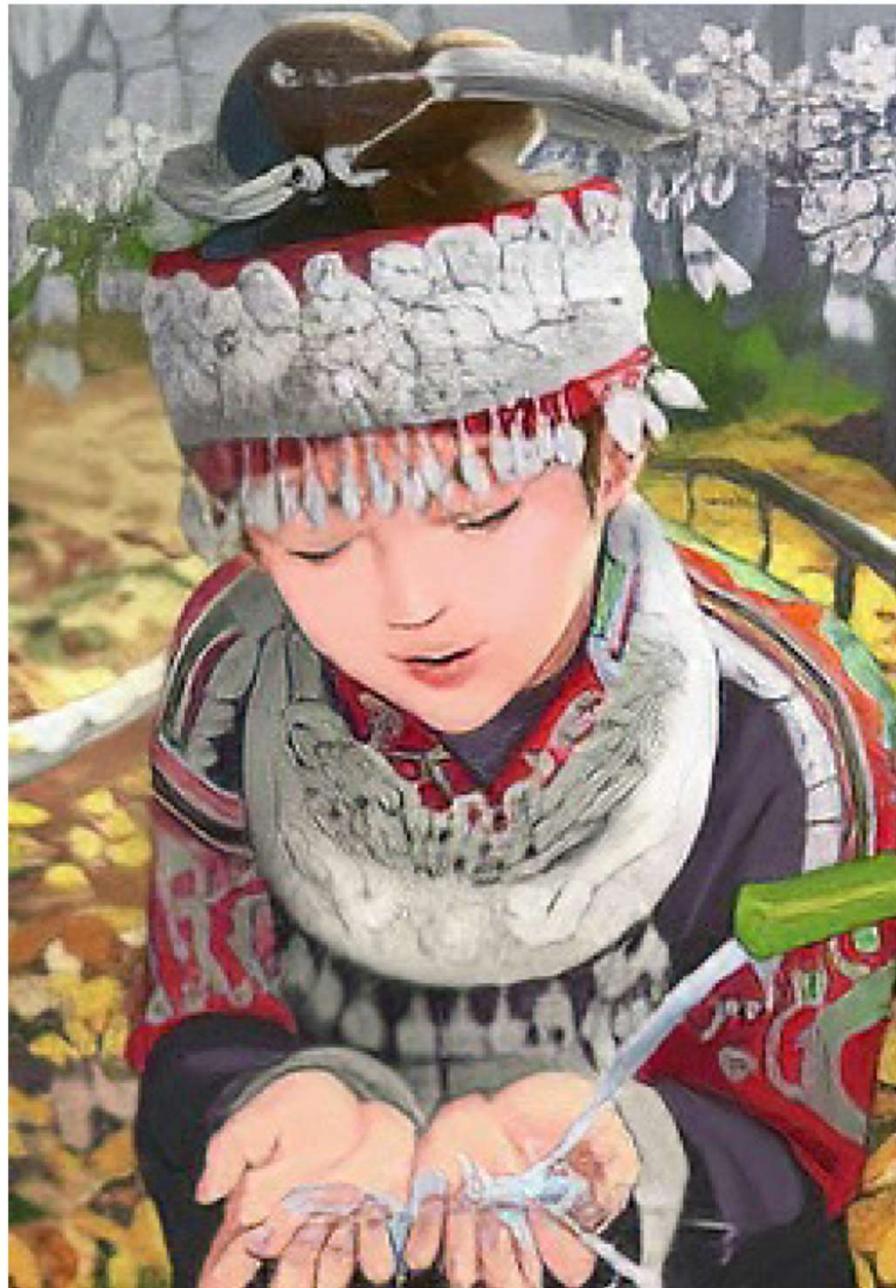
original



Adding GAN term into the loss

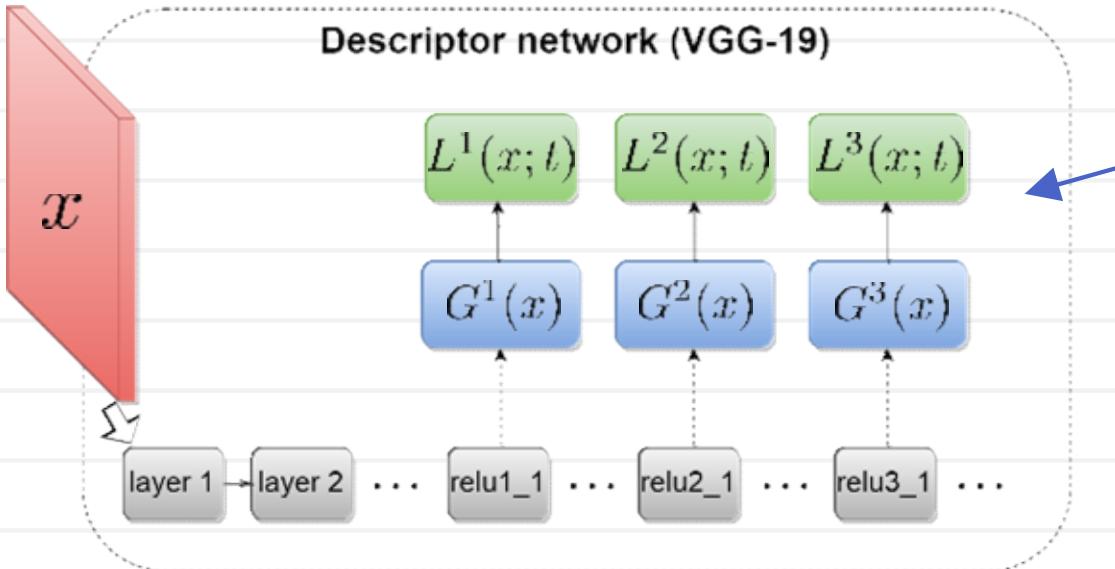
[Ledig et al. 2016]

# Superresolution with GAN



[Ledig et al. 2016]

# Gatys et al. texture loss



Activations at layer  $l$ :  $F^l(t)$

Gram-matrix statistics:  $G_{ij}^l(t) = \sum_{k=1}^{M_l N_l} F_{ik}^l(t) F_{jk}^l(t)$

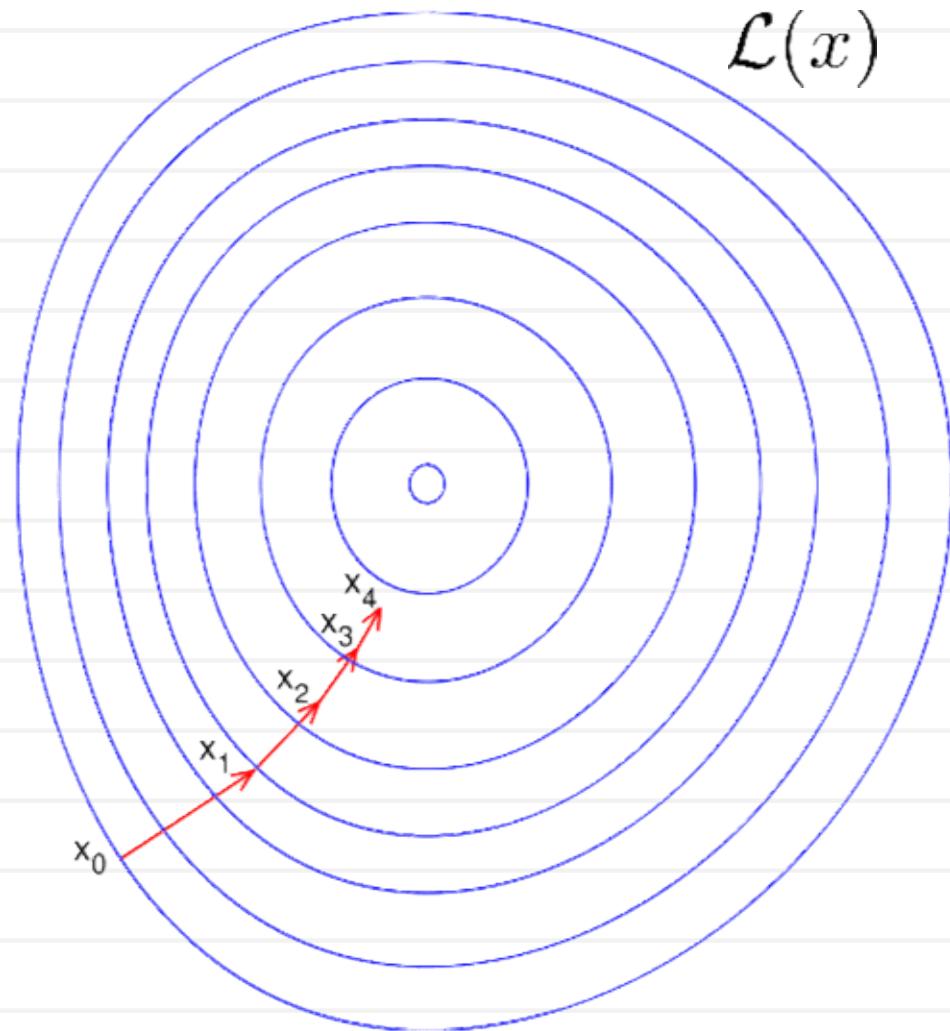
Texture loss (dissimilarity between  $t$  and  $x$ ):

$$\mathcal{L}_{texture}(x; t) = \sum_l ||G^l(t) - G^l(x)||_2^2$$

[Gatys, Ecker, Bethge 2015]

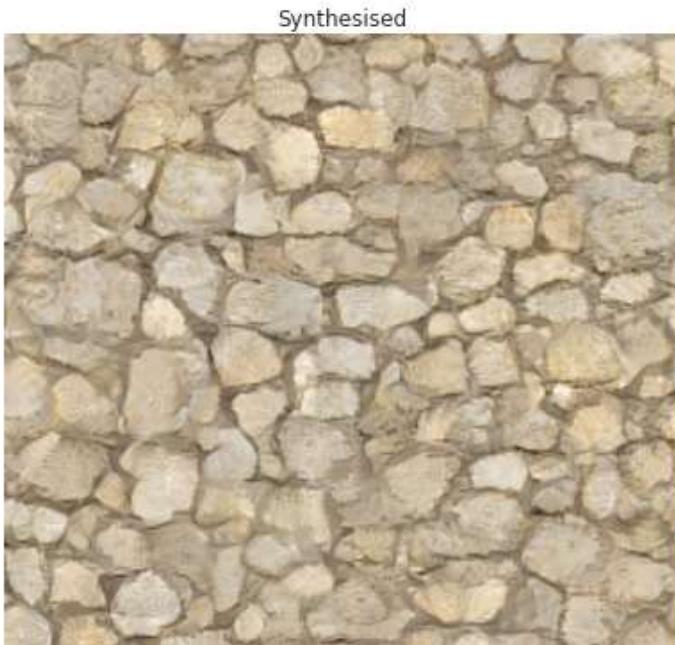
# Texture generation by optimization

$$x^* = \arg \min_x \mathcal{L}(x)$$



[Gatys, Ecker, Bethge 2015]

# Pre-image texture synthesis



[Gatys, Ecker, Bethge 2015]

# Texture network

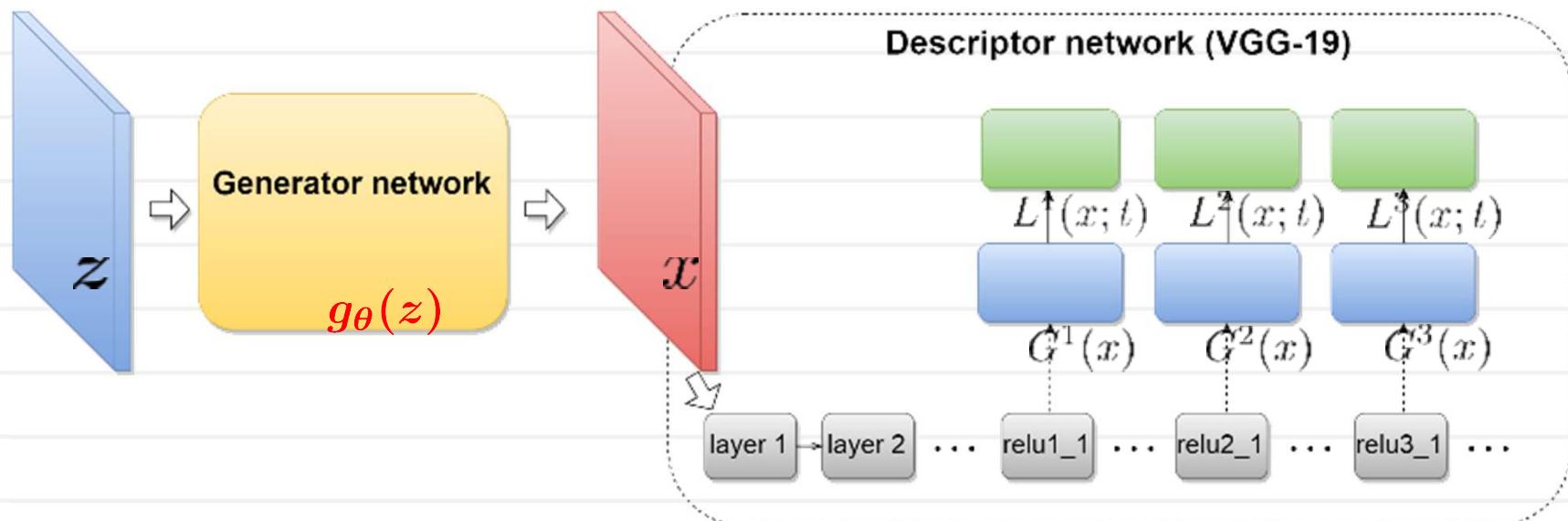


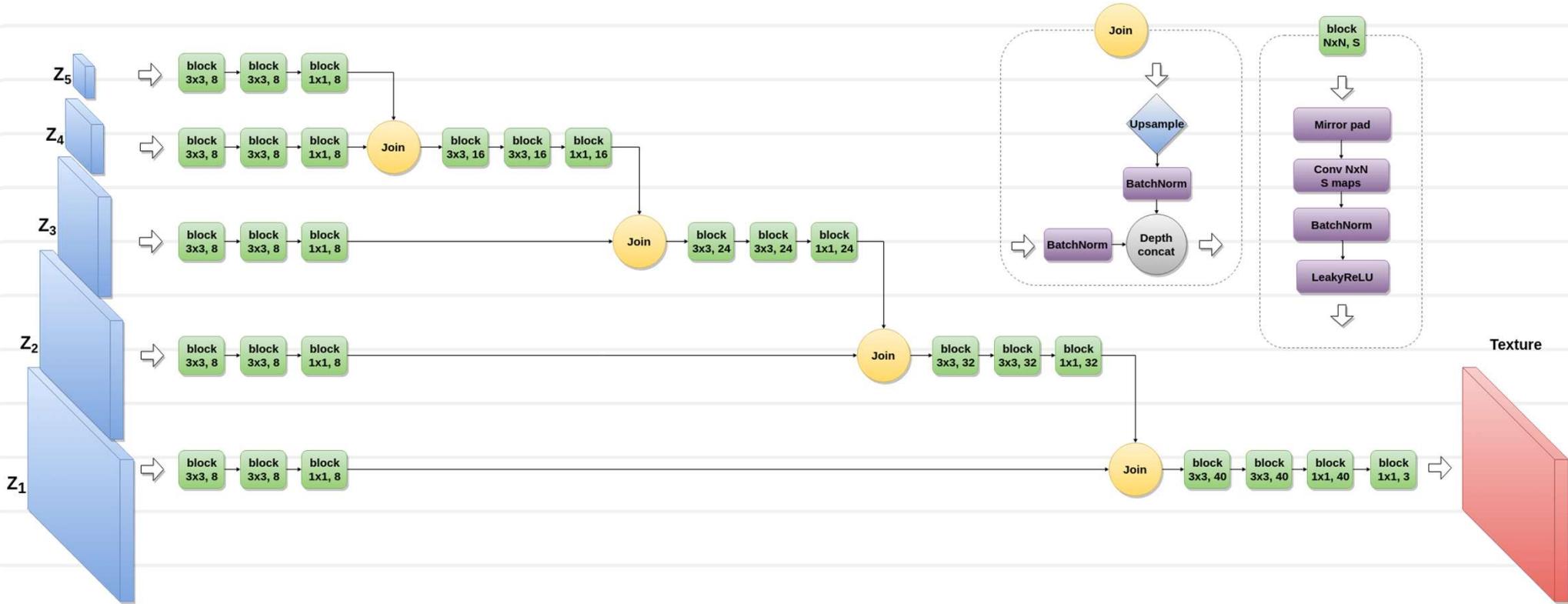
Image generation:  $x = g_\theta(z), \quad z \sim U(0, 1)$

Optimization task becomes:

$$\min_{\theta} \mathbb{E} \mathcal{L}_{texture}(g_\theta(z); t), \quad z \sim U(0, 1)$$

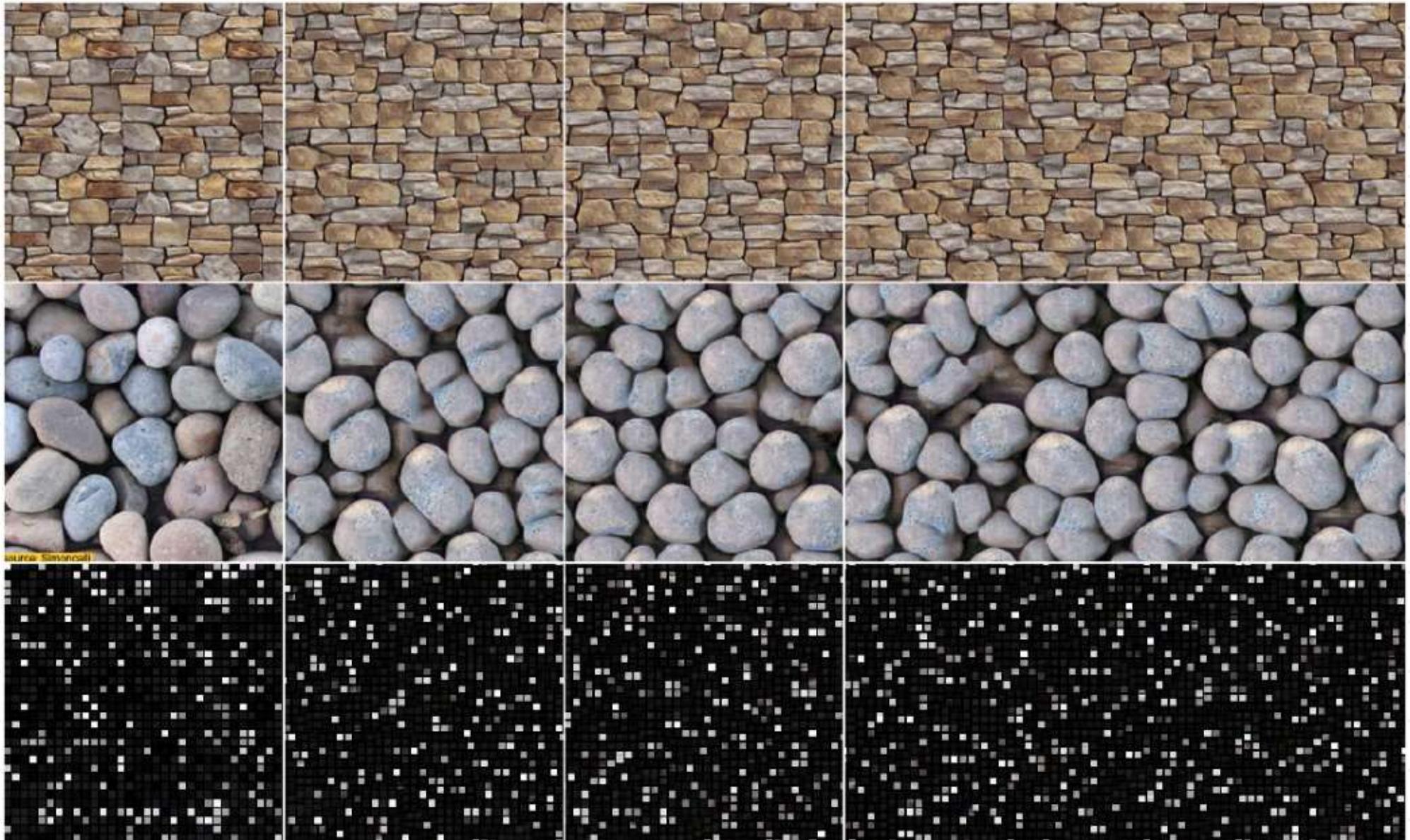
$$\theta^{k+1} = \theta^k - \alpha \frac{\partial \mathcal{L}(g_\theta(z); t)}{\partial \theta}$$

# Generator networks details



- Multi-scale approach simplifies learning
- Looks scary, but still feedforward and fully convolutional
- Bias terms removed to avoid overfitting

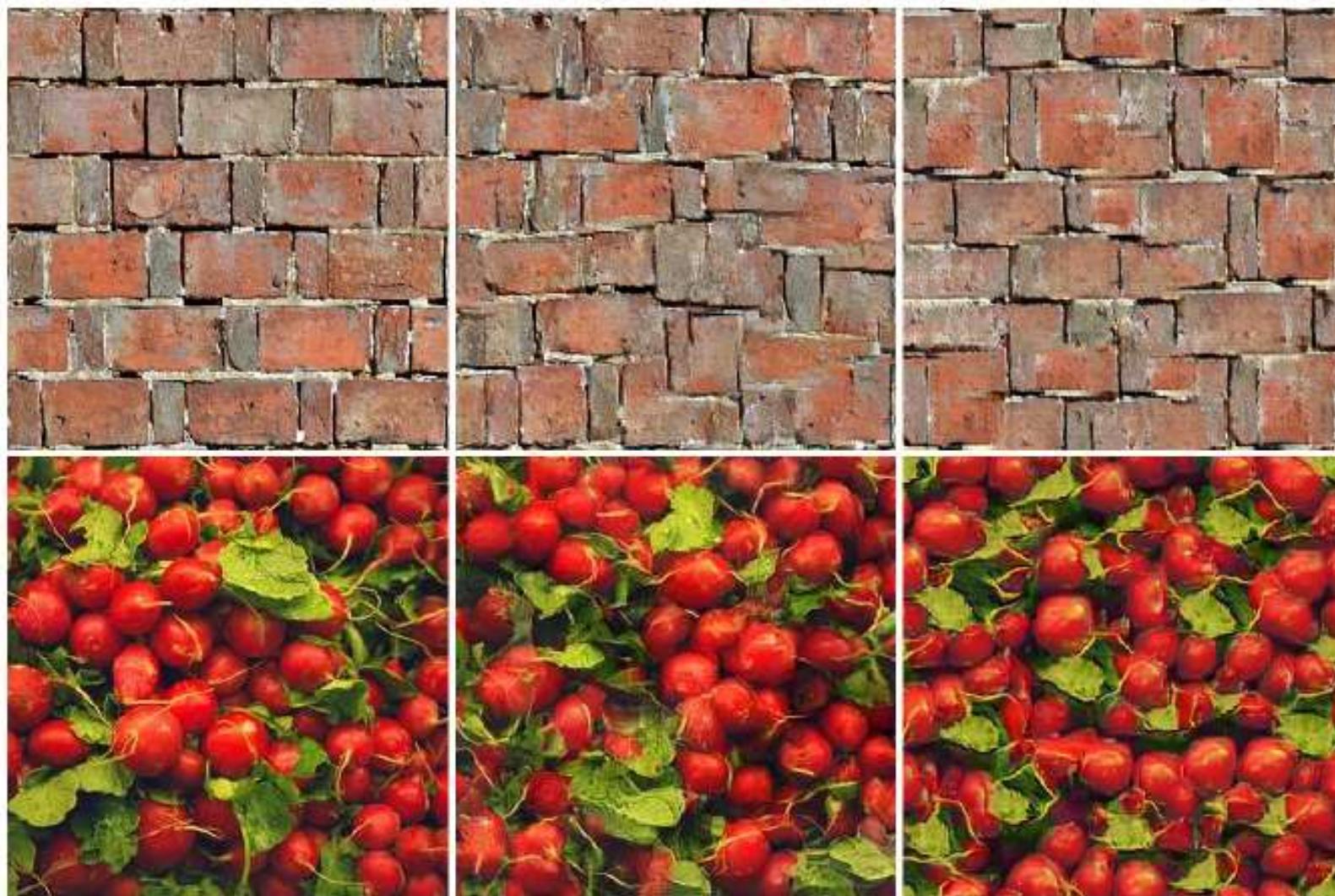
# Feed-forward texture synthesis



sample

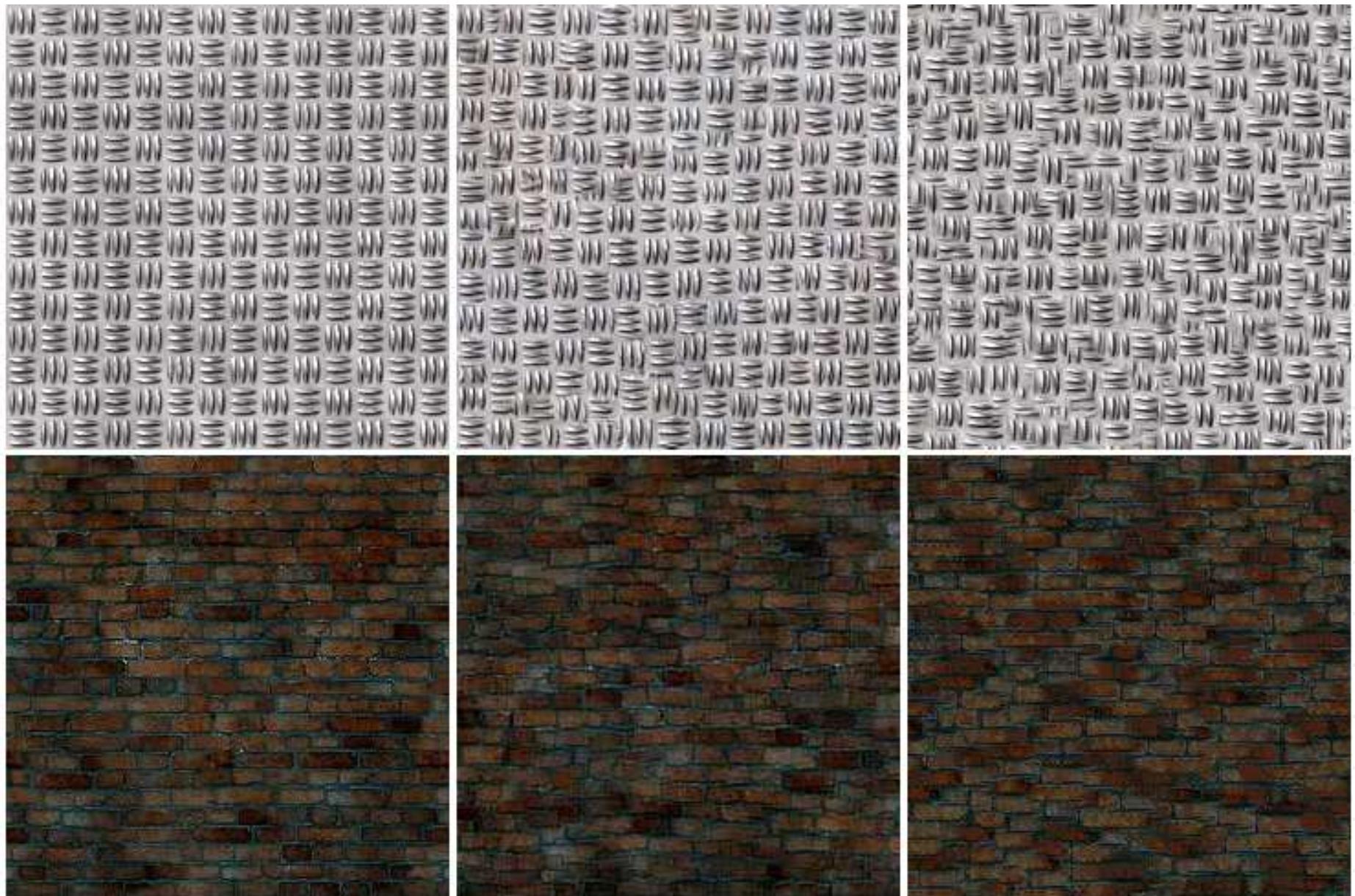
Results

# Pre-image vs. feed-forward



sample      Optimized      Feed-forward

# Pre-image vs. feed-forward

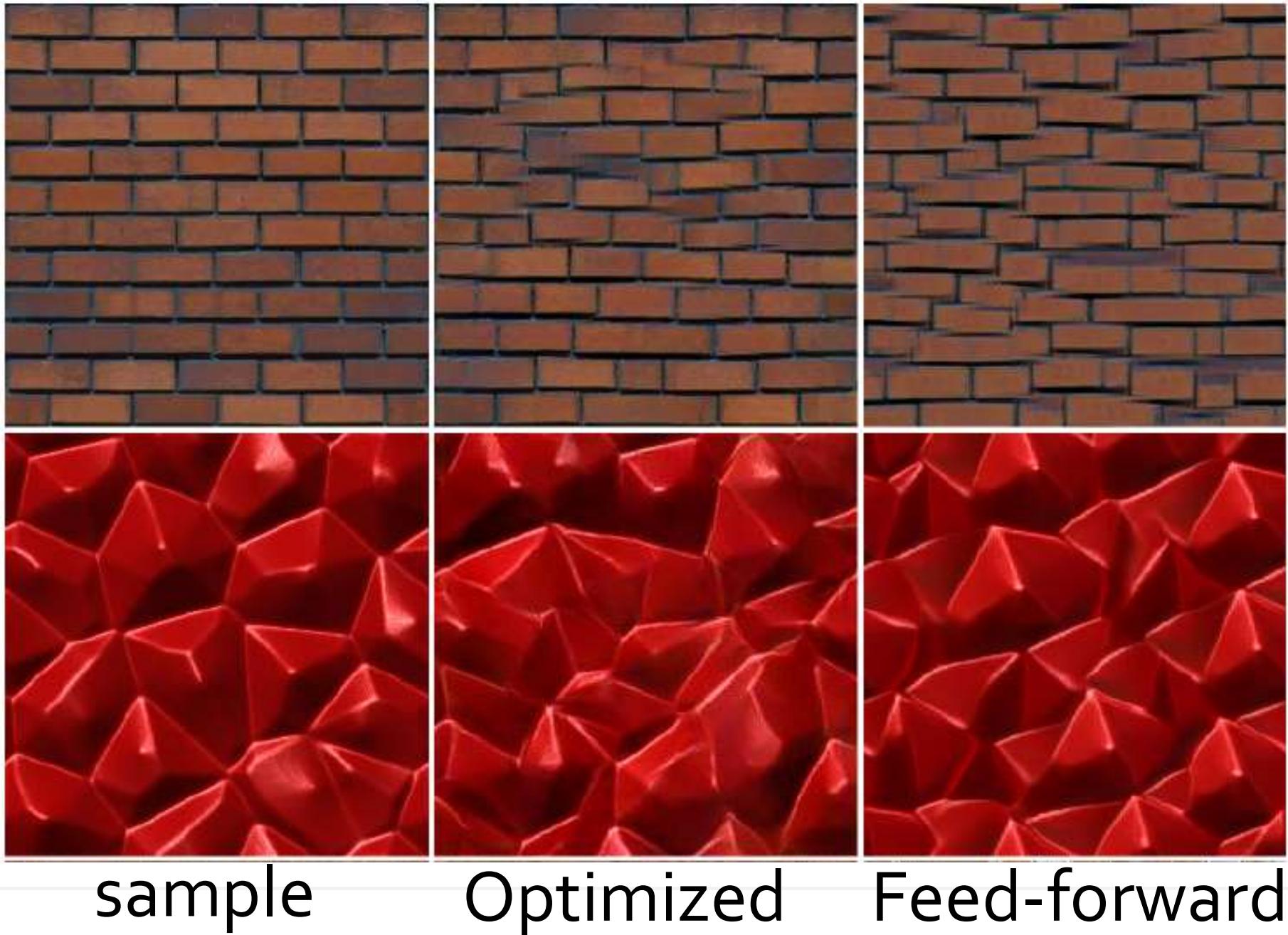


sample

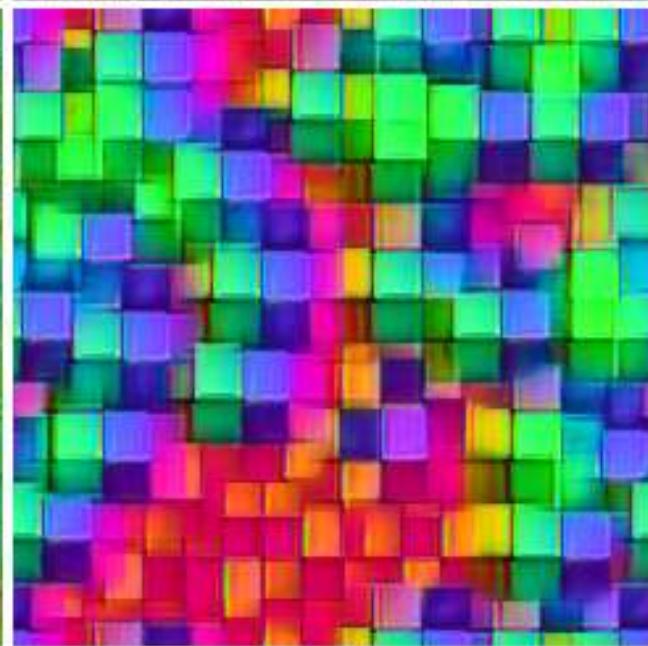
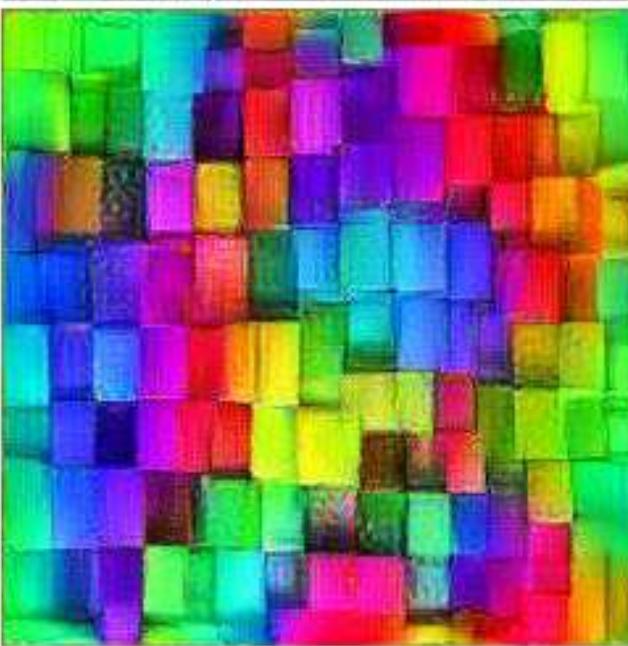
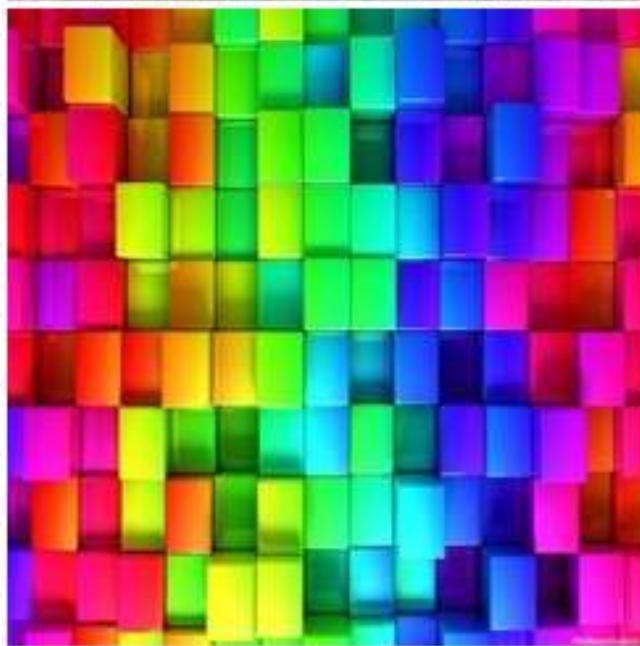
Optimized

Feed-forward

# Pre-image vs. feed-forward



# Pre-image vs. feed-forward

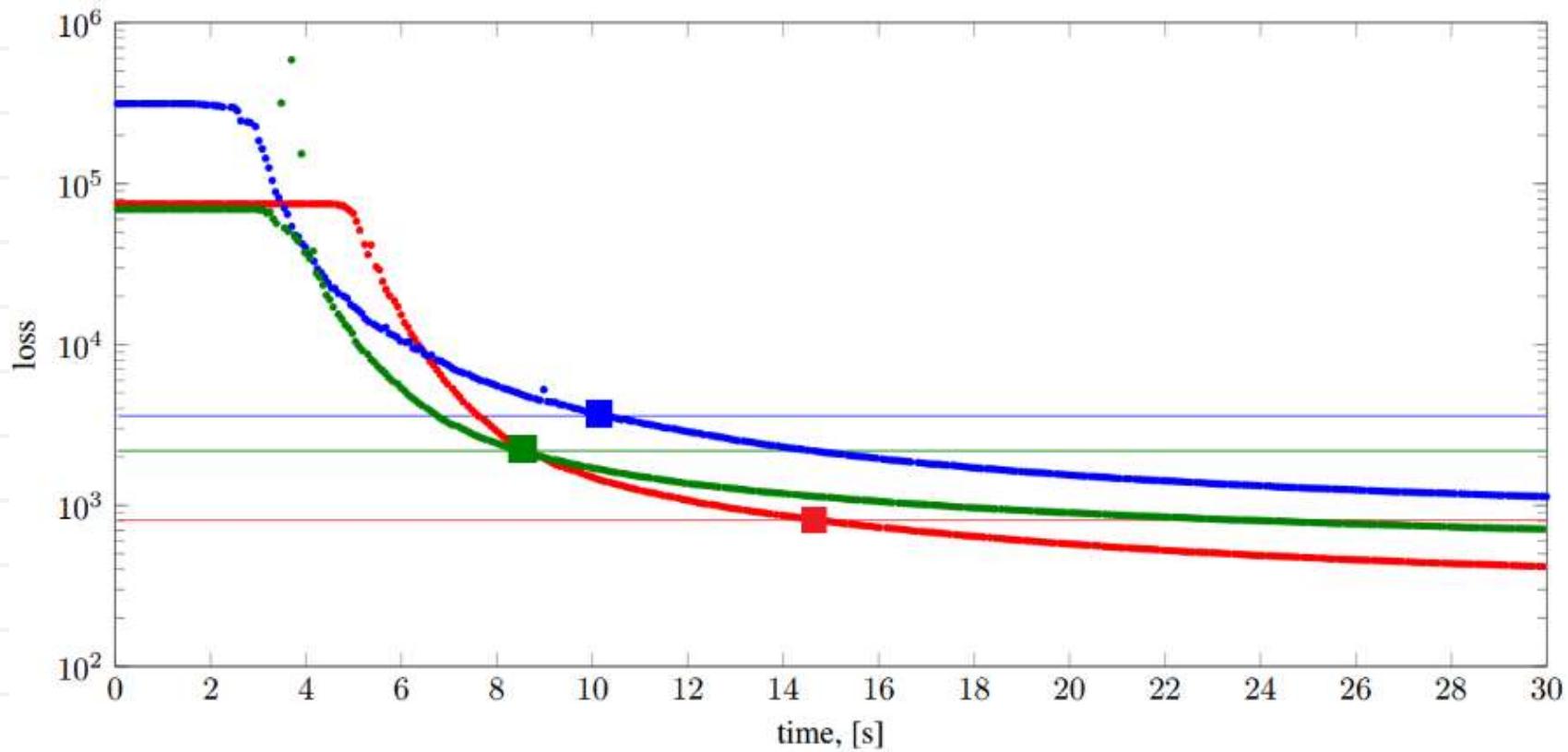


sample

Optimized

Feed-forward

# Quantitative evaluation



- Horizontal lines: average sample loss
  - **0.06 second** to generate sample
- Dotted lines: *optimization-based* loss as a function of time
  - **10 seconds** to achieve the same loss values

# “Style” network

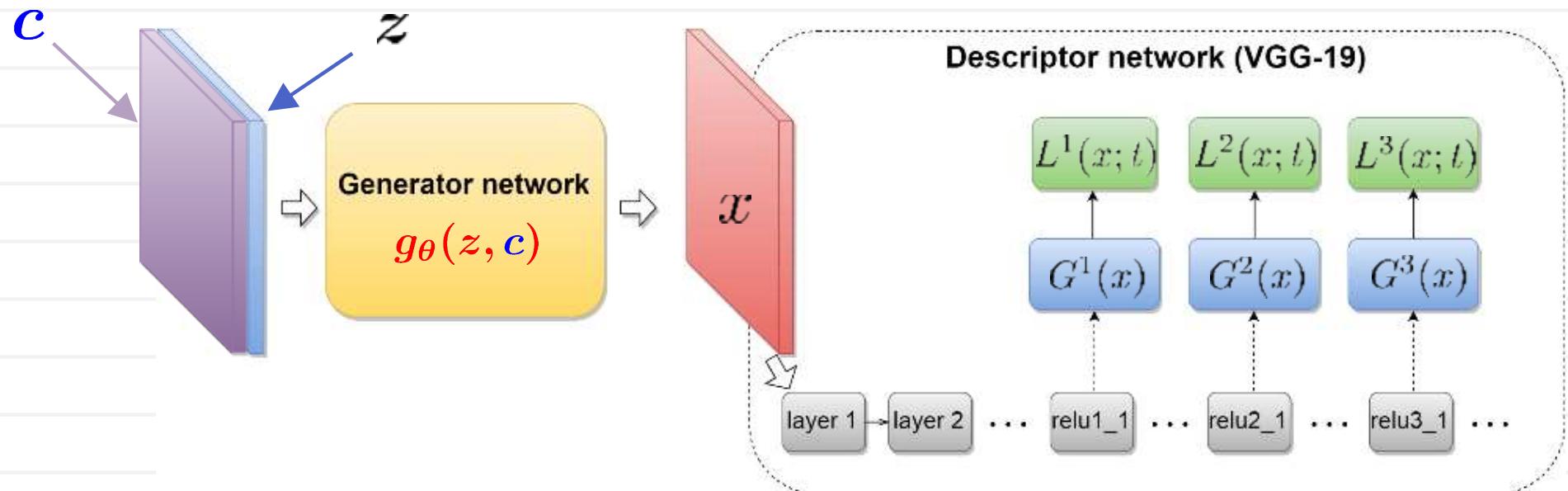


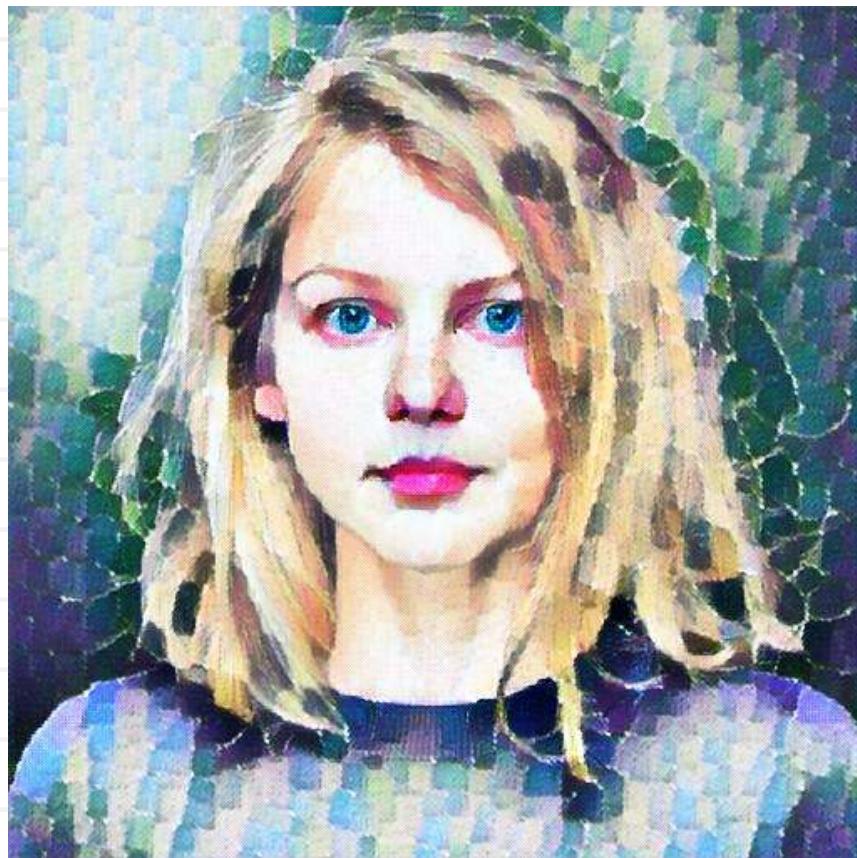
Image generation:

$$x = g_{\theta}(z, \mathbf{c}), \quad z \sim U(0, 1)$$

Optimization task:

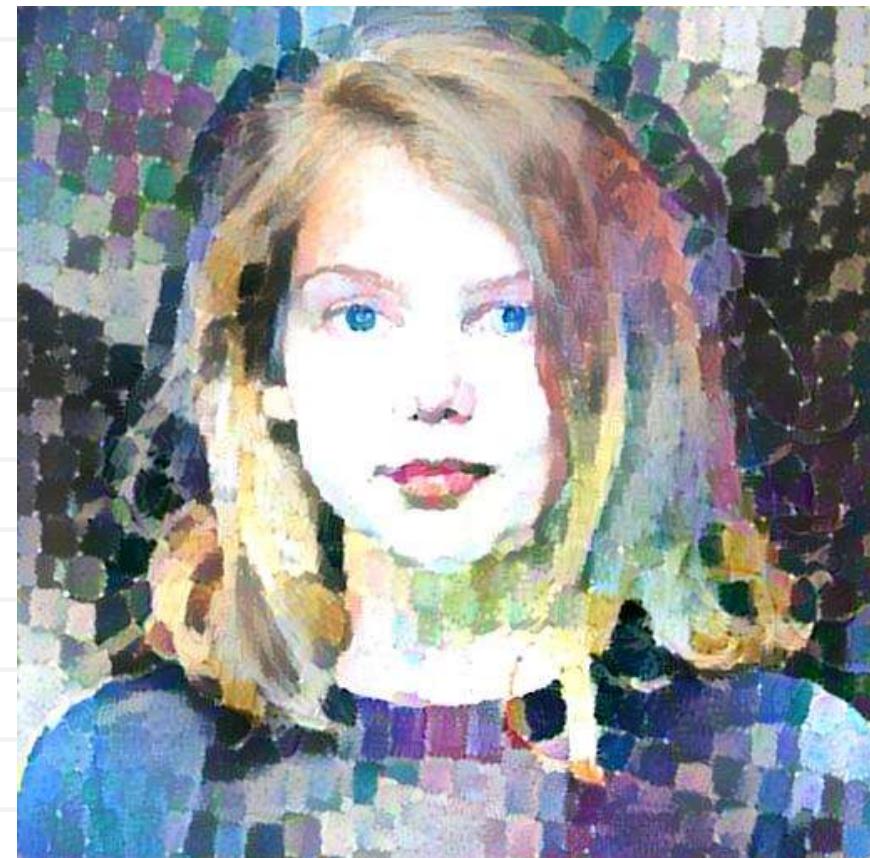
$$\min_{\theta} \mathbb{E} \mathcal{L}(g_{\theta}(z, \mathbf{c}); c, t), \quad z \sim U(0, 1)$$

# Feed-forward stylization

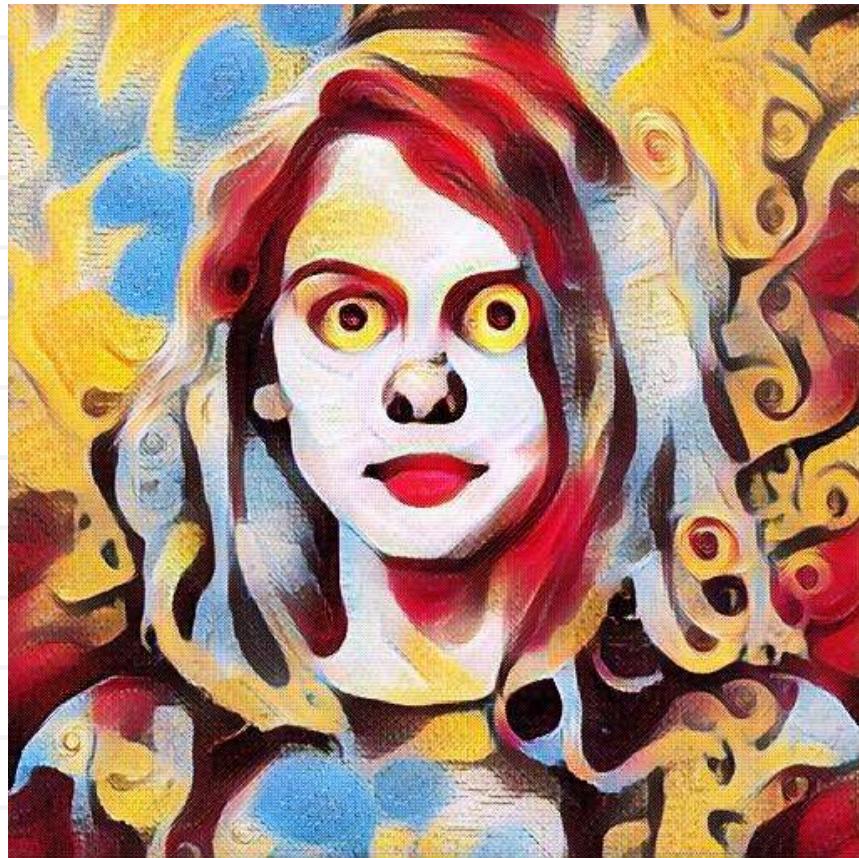


feedforward

optimization



# Feed-forward stylization

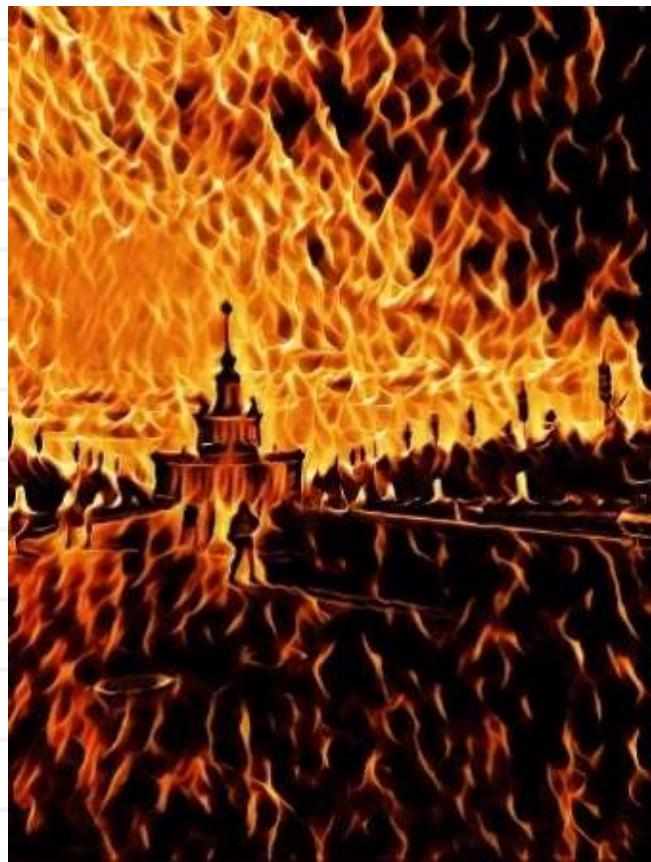
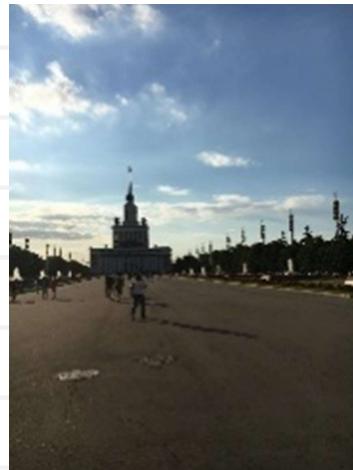


feedforward

optimization

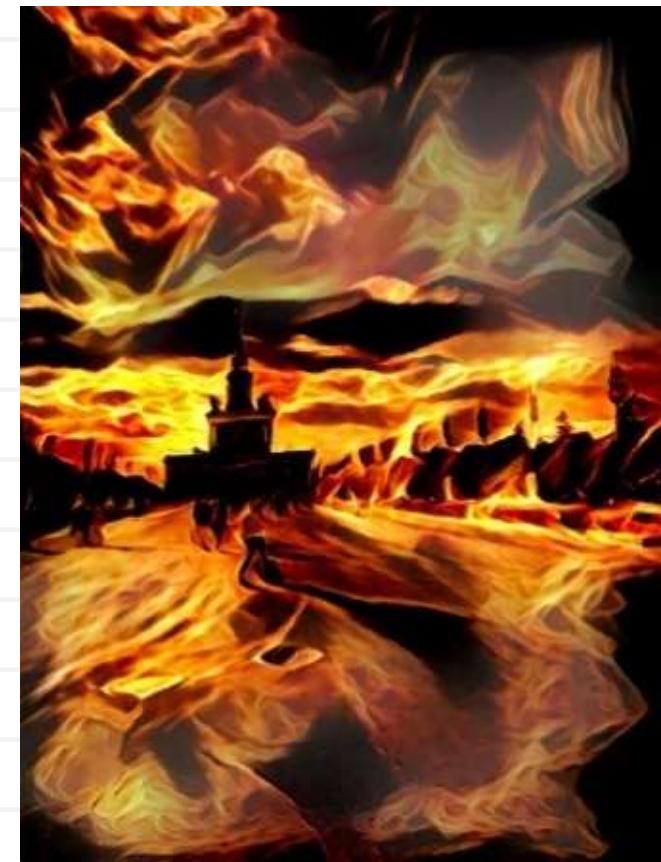


# Feed-forward stylization

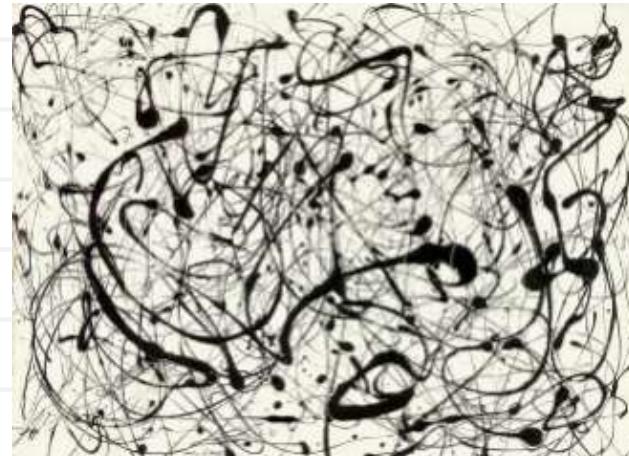


feedforward

optimization



# Feed-forward stylization



feedforward



optimization

# Recap

- *Synthesizing* realistic images is hard
- Better loss functions are key
- With good loss functions, we can train feedforward generators
- Big promise/hype about GANs, some fundamental challenges remain

# Bibliography

Alexey Dosovitskiy, Jost Tobias Springenberg, Thomas Brox:  
Learning to generate chairs with convolutional neural networks. CVPR 2015:  
**1538-1546**

Alexey Dosovitskiy, Thomas Brox:  
Inverting Convolutional Networks with Convolutional Networks. CoRR  
abs/1506.02753 (2015)

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, Yoshua Bengio:  
Generative Adversarial Nets. NIPS 2014

Peter J. Burt, Edward H. Adelson:  
A Multiresolution Spline with Application to Image Mosaics. ACM Trans. Graph.  
2(4): 217-236 (1983)

Emily L. Denton, Soumith Chintala, Arthur Szlam, Robert Fergus:  
Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks. ICLR 2016

# Bibliography

Burt, Peter J., and Edward H. Adelson. "The Laplacian pyramid as a compact image code." *Communications, IEEE Transactions on* 31.4 (1983): 532-540.

J Kim, JK Lee, KM Lee, Accurate Image Super-Resolution Using Very Deep Convolutional Networks, CVPR 2016

Leon A. Gatys, Alexander S. Ecker, Matthias Bethge:  
Texture Synthesis Using Convolutional Neural Networks. NIPS 2015: 262-270

Leon A. Gatys, Alexander S. Ecker, Matthias Bethge:  
A Neural Algorithm of Artistic Style. CoRR abs/1508.06576 (2015)

Dmitry Ulyanov, Vadim Lebedev, Andrea Vedaldi, Victor S. Lempitsky:  
Texture Networks: Feed-forward Synthesis of Textures and Stylized Images.  
CoRR abs/1603.03417 (2016)