

# Занятие 8

## Распознавание людей, Wasserstein GAN

Дмитрий Яшунин, к.ф.-м.н  
IntelliVision

# Организационные моменты

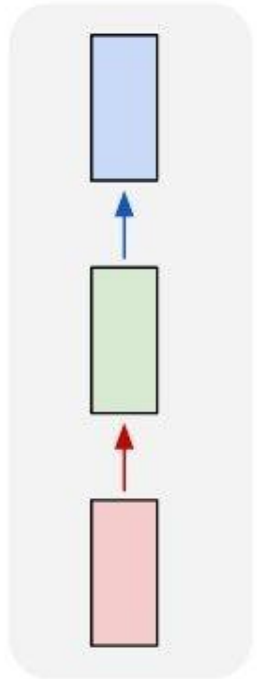
Сессия:

на следующей неделе

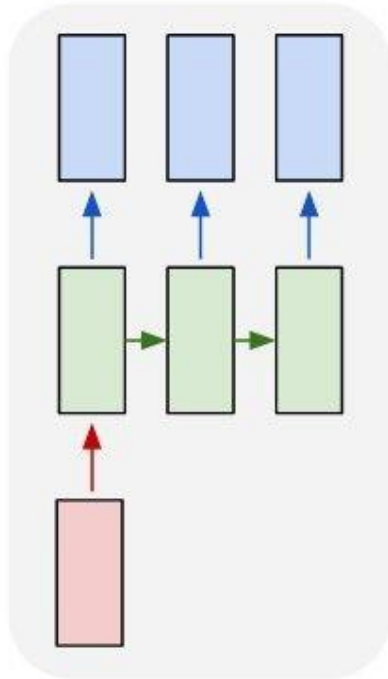
27 декабря, в этой же аудитории, с 18:00 до 21:00

# В прошлый раз: Recurrent Neural Networks

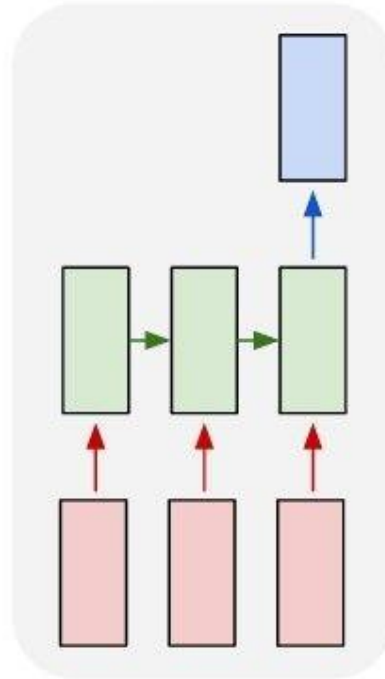
one to one



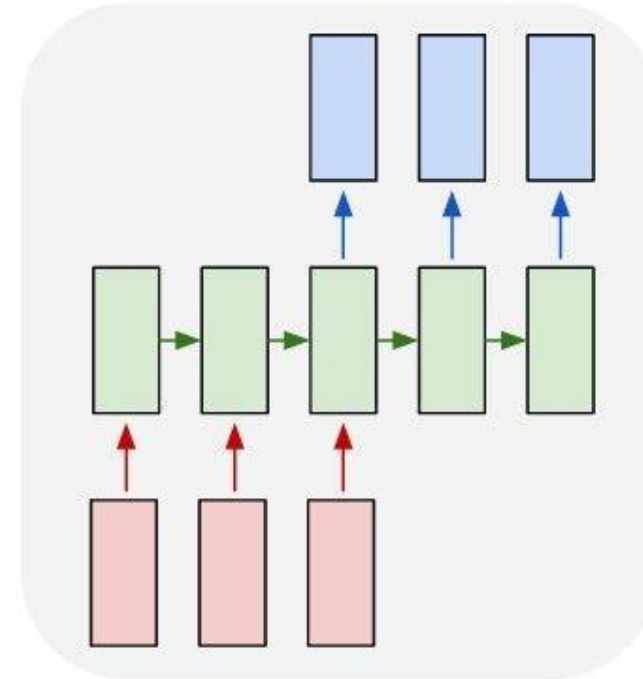
one to many



many to one



many to many



many to many

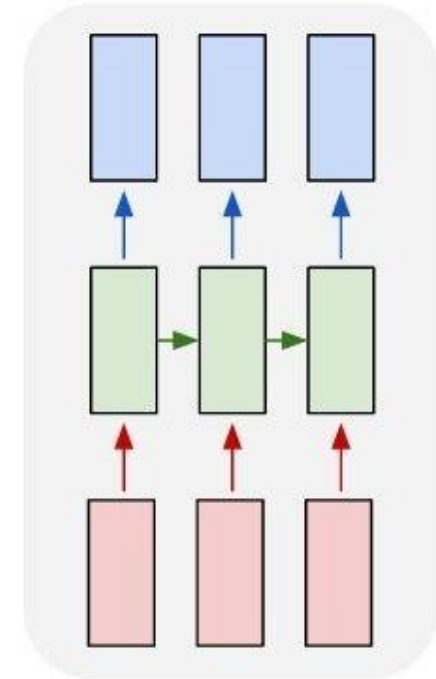


Image  
Classification

# В прошлый раз: Recurrent Neural Networks

one to one

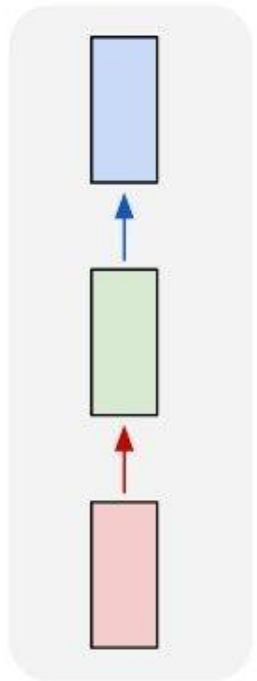


Image  
Classification

one to many

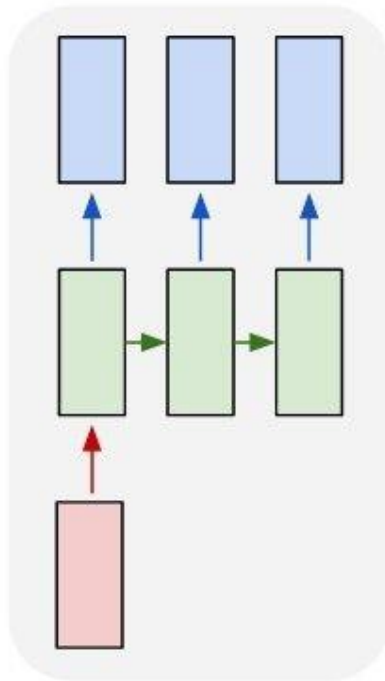
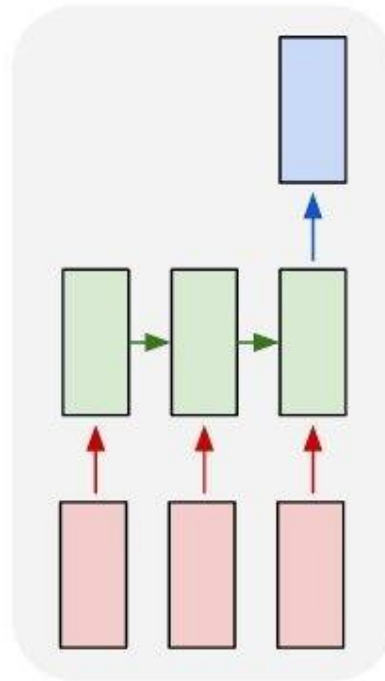
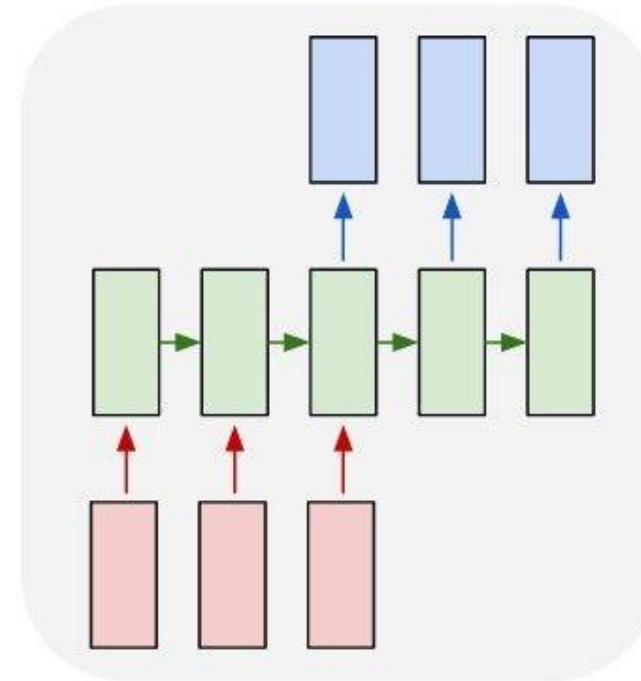


Image  
Captioning

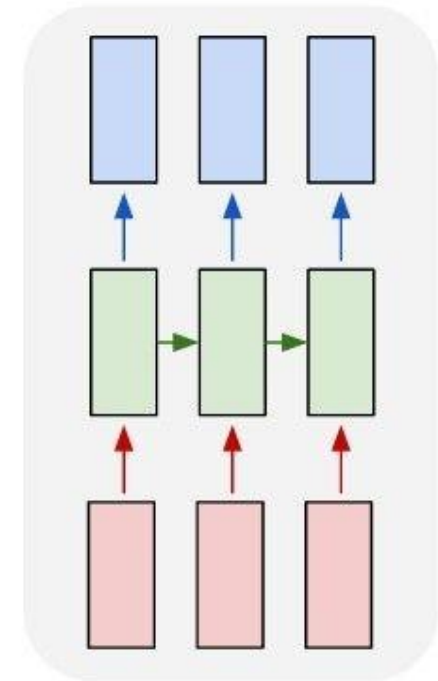
many to one



many to many

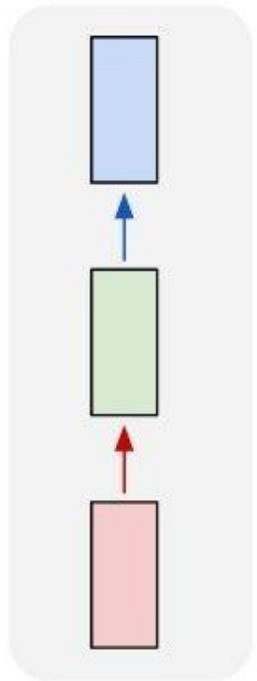


many to many



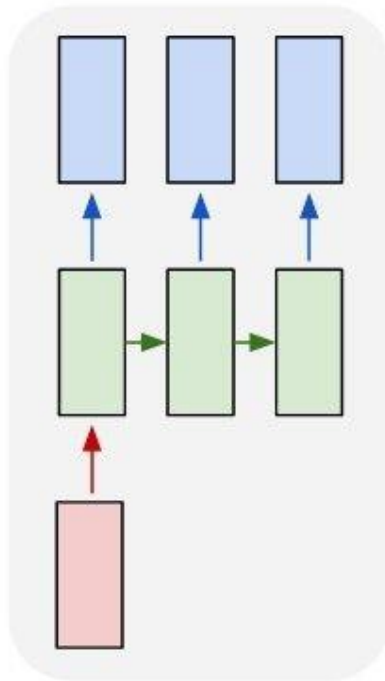
# В прошлый раз: Recurrent Neural Networks

one to one



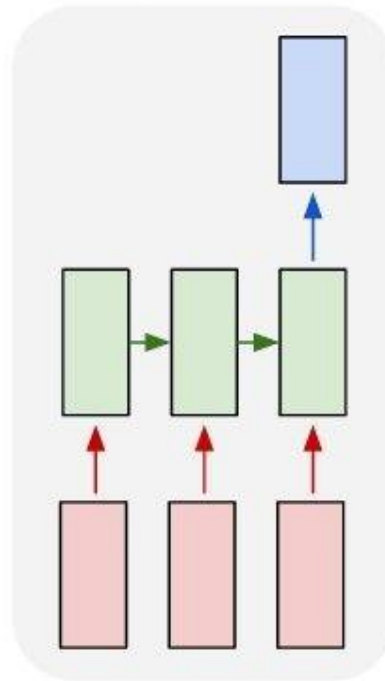
**Image  
Classification**

one to many



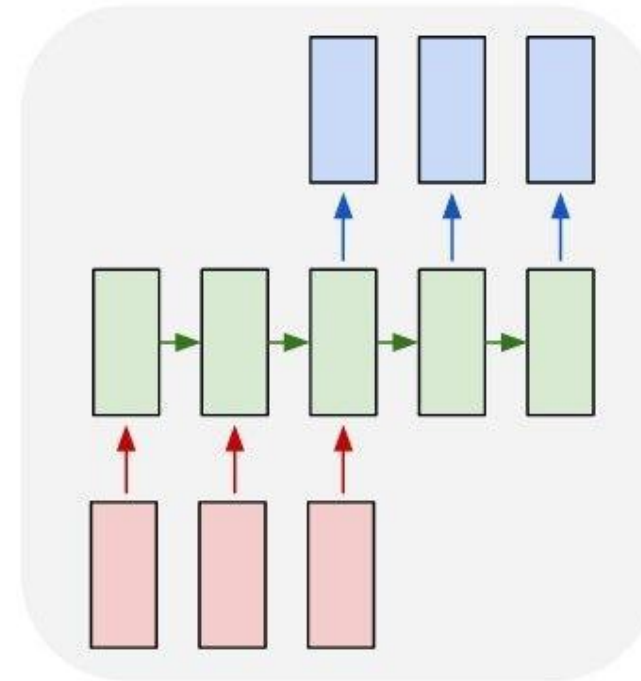
**Image  
Captioning**

many to one

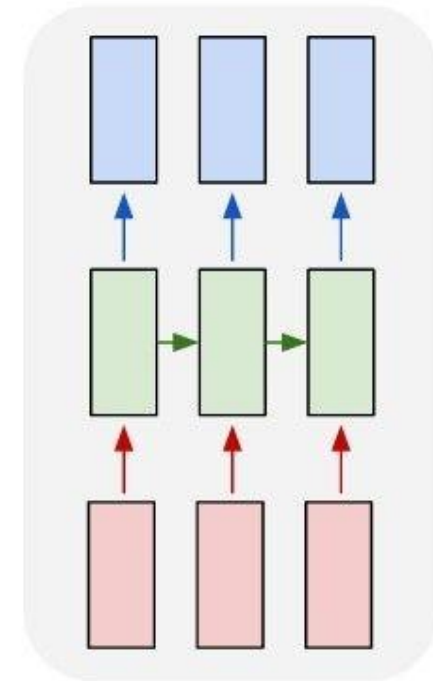


**Person reidentification by  
sequence of images**

many to many

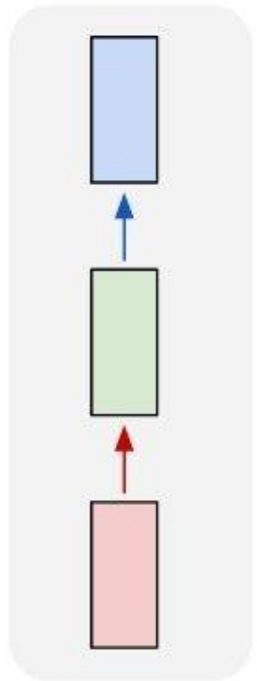


many to many



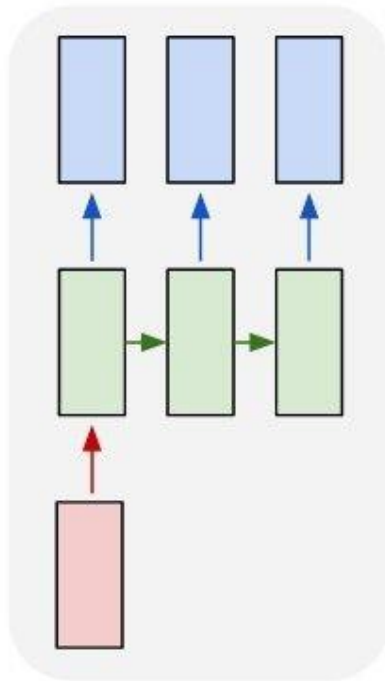
# В прошлый раз: Recurrent Neural Networks

one to one



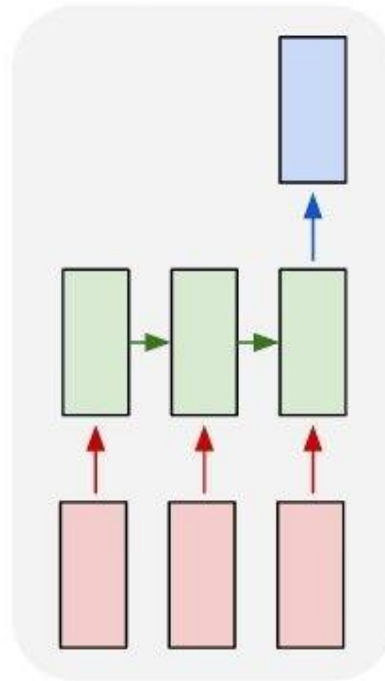
**Image  
Classification**

one to many



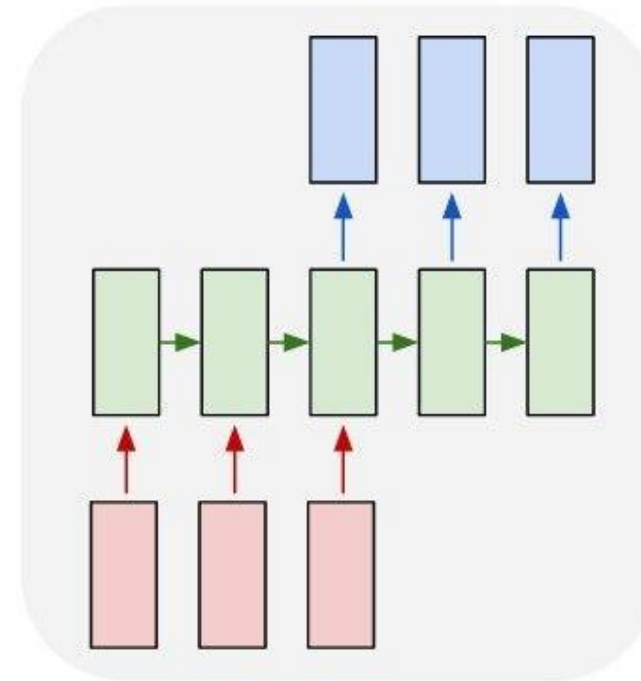
**Image  
Captioning**

many to one



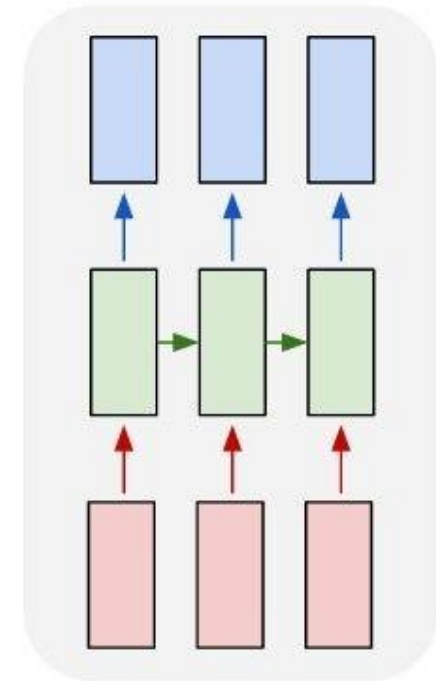
**Person  
reidentification**

many to many



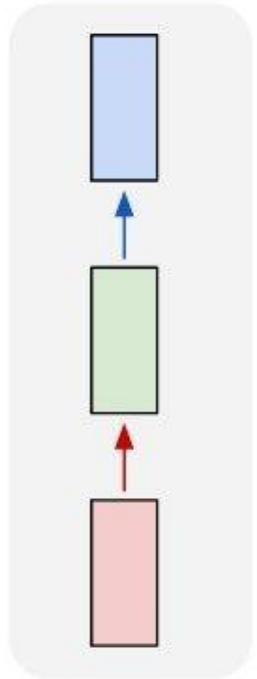
**Machine Translation**

many to many



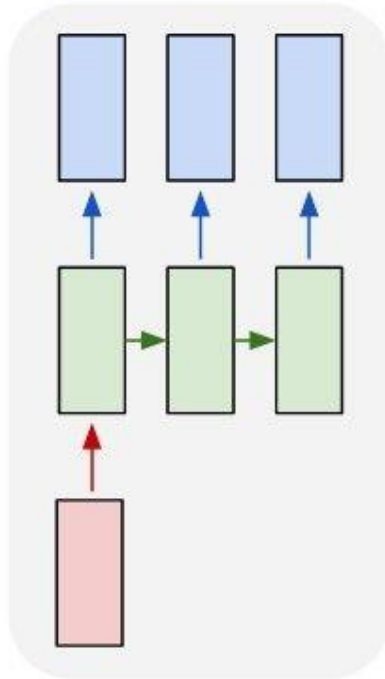
# В прошлый раз: Recurrent Neural Networks

one to one



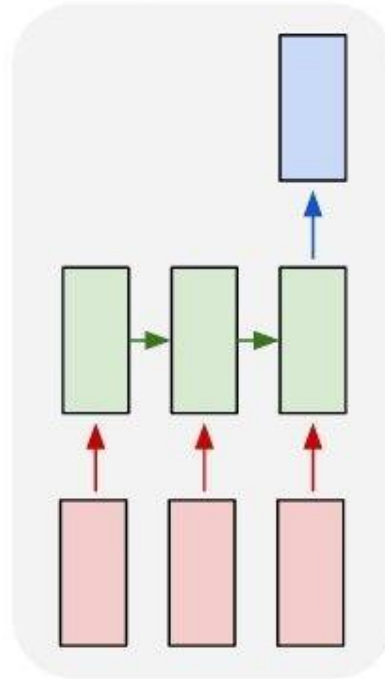
**Image  
Classification**

one to many



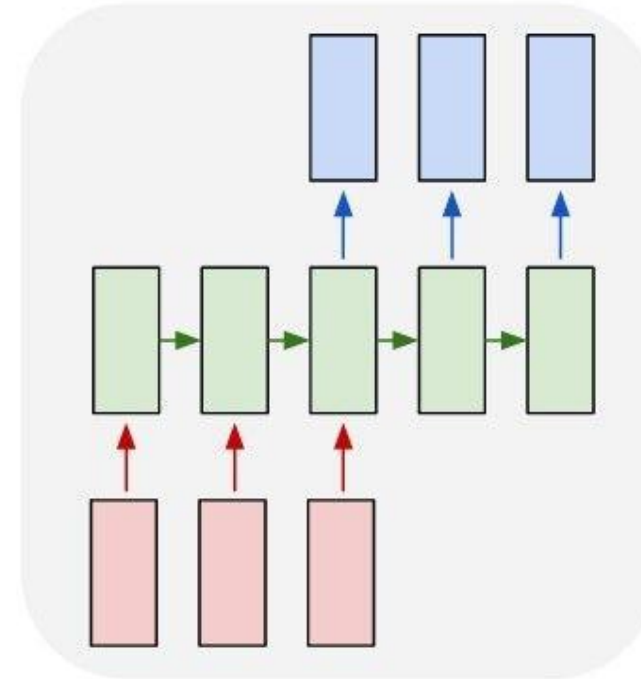
**Image  
Captioning**

many to one



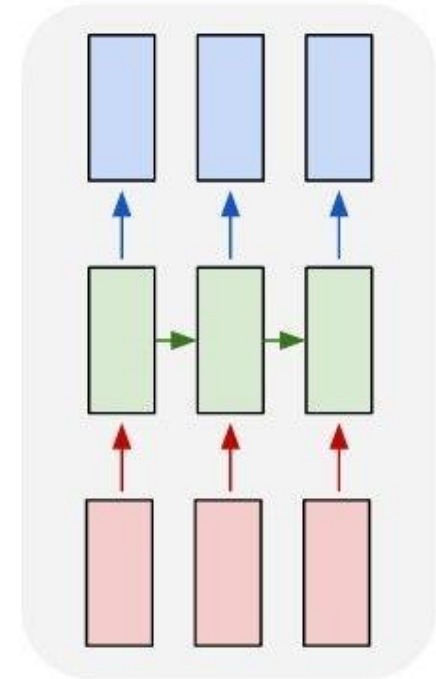
**Person  
reidentification**

many to many



**Machine Translation**

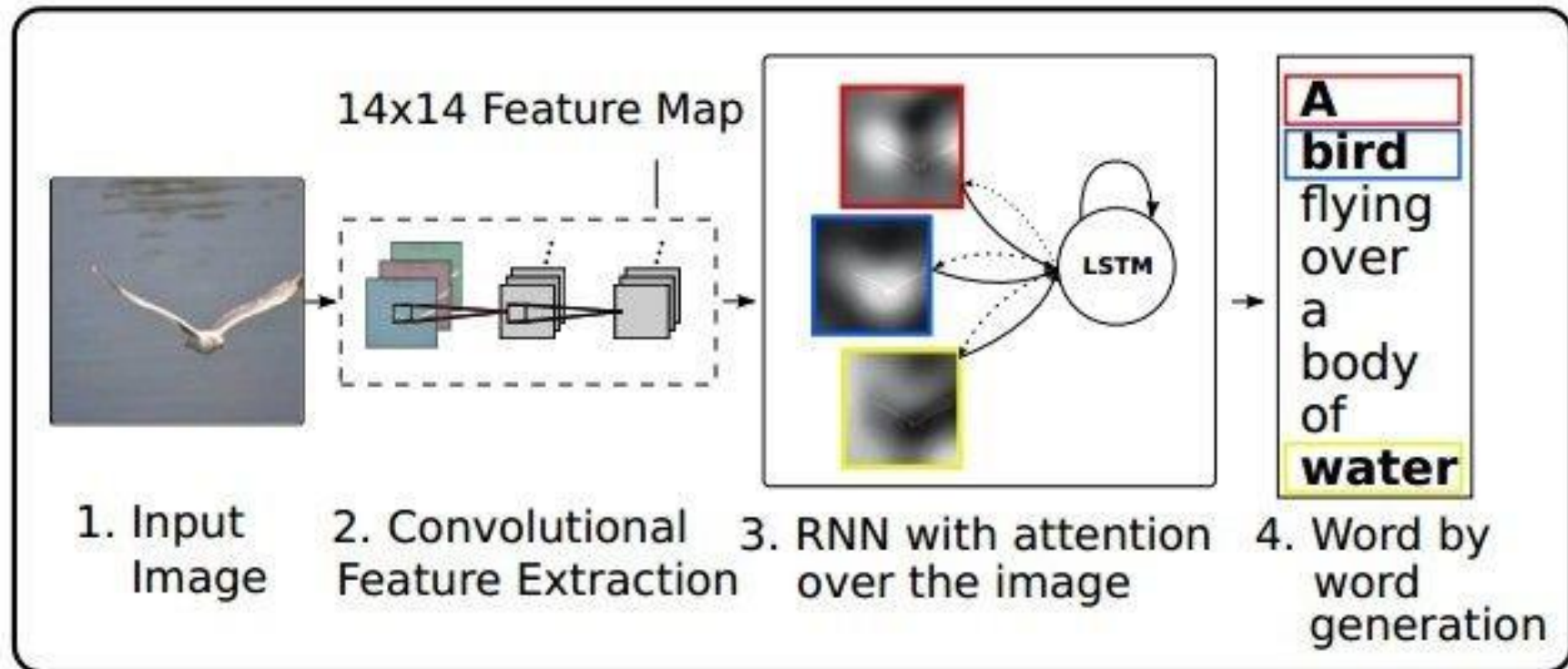
many to many



**Video  
classification on  
frame level**

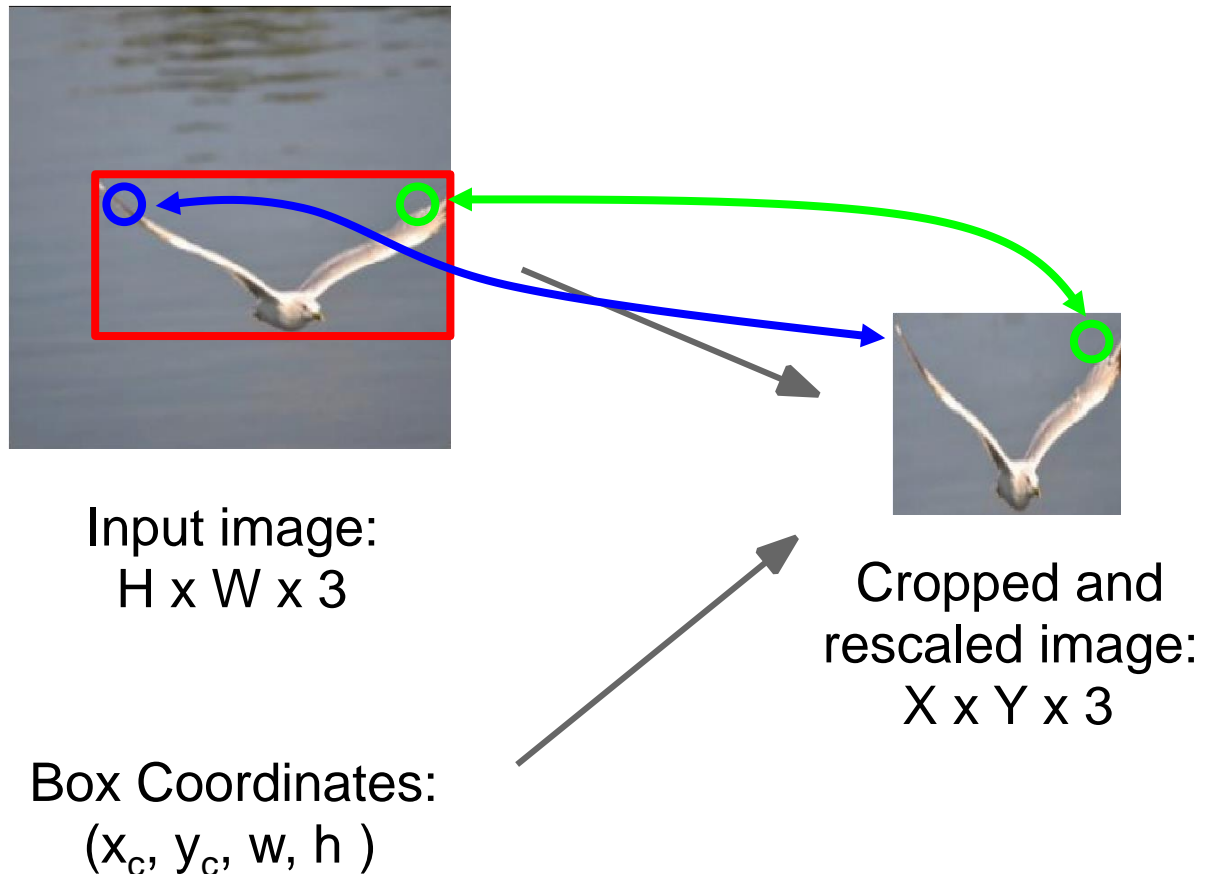
# В прошлый раз: Image Captioning with Attention

RNN focuses its attention at a different spatial location when generating each word



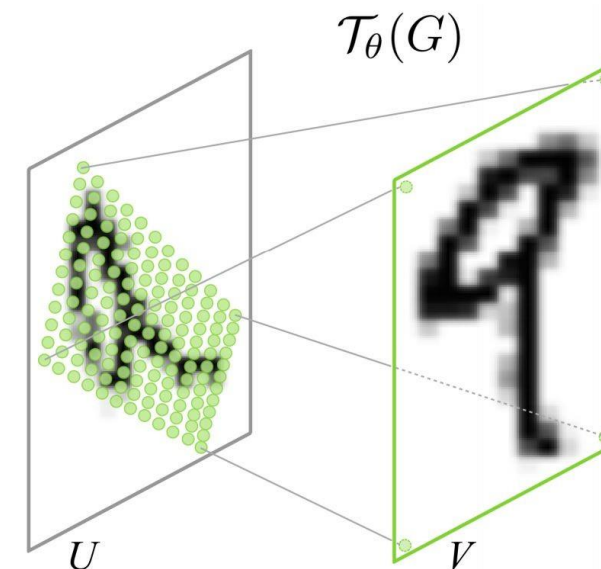


# В прошлый раз: Spatial Transformer Networks



**Idea:** Function mapping  
*pixel coordinates*  $(x^t, y^t)$  of  
output to *pixel coordinates*  
 $(x^s, y^s)$  of input

$$\begin{pmatrix} x_i^s \\ y_i^s \end{pmatrix} = \begin{bmatrix} \theta_{11} & \theta_{12} & \theta_{13} \\ \theta_{21} & \theta_{22} & \theta_{23} \end{bmatrix} \begin{pmatrix} x_i^t \\ y_i^t \\ 1 \end{pmatrix}$$

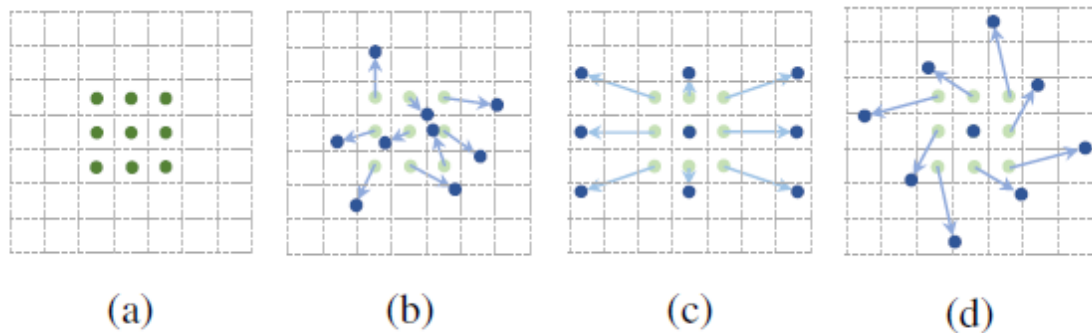


Repeat for all pixels  
in *output* to get a  
**sampling grid**

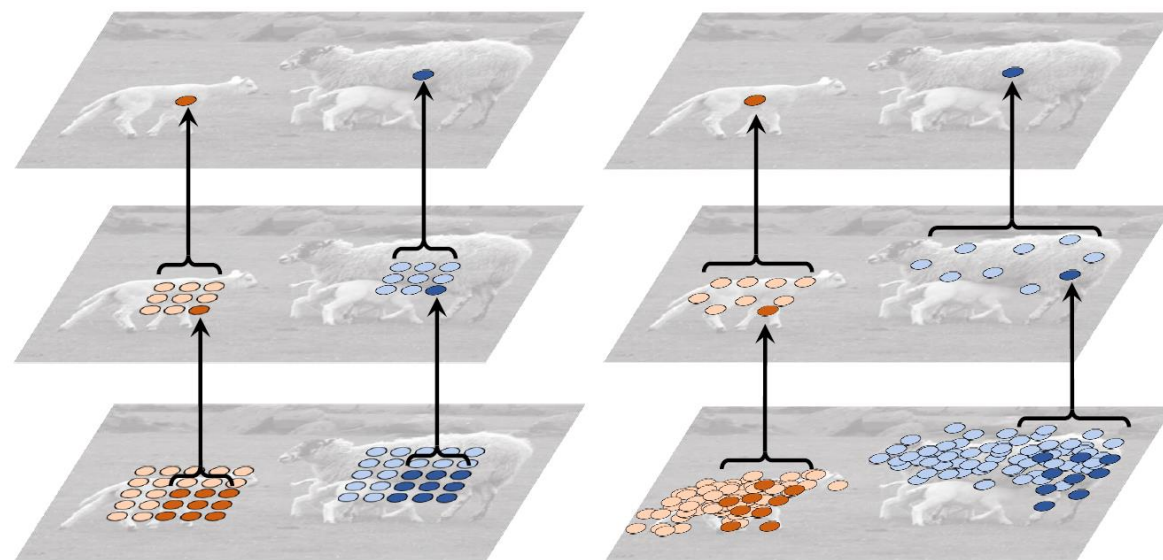
Then use **bilinear  
interpolation** to  
compute output

# В прошлый раз: Deformable Convolutions

Dynamic & learnable receptive field



- (a) regular sampling grid of standard convolution
- (b) deformed sampling locations with augmented offsets in deformable convolution
- (c)(d) show that the deformable convolution generalizes various transformations for scale, (anisotropic) aspect ratio and rotation



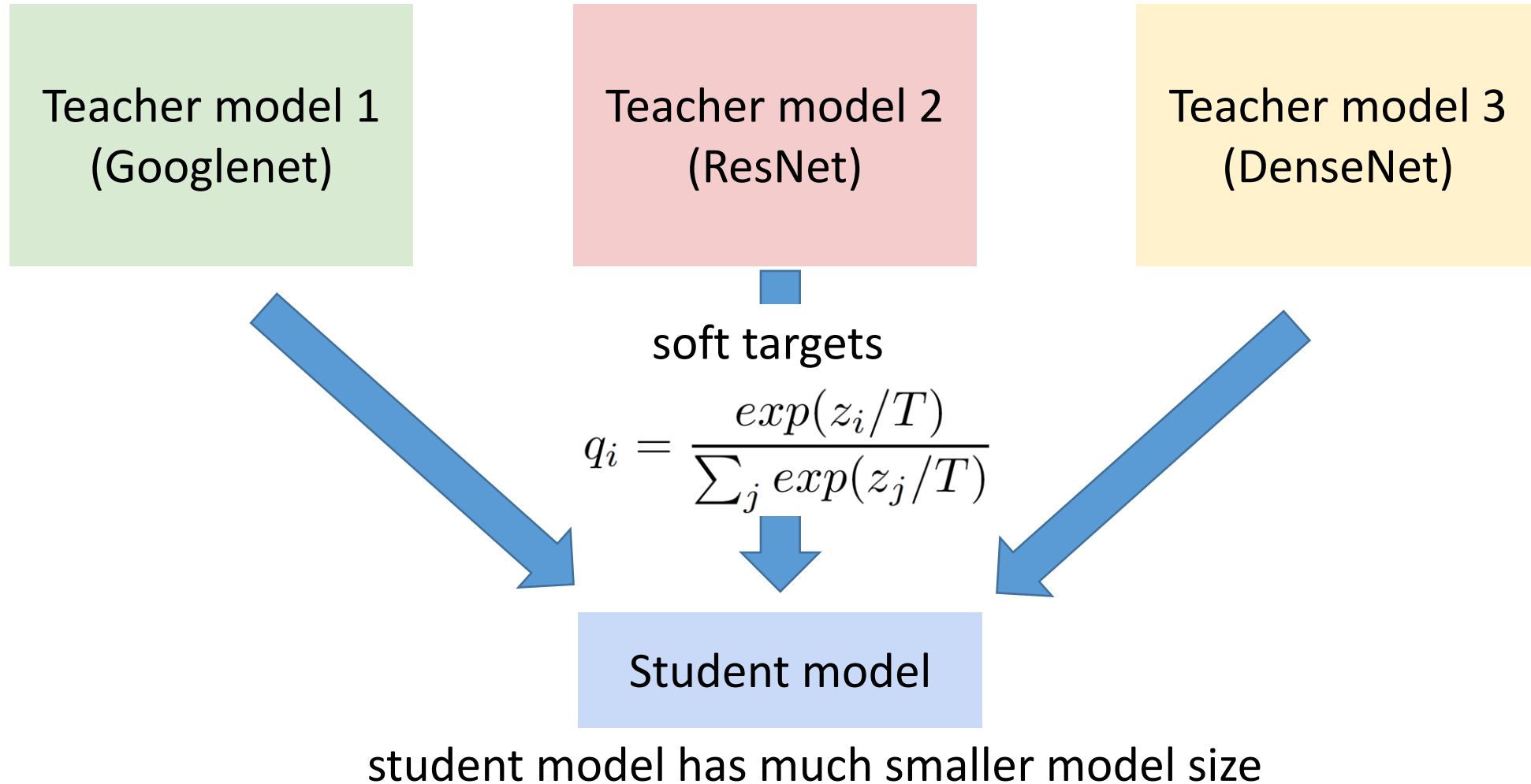
(a) standard convolution

(b) deformable convolution

- (a) fixed receptive field in standard convolution
- (b) adaptive receptive field in deformable convolution

# Network Distillation

Transfer knowledge from an ensemble of models to a small single model

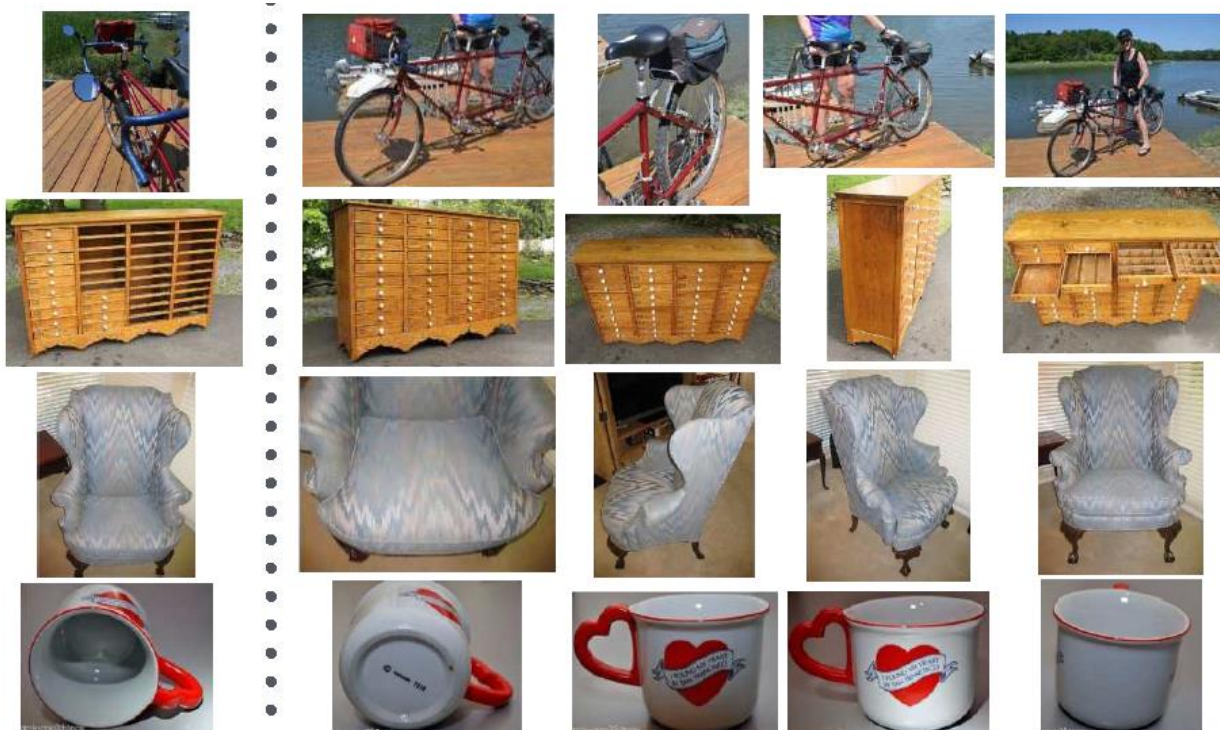


Сегодня:  
Распознавание людей,  
Wasserstein GAN

# Image similarity: Image Retrieval

Query

Retrieval



Example retrieval results on Stanford Online Products



# Image similarity: Face recognition



Example face images in MegaFace dataset

Wen et al, "A Discriminative Feature Learning Approach for Deep Face Recognition", 2016

# Image similarity: Person Reidentification



# Image similarity

Image similarity

Image retrieval

Zero/one shot learning

Face recognition

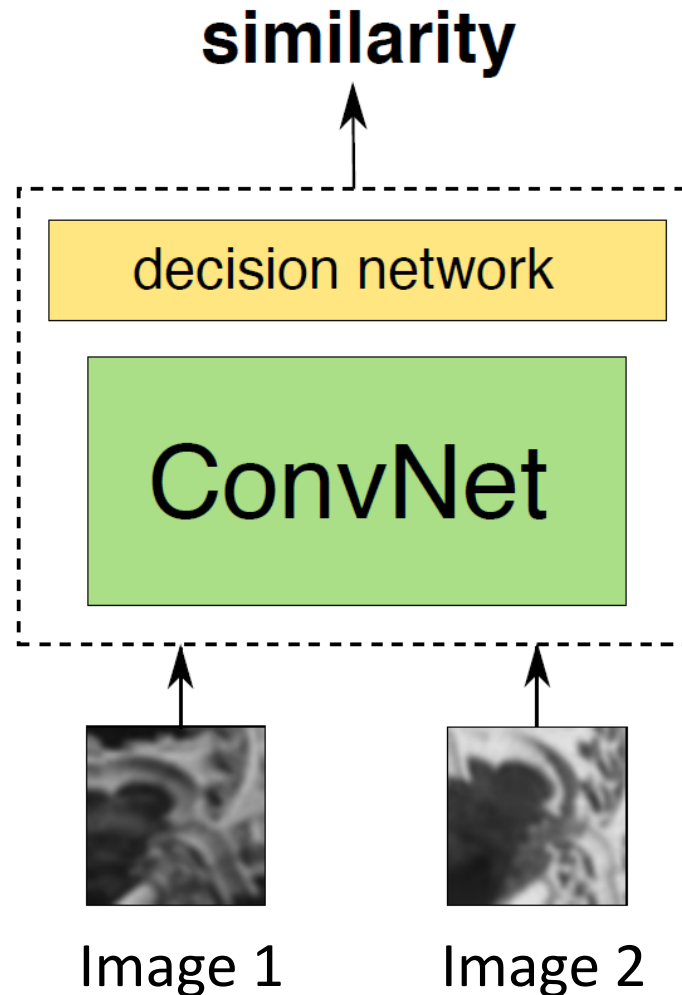
Person reidentification (ReID)



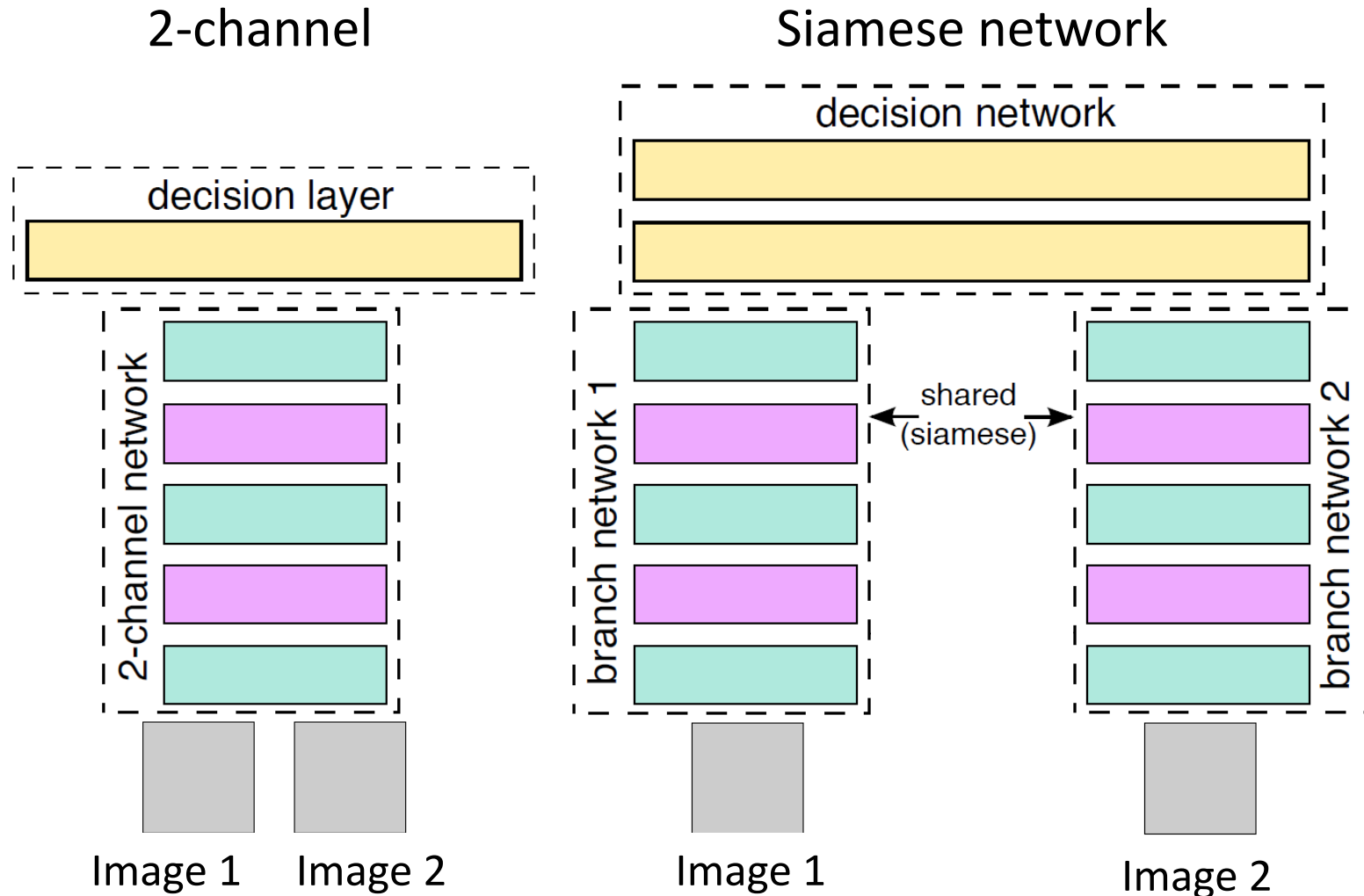
The same task:  
find “similar” images



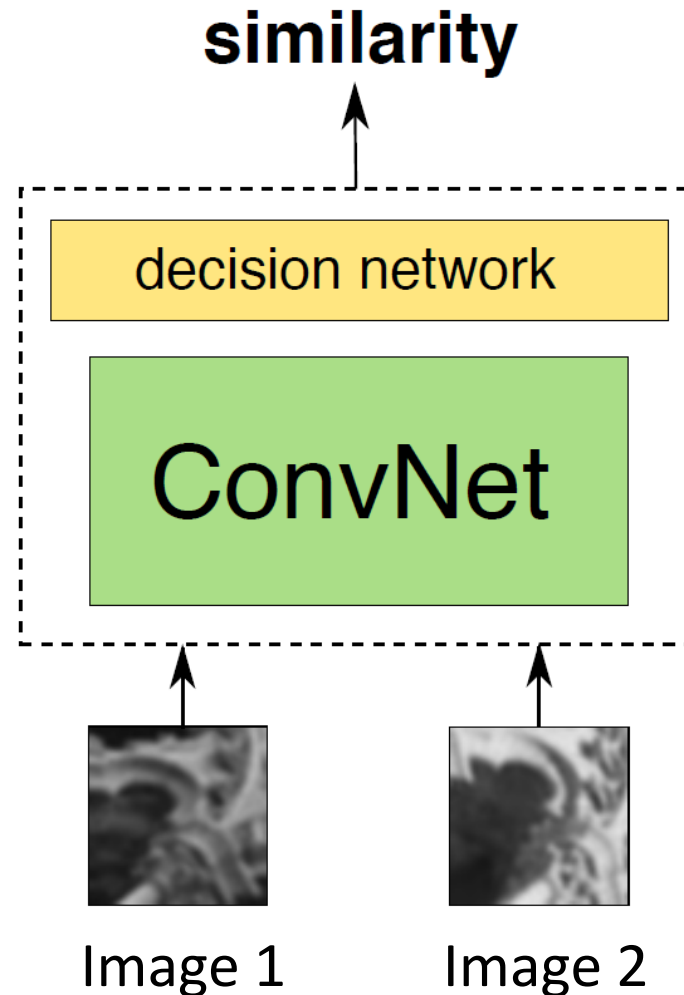
# Image similarity: Learning Similarity Function



# Image similarity: Learning Similarity Function



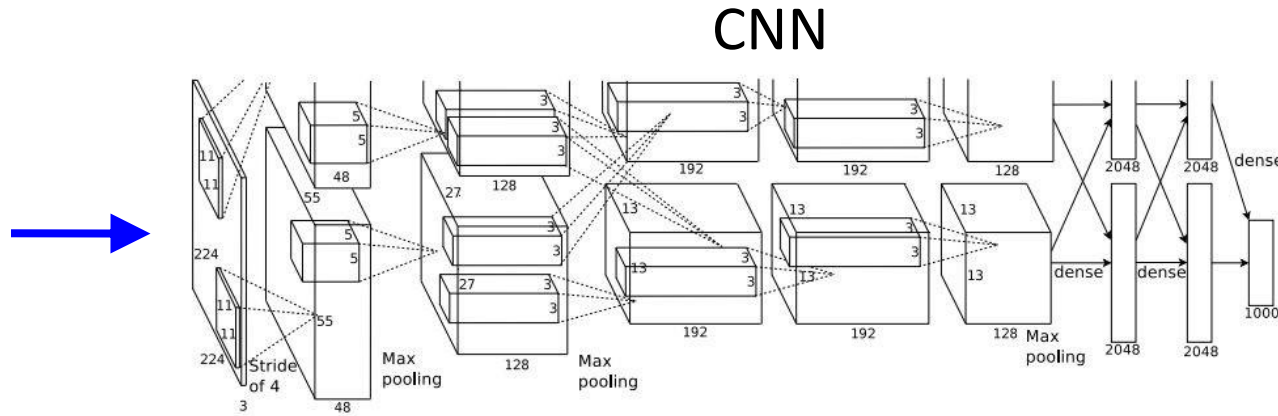
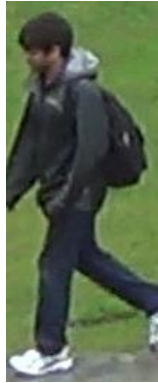
# Image similarity: Learning Similarity Function



Very slow  
for many  
pairs

# Image Similarity: Embeddings

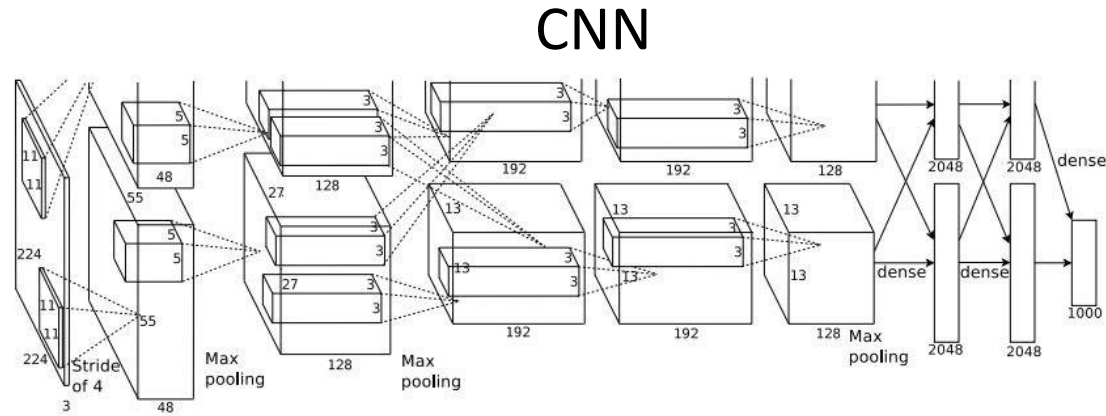
Image



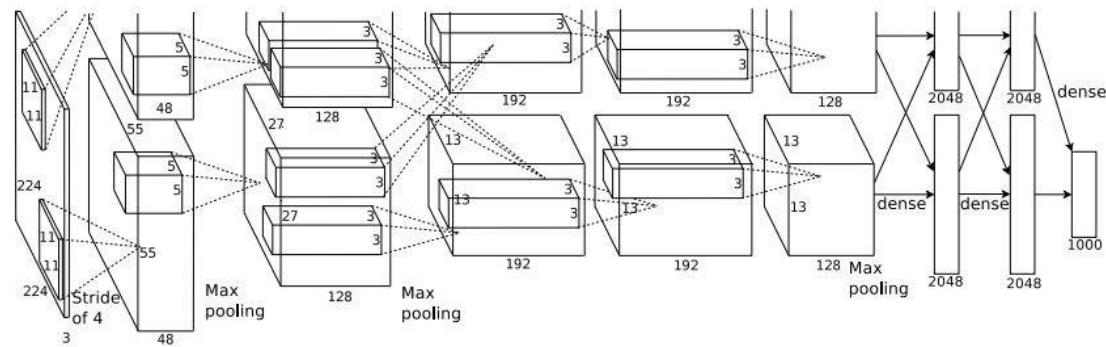
Embedding  
vector of fixed size  
typical  $2^n$ : 128, 256, 512

# Image Similarity: Embeddings

Image



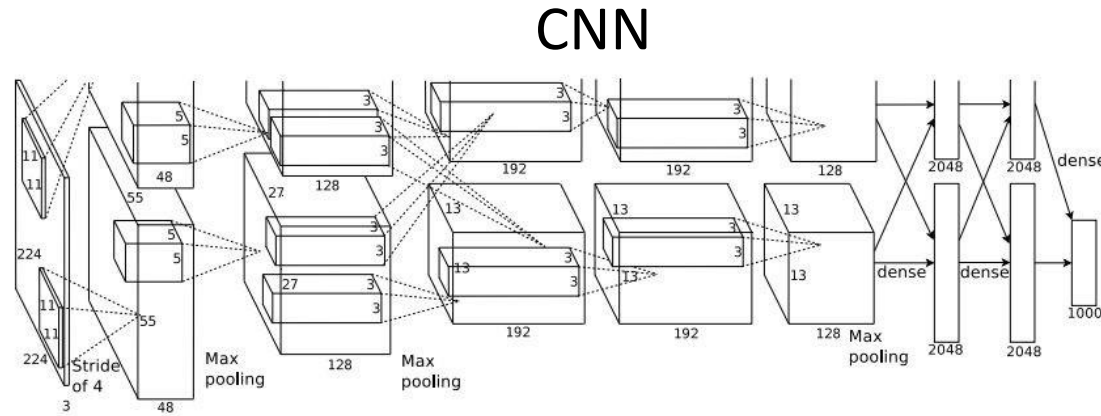
## Embedding 1



## Embedding 2

# Image Similarity: Embeddings

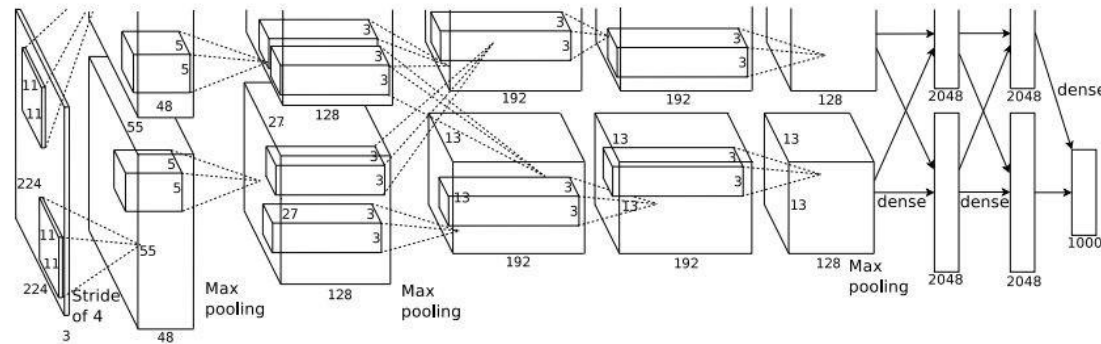
Image



Embedding 1



compare l2 distance  
between embeddings



Embedding 2





# Image Similarity: Loss functions: **Softmax**

Test image    L2 Nearest neighbors in feature space

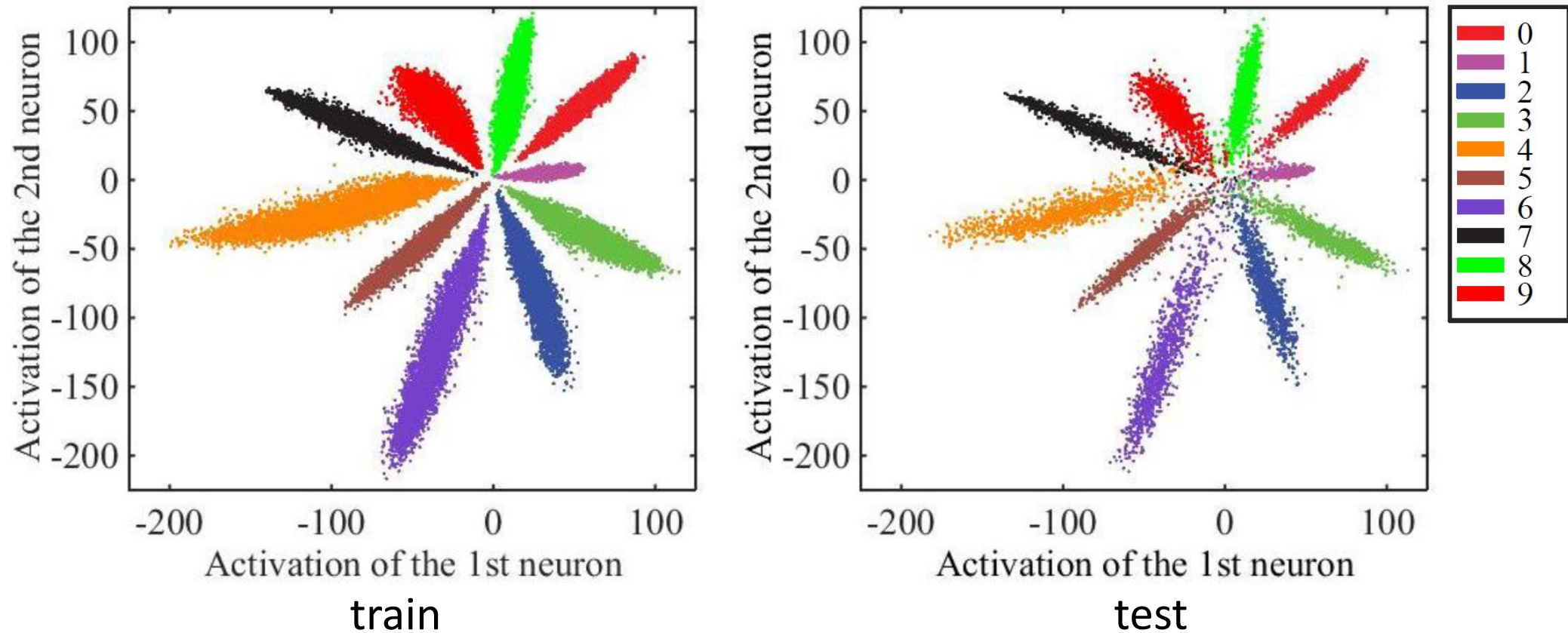
Solve classification task

Use output of last layer  
before softmax as  
embedding



# Image Similarity: Loss functions: Softmax

MNIST handwritten digits 0-9





# Image Similarity: Loss functions: **Center Loss**

$$\mathcal{L}_S = - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}} \quad \text{Softmax loss}$$

# Image Similarity: Loss functions: Center Loss

$$\mathcal{L}_S = - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}}$$

Softmax loss

$\mathbf{x}_i$  – output features (before softmax layer)

# Image Similarity: Loss functions: **Center Loss**

$$\mathcal{L}_S = - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}}$$

Softmax loss

$\mathbf{x}_i$  – output features (before softmax layer)

$$\mathcal{L}_C = \frac{1}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

Center loss: direct clustering of features around centers

  $\mathbf{c}_{y_i}$  – learnable centers for classes  $y_i$

# Image Similarity: Loss functions: **Center Loss**

$$\mathcal{L}_S = - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}}$$

Softmax loss

$\mathbf{x}_i$  – output features (before softmax layer)

$$\mathcal{L}_C = \frac{1}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

Center loss: direct clustering of features around centers

$$\mathcal{L} = \mathcal{L}_S + \lambda \mathcal{L}_C$$

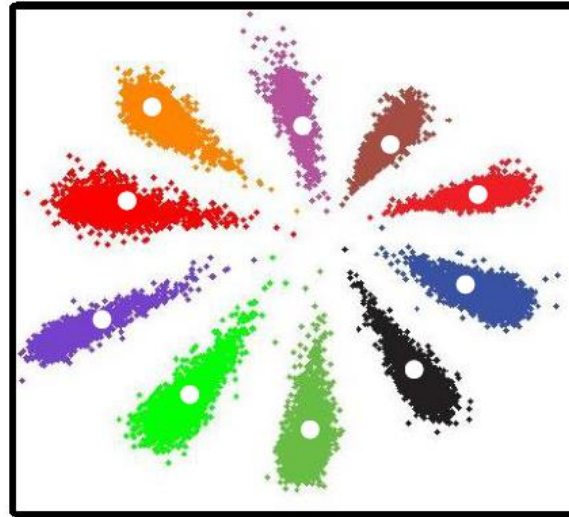
$$= - \sum_{i=1}^m \log \frac{e^{W_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^n e^{W_j^T \mathbf{x}_i + b_j}} + \frac{\lambda}{2} \sum_{i=1}^m \|\mathbf{x}_i - \mathbf{c}_{y_i}\|_2^2$$

Total loss

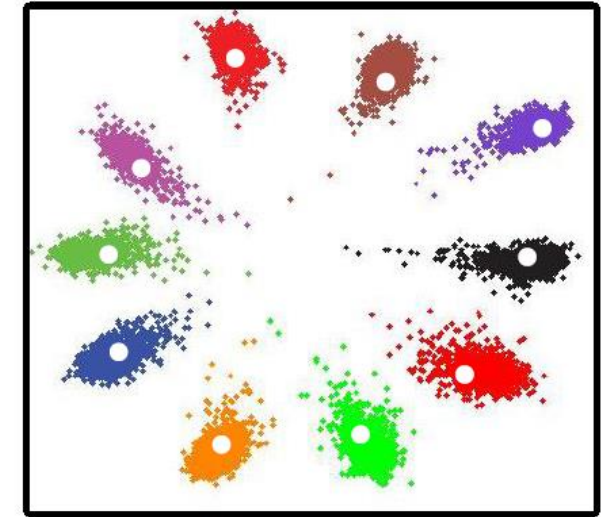
# Image Similarity: Loss functions: **Center Loss**

MNIST handwritten digits 0-9

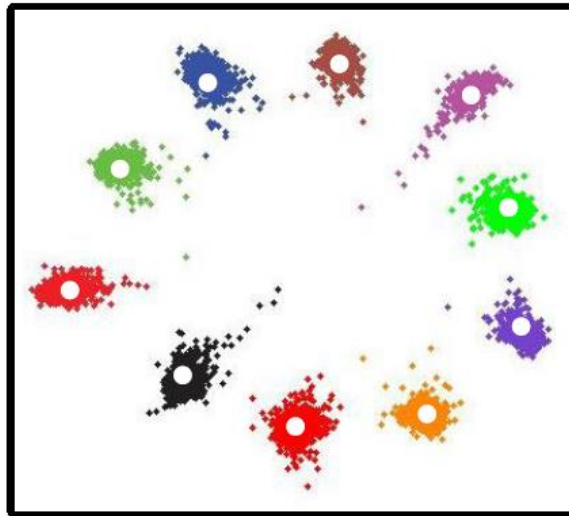
Feature clustering for softmax+center loss



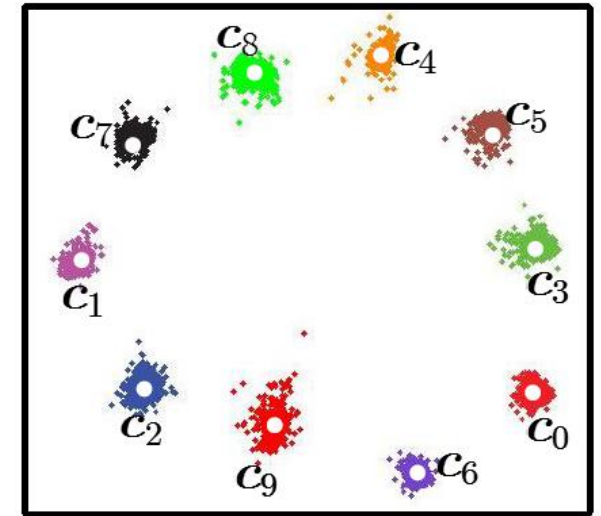
(a)  $\lambda = 0.001$



(b)  $\lambda = 0.01$



(c)  $\lambda = 0.1$



(d)  $\lambda = 1$

# Image Similarity: Loss functions: **Contrastive loss**

$$L_i(X_1, X_2) = (1 - Y)D_{X_1, X_2}^2 + Y\max(0, m - D_{X_1, X_2})^2$$

# Image Similarity: Loss functions: **Contrastive loss**

$X_1$  and  $X_2$  images to compare

$$L_i(\boxed{X_1, X_2}) = (1 - Y)D_{X_1, X_2}^2 + Y\max(0, m - \boxed{D_{X_1, X_2}})^2$$

$D$  – distance between  
image embeddings

# Image Similarity: Loss functions: **Contrastive loss**

$X_1$  and  $X_2$  images to compare

$$L_i(X_1, X_2) = (1 - Y)D_{X_1, X_2}^2 + Y\max(0, m - D_{X_1, X_2})^2$$

$Y=0$  if  $X_1$  and  $X_2$  are similar  
 $Y=1$  if  $X_1$  and  $X_2$  are dissimilar

margin

$D$  – distance between image embeddings



# Image Similarity: Loss functions: **Contrastive loss**

$X_1$  and  $X_2$  images to compare

margin

$$L_i(X_1, X_2) = (1 - Y)D_{X_1, X_2}^2 + Y\max(0, m - D_{X_1, X_2})^2$$

$Y=0$  if  $X_1$  and  $X_2$  are similar

$Y=1$  if  $X_1$  and  $X_2$  dissimilar

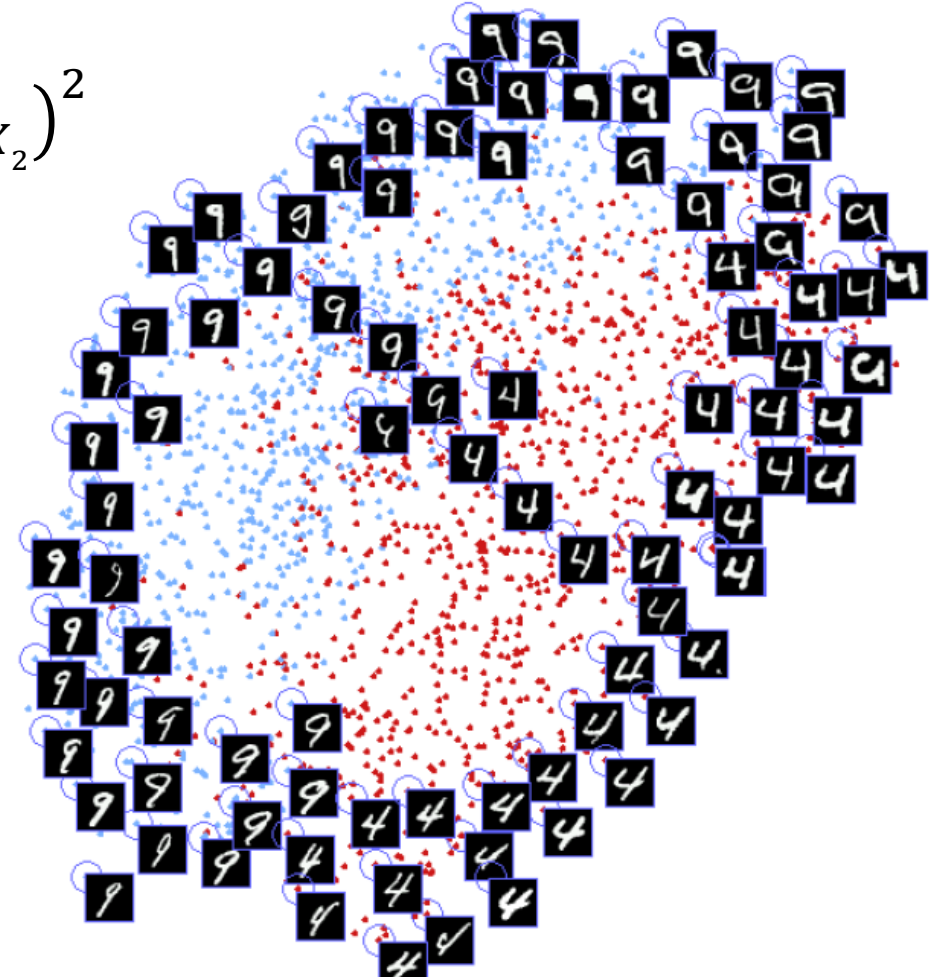
$D$  – distance between image embeddings

$$L_i(X_1, X_2) = D^2$$
$$L_i(X_1, X_2) = \max(0, m - D)^2$$

# Image Similarity: Loss functions: Contrastive loss

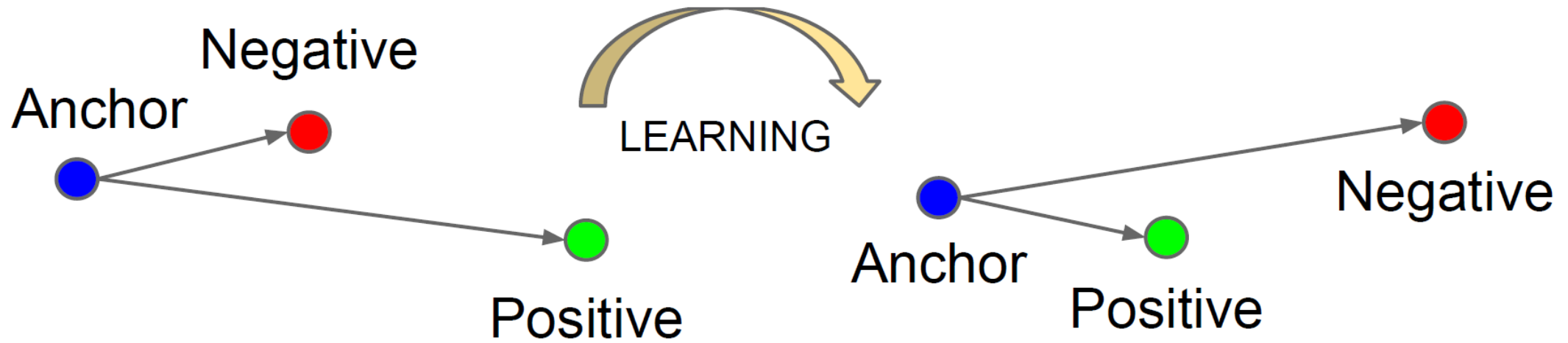
MNIST handwritten digits 0-9

$$L_i(X_1, X_2) = (1 - Y)D_{X_1, X_2}^2 + Y\max(0, m - D_{X_1, X_2})^2$$



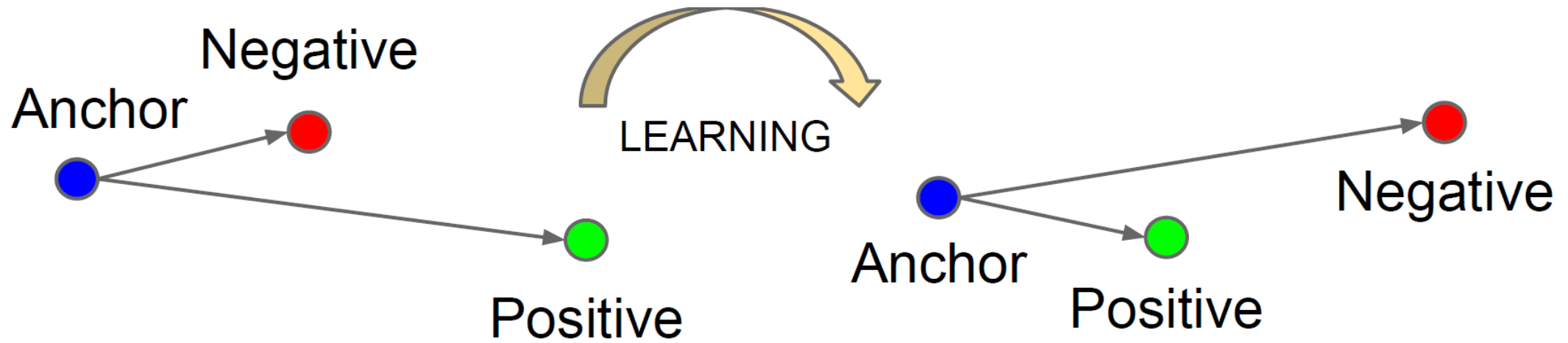
test

# Image Similarity: Loss functions: Triplet Loss



$$L_i(a, p, n) = \max(0, D_{a,p}^2 - D_{a,n}^2 + m)$$

# Image Similarity: Loss functions: Triplet Loss



$$L_i(a, p, n) = \max(0, D_{a,p}^2 - D_{a,n}^2 + m)$$

margin

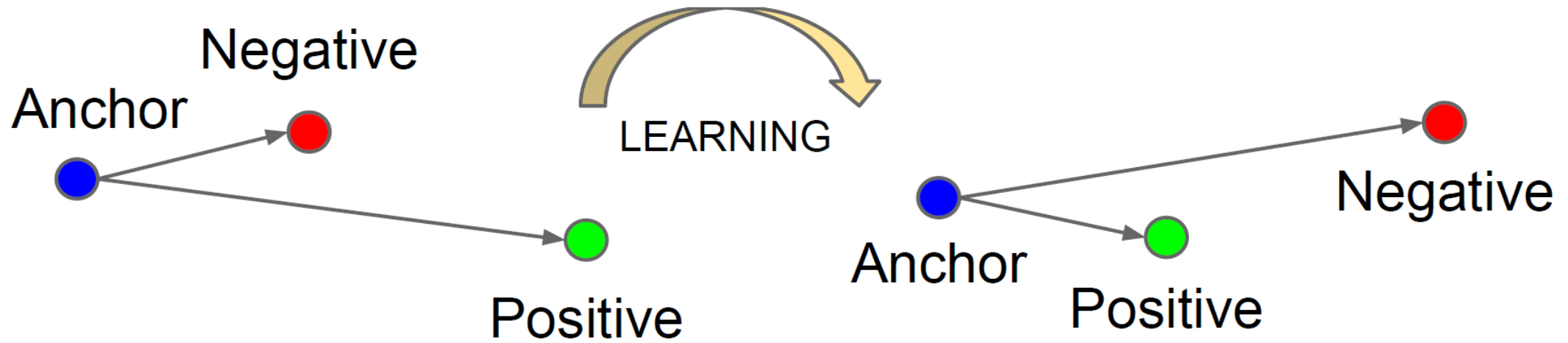
$a$  – anchor

$p$  – positive

$n$  – negative

distances between  
image embeddings

# Image Similarity: Loss functions: Triplet Loss



$$L_i(a, p, n) = \max(0, D_{a,p}^2 - D_{a,n}^2 + m)$$

Hard negative mining:  
for anchor-positive pair  
search for difficult negative:

$$\min_n (D_{a,n}^2)$$

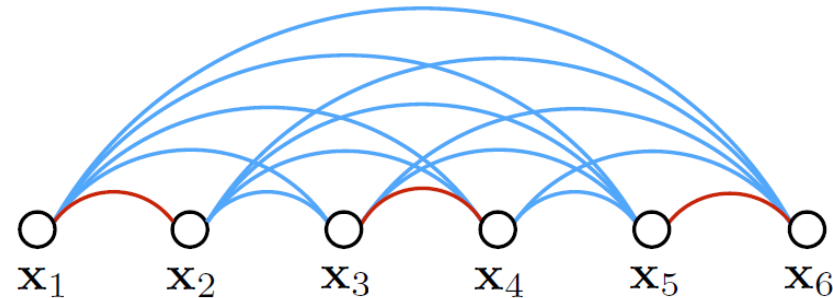
# Image Similarity: Loss functions: **Lifted Structured Loss**



(a) Contrastive embedding



(b) Triplet embedding



(c) Lifted structured embedding

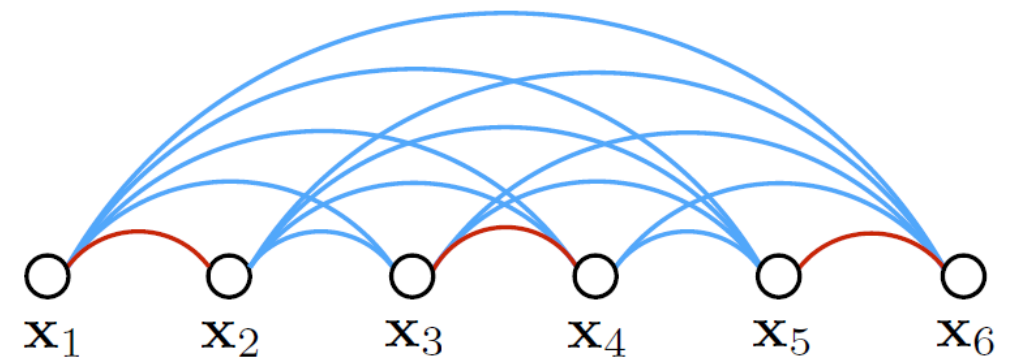
Automatic hard  
negative  
mining within  
batch

Training batch with six examples. Red edges and blue edges represent similar and dissimilar examples respectively.

# Image Similarity: Loss functions: **Lifted Structured Loss**

$$\tilde{J}_{i,j} = \log \left( \sum_{(i,k) \in N} \exp(m - D_{i,k}) + \sum_{(j,l) \in N} \exp(m - D_{j,l}) \right) + D_{i,j}$$

$$L_{batch} = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{J}_{i,j})^2$$



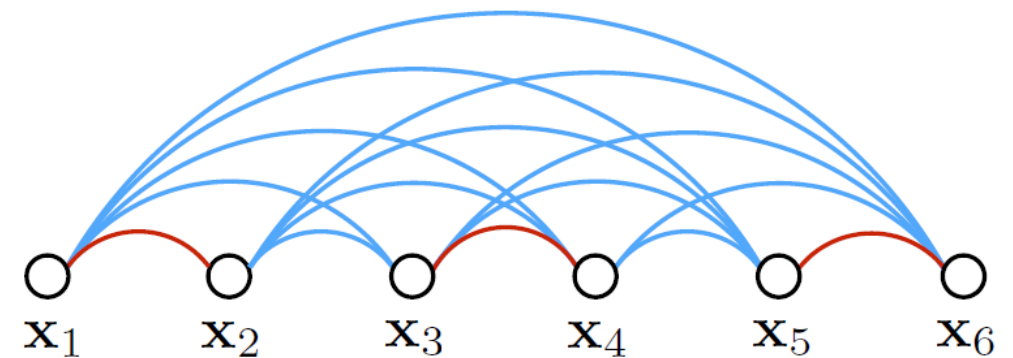
(c) Lifted structured embedding

# Image Similarity: Loss functions: **Lifted Structured Loss**

$$\tilde{J}_{i,j} = \log \left( \sum_{(i,k) \in N} \exp(m - D_{i,k}) + \sum_{(j,l) \in N} \exp(m - D_{j,l}) \right) + D_{i,j}$$

$$L_{batch} = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{J}_{i,j})^2$$

all positive pairs  
within batch



(c) Lifted structured embedding

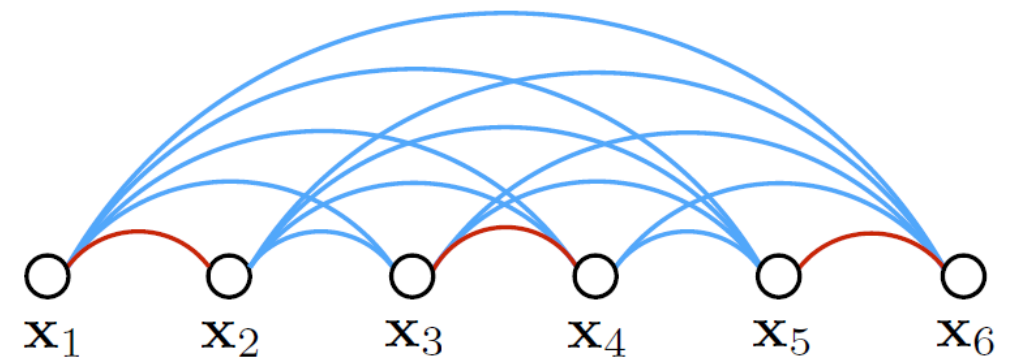


# Image Similarity: Loss functions: **Lifted Structured Loss**

$$\tilde{J}_{i,j} = \log \left( \sum_{(i,k) \in N} \exp(\overset{\text{margin}}{\boxed{m}} - D_{i,k}) + \sum_{(j,l) \in N} \exp(m - D_{j,l}) \right) + D_{i,j}$$

$$L_{batch} = \frac{1}{2|P|} \sum_{\boxed{(i,j) \in P}} \max(0, \tilde{J}_{i,j})^2$$

all positive pairs  
within batch



(c) Lifted structured embedding

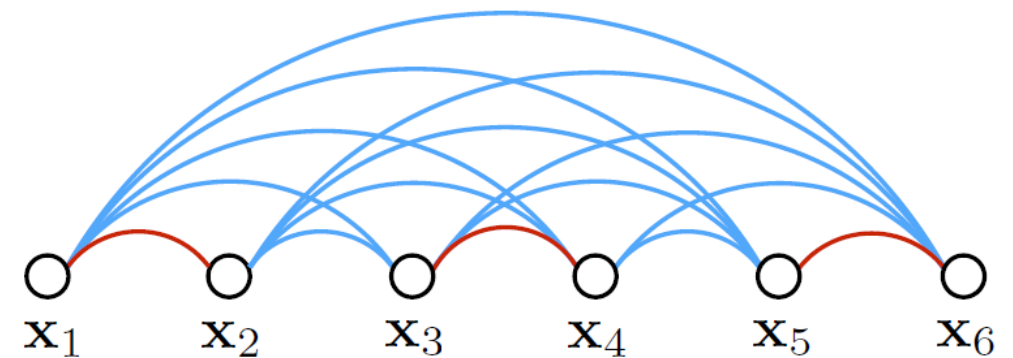
# Image Similarity: Loss functions: **Lifted Structured Loss**

$$\tilde{J}_{i,j} = \log \left( \sum_{(i,k) \in N} \exp(m - D_{i,k}) + \sum_{(j,l) \in N} \exp(m - D_{j,l}) \right) + D_{i,j}$$

all negative pairs for  $i$                       all negative pairs for  $j$

$$L_{batch} = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{J}_{i,j})^2$$

all positive pairs  
within batch



(c) Lifted structured embedding

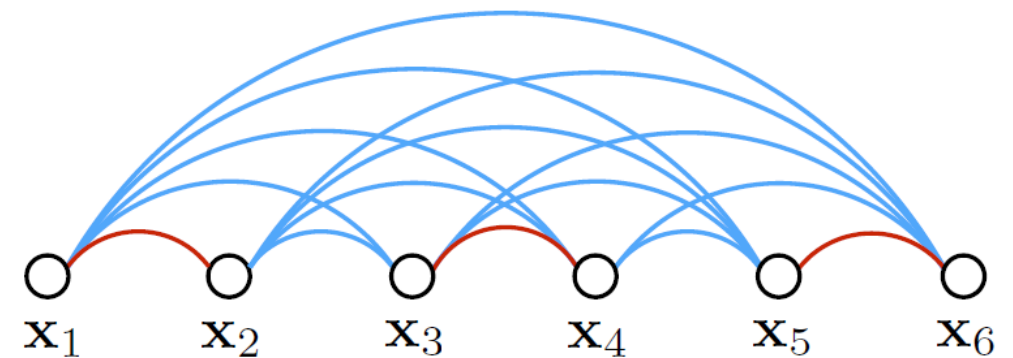
# Image Similarity: Loss functions: **Lifted Structured Loss**

softmax ☺ gives hardest negative for  $i$  and  $j$

$$\tilde{J}_{i,j} = \log \left( \sum_{(i,k) \in N} \exp(m - D_{i,k}) + \sum_{(j,l) \in N} \exp(m - D_{j,l}) \right) + D_{i,j}$$

$$L_{batch} = \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, \tilde{J}_{i,j})^2$$

all positive pairs  
within batch

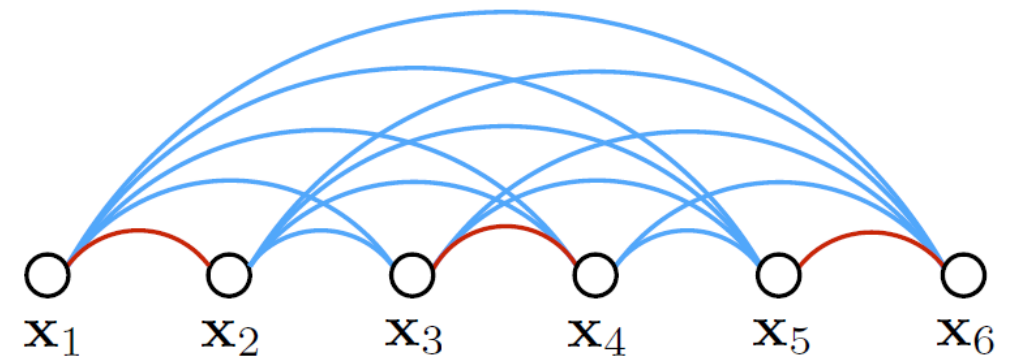


(c) Lifted structured embedding

# Image Similarity: Loss functions: **Lifted Structured Loss**

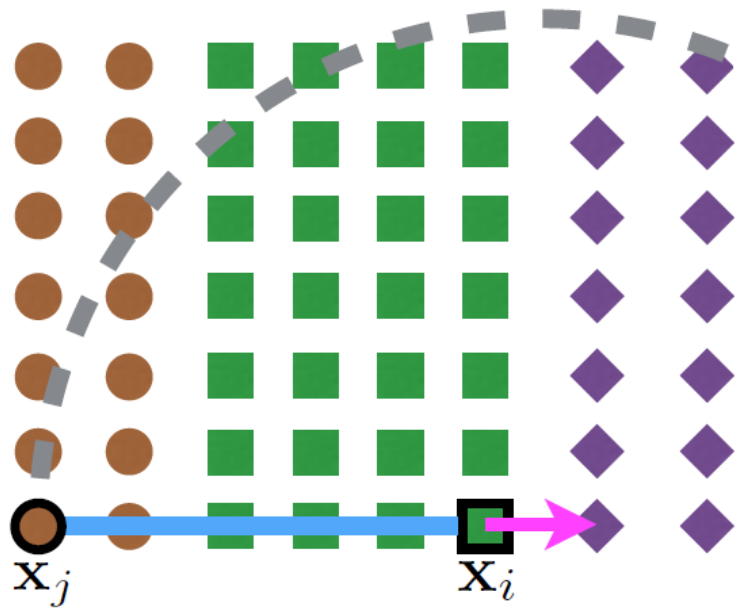
$$L_{batch} \sim \frac{1}{2|P|} \sum_{(i,j) \in P} \max(0, D_{i,j} - D_{i||j,hardest\ negative} + m)^2$$

Similar to triplet loss but with automatic hard negative mining within batch!

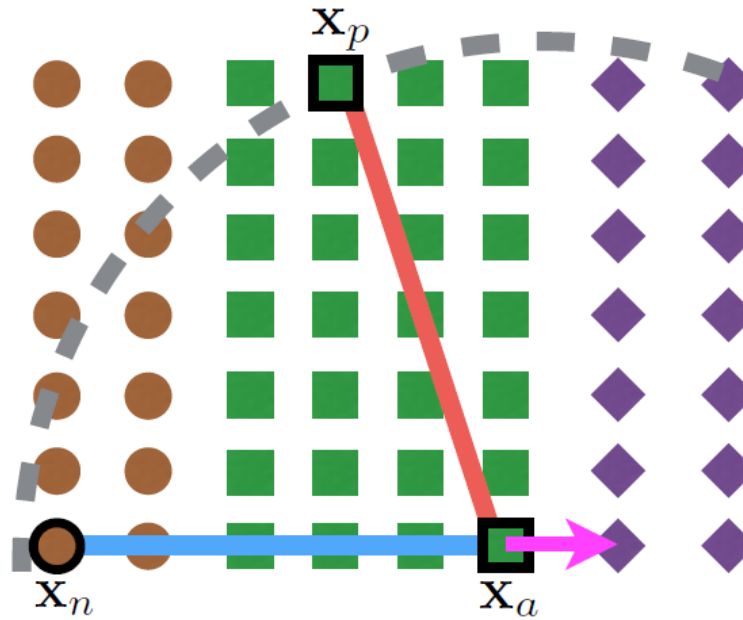


(c) Lifted structured embedding

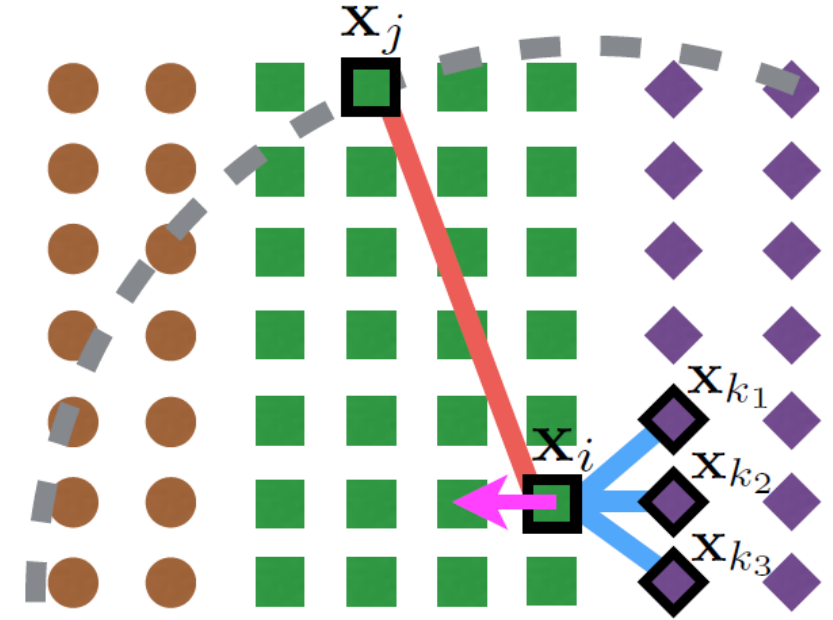
# Image Similarity: Loss functions: **Lifted Structured Loss**



(a) Contrastive embedding



(b) Triplet embedding



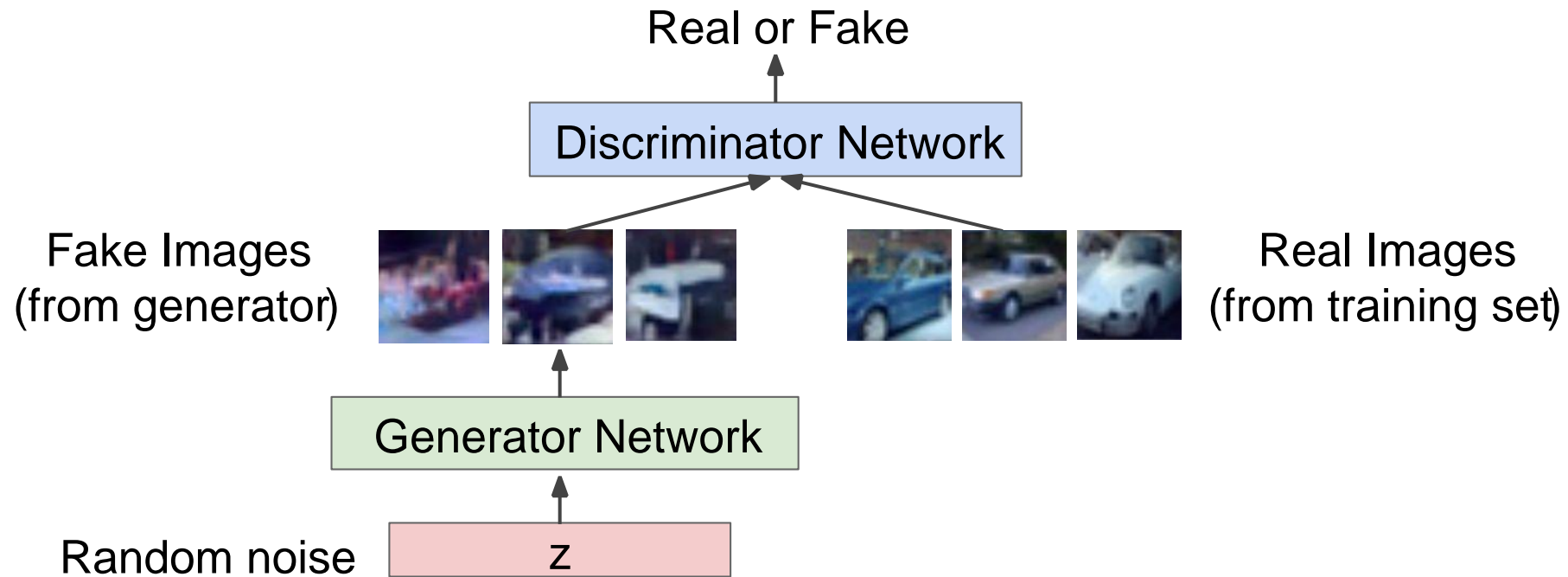
(c) Lifted structured similarity

Illustration of failure modes of contrastive and triplet loss with randomly sampled training batch.

# GANs

**Generator network:** try to fool the discriminator by generating real-looking images

**Discriminator network:** try to distinguish between real and fake images



# Training GANs: Two-player game

**Generator network:** try to fool the discriminator by generating real-looking images

**Discriminator network:** try to distinguish between real and fake images

Train jointly in **minimax game**

Minimax objective function: Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[ \underbrace{\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x)}_{\substack{\text{Discriminator output} \\ \text{for real data } x}} + \mathbb{E}_{z \sim p(z)} \log \underbrace{(1 - D_{\theta_d}(G_{\theta_g}(z)))}_{\substack{\text{Discriminator output for} \\ \text{generated fake data } G(z)}} \right]$$

# Wasserstein GAN

Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$

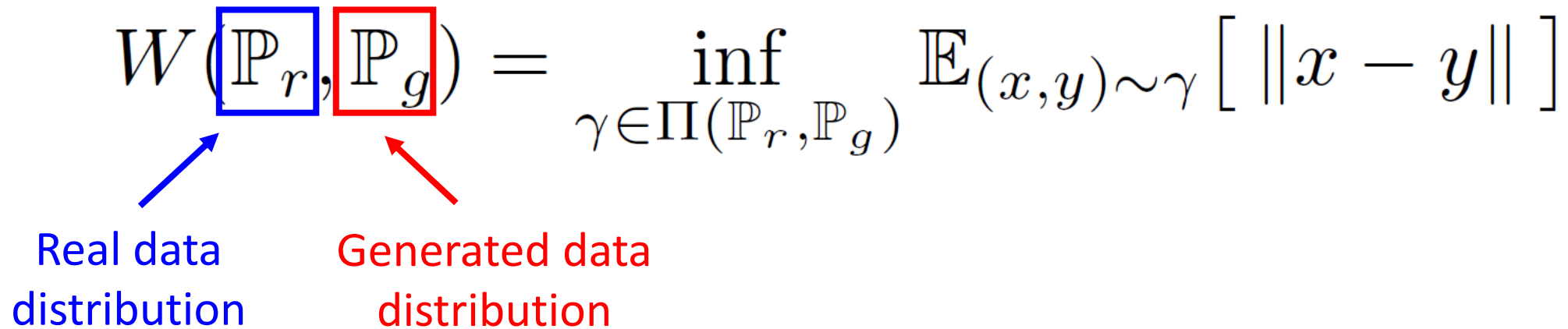


# Wasserstein GAN

Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$

Real data distribution      Generated data distribution



# Wasserstein GAN

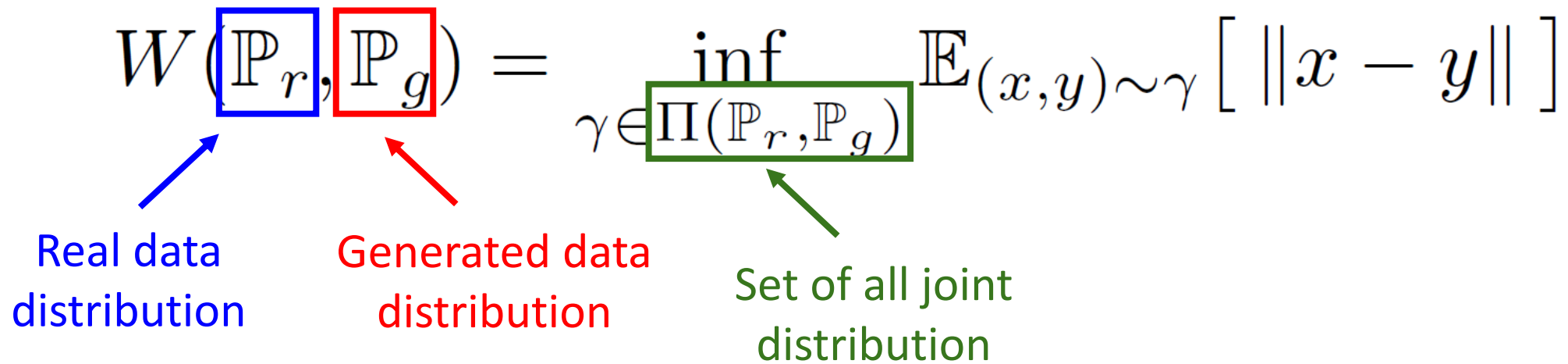
Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$

Real data distribution

Generated data distribution

Set of all joint distribution



# Wasserstein GAN

Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$

Real data distribution      Generated data distribution      Set of all joint distribution

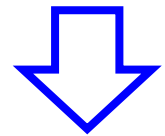
Intuitively,  $\gamma(x, y)$  indicates how much “mass” must be transported from  $x$  to  $y$  in order to transform the distributions  $P_r$  into the distribution  $P_g$ .

The EM distance is the “cost” of the optimal transport plan.

# Wasserstein GAN: CNN approximation

Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$



Kantorovich-Rubinstein theorem

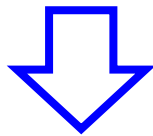
$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)]$$

supremum is over all the 1-Lipschitz functions

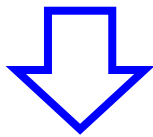
# Wasserstein GAN: CNN approximation

Minimize Earth-Mover (EM) distance or Wasserstein distance between  $P_r$  **real data distribution** and  $P_g$  **generated data distribution**

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [ \|x - y\| ]$$



$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)]$$



approximate  $f(x)$  by critic CNN  $f_w(x)$

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim \mathbb{P}_r} [f_w(x)] - \mathbb{E}_{z \sim p(z)} [f_w(g_\theta(z))]$$

# Wasserstein GAN: Training

---

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

---

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

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

```

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

```

---

# Wasserstein GAN: Training

---

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

---

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

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

```

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

```

---

train critic



# Wasserstein GAN: Training

---

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

---

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

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

```

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

train generator

# Wasserstein GAN: Training

---

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

---

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

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

```

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

```

---

gradient of critic parameters  
 using Wasserstein distance  
 approximation with neural  
 network

# Wasserstein GAN: Training

---

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

---

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

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

```

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

```

---

update of critic parameters

# Wasserstein GAN: Training

---

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

---

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

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

```

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

---

clip parameters of critic

# Wasserstein GAN: Training

---

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

---

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

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

```

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

```

---

gradient of generator parameters

# Wasserstein GAN: Training

---

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

---

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

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

```

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

update of generator parameters

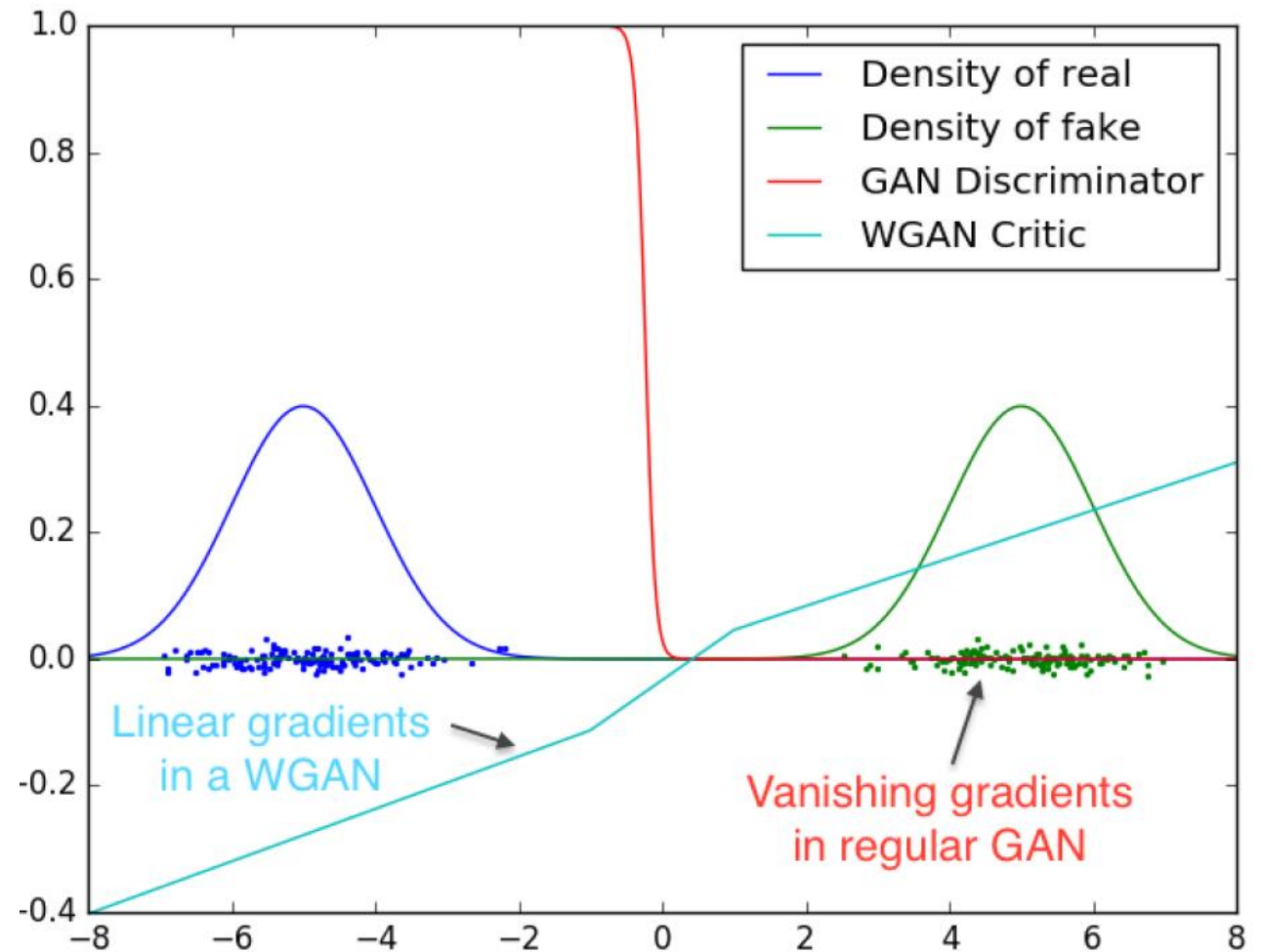
---



# Wasserstein GAN: No Vanishing Gradients

Optimal discriminator and critic when learning to differentiate two Gaussians. Traditional GAN discriminator saturates and results in vanishing gradients.

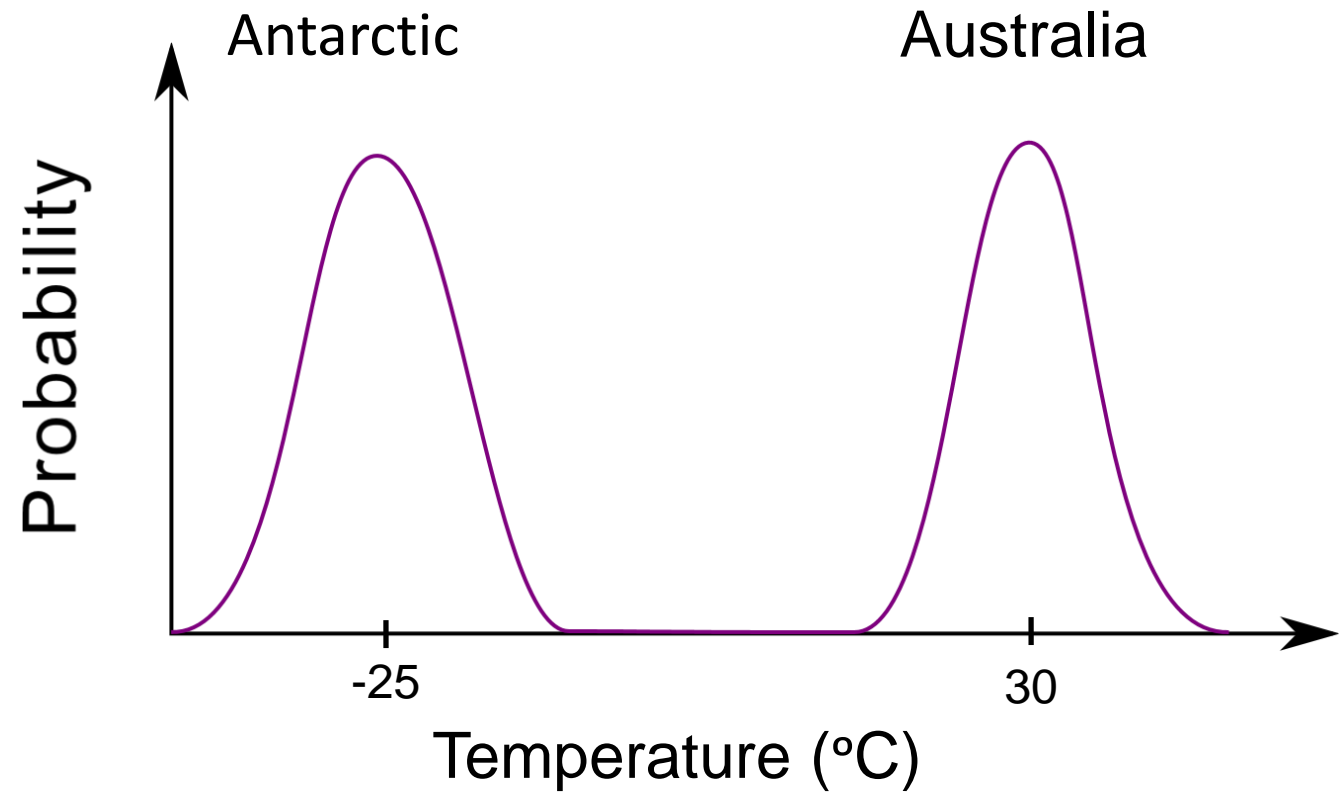
Wasserstein GAN critic provides very clean gradients on all parts of the space.





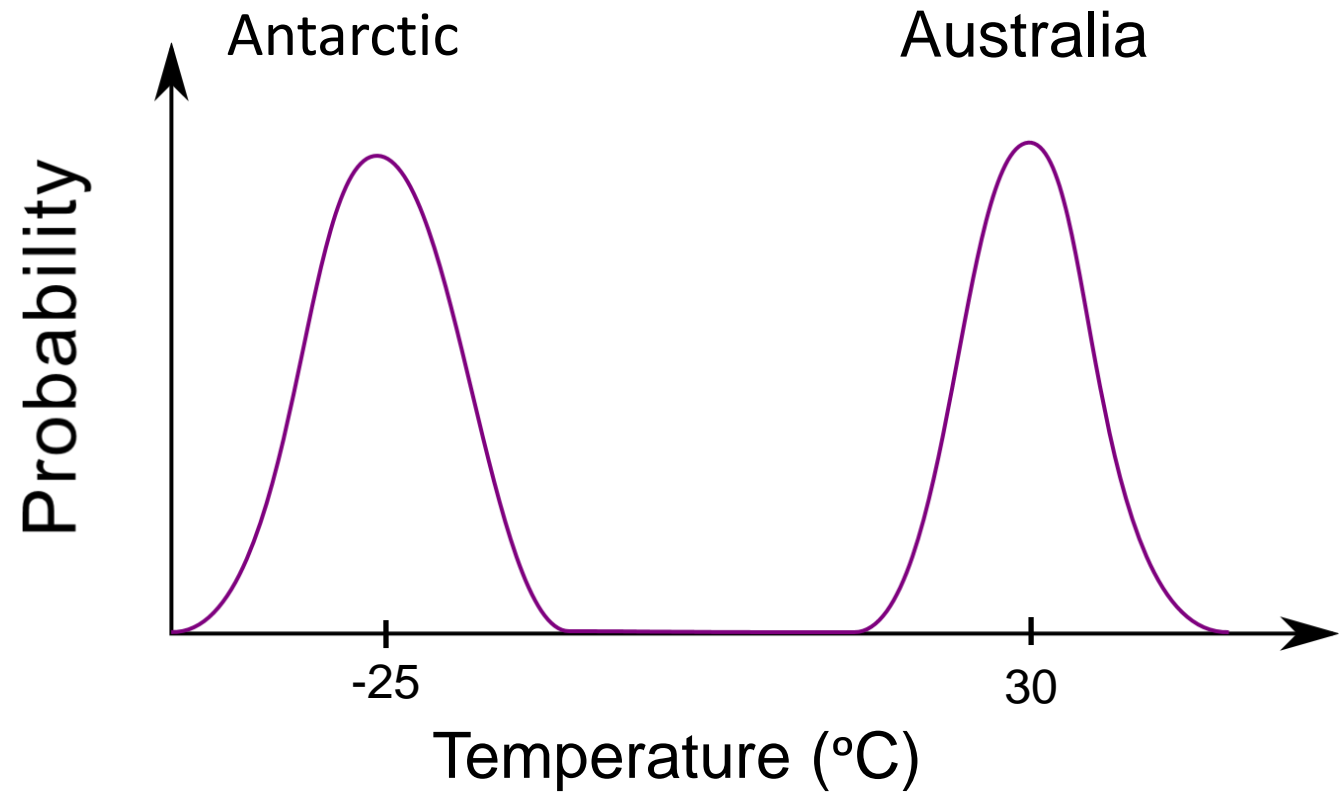
# Mode collapse in GANs

1. Generator learns that it can fool the discriminator by producing values close to Antarctic temperatures
2. The discriminator counters by learning that all Australian temperatures are real, and guesses whether Antarctic temperatures are real or fake
3. The generator exploits the discriminator by switching modes to produce values close to Australian temperatures instead, abandoning the Antarctic mode
4. The discriminator now assumes that all Australian temperatures are fake and Antarctic temperatures are real
5. Return to step 1



# Mode collapse in GANs

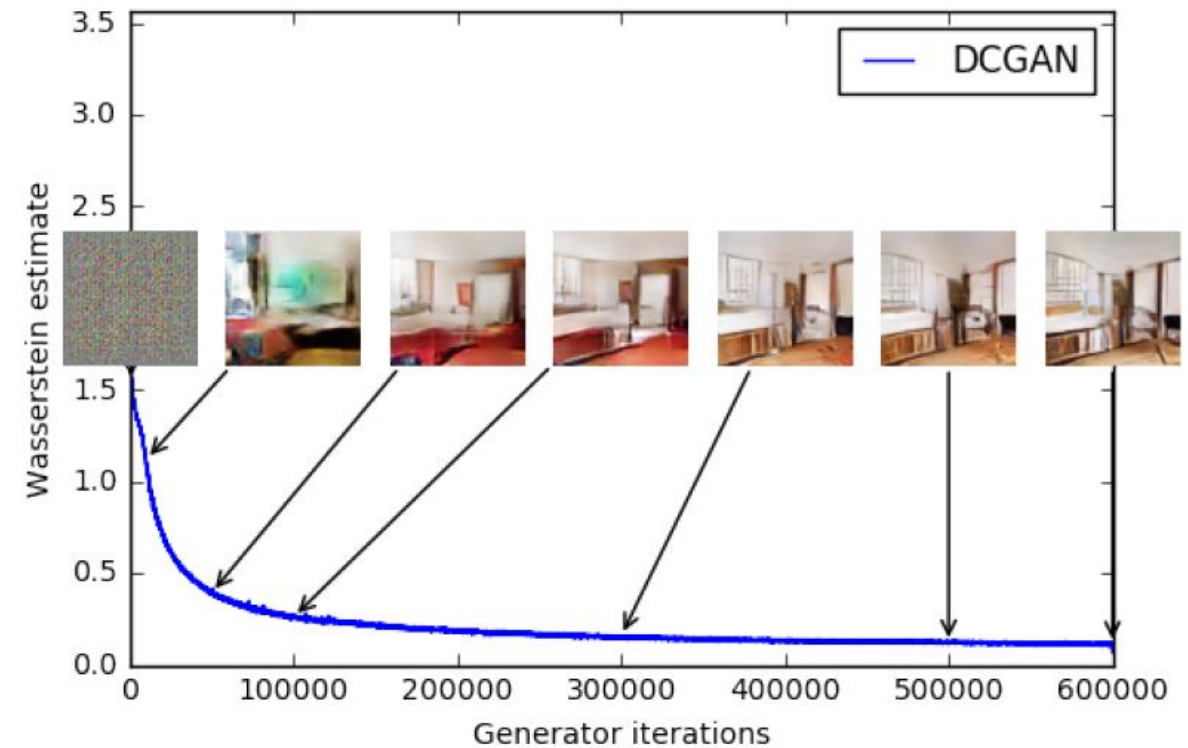
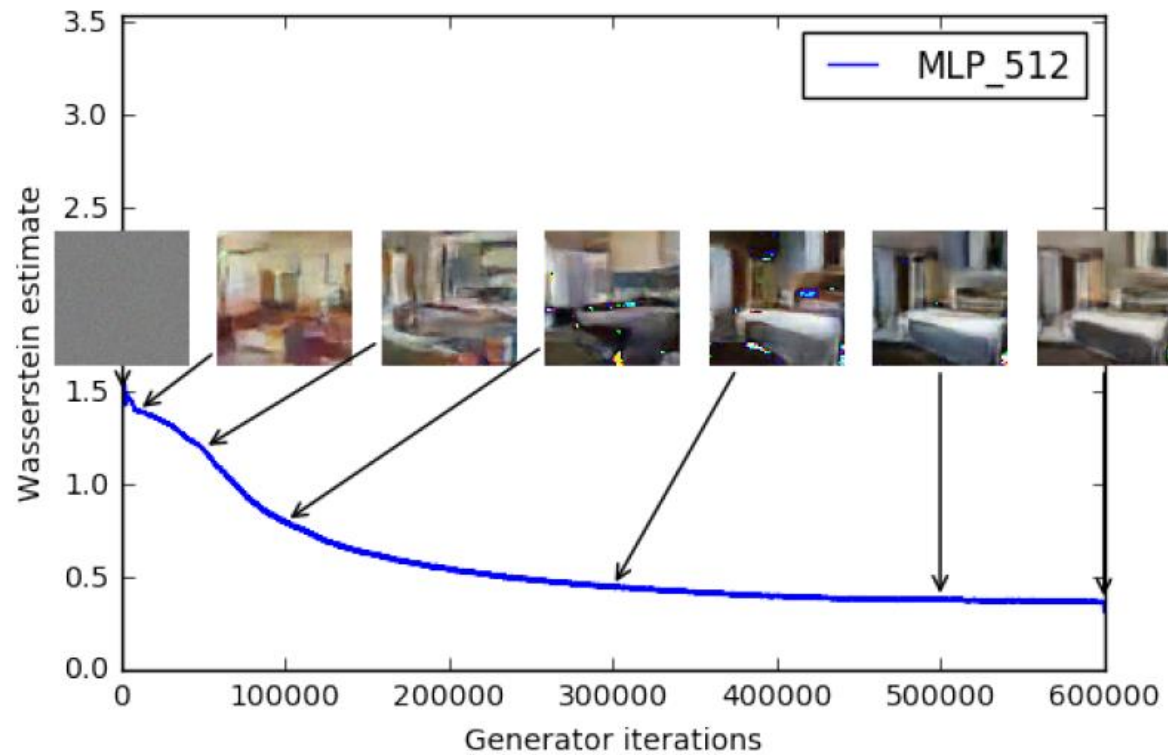
1. Generator learns that it can fool the discriminator by producing values close to Antarctic temperatures
2. The discriminator counters by learning that all Australian temperatures are real, and guesses whether Antarctic temperatures are real or fake
3. The generator exploits the discriminator by switching modes to produce values close to Australian temperatures instead, abandoning the Antarctic mode
4. The discriminator now assumes that all Australian temperatures are fake and Antarctic temperatures are real
5. Return to step 1



**Wasserstein GANs do not have mode collapse**

# Wasserstein GAN

Nice property lower loss better sample quality



# ICCV 2017: Venice, Italy, Oct 22-29



Growth of paper  
submissions  
~30%

Growth of participants  
~200%



# ICCV 2017: We are hiring



We are hiring 😊