

# Learning and forgetting in image classification and generation

Luis Herranz  
Learning and machine perception (LAMP) group  
Computer Vision Center (Barcelona)

# About me



ETS Ing. Telecomunicación  
Universidad Politécnica de Madrid (Spain)



(Ph.D) Escuela Politécnica Superior  
Universidad Autónoma de Madrid (Spain)



Misubishi Electric R&D (UK)



Institute of Computing Technology,  
Chinese Academy of Sciences (China)



Computer Vision Center,  
Universitat Autònoma de Barcelona (Spain)

Research interests:  
Multimedia  
Computer vision  
Deep learning  
Multimodal  
representations  
Lifelong learning

# Computer Vision Center (UAB campus)



**Only Center in Europe fully devoted to Computer Vision**

**23** Years

**+130** Staff

**+20** Nationalities

**M€2,3** Income /  
year

**8** PhD thesis  
/year

**+100** Intl  
publications /  
year

# Learning and Machine Perception (LAMP) group

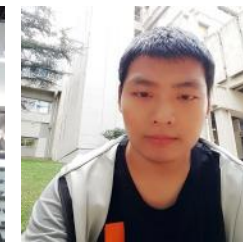
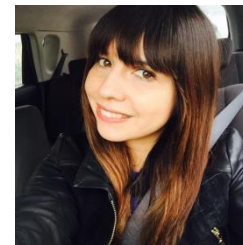
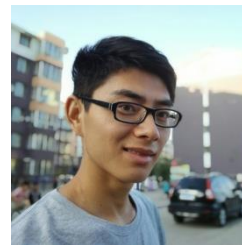
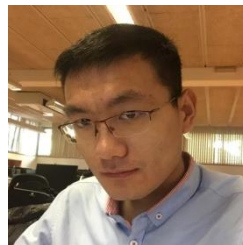
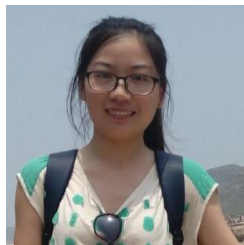
## Senior PhDs

## Postdocs

**Joost van de Weijer**  
Leader



## PhD students



# Outline

- Introduction
- Transferring GANs (ECCV 2018)
- Rotated elastic weight consolidation (ICPR 2018)
- Memory Replay GANs (NIPS 2018)
- Mix and match networks (CVPR 2018)

# Outline

- Introduction
- Transferring GANs (ECCV 2018)
- Rotated elastic weight consolidation (ICPR 2018)
- Memory Replay GANs (NIPS 2018)
- Mix and match networks (CVPR 2018)

# Transfer learning and lifelong learning

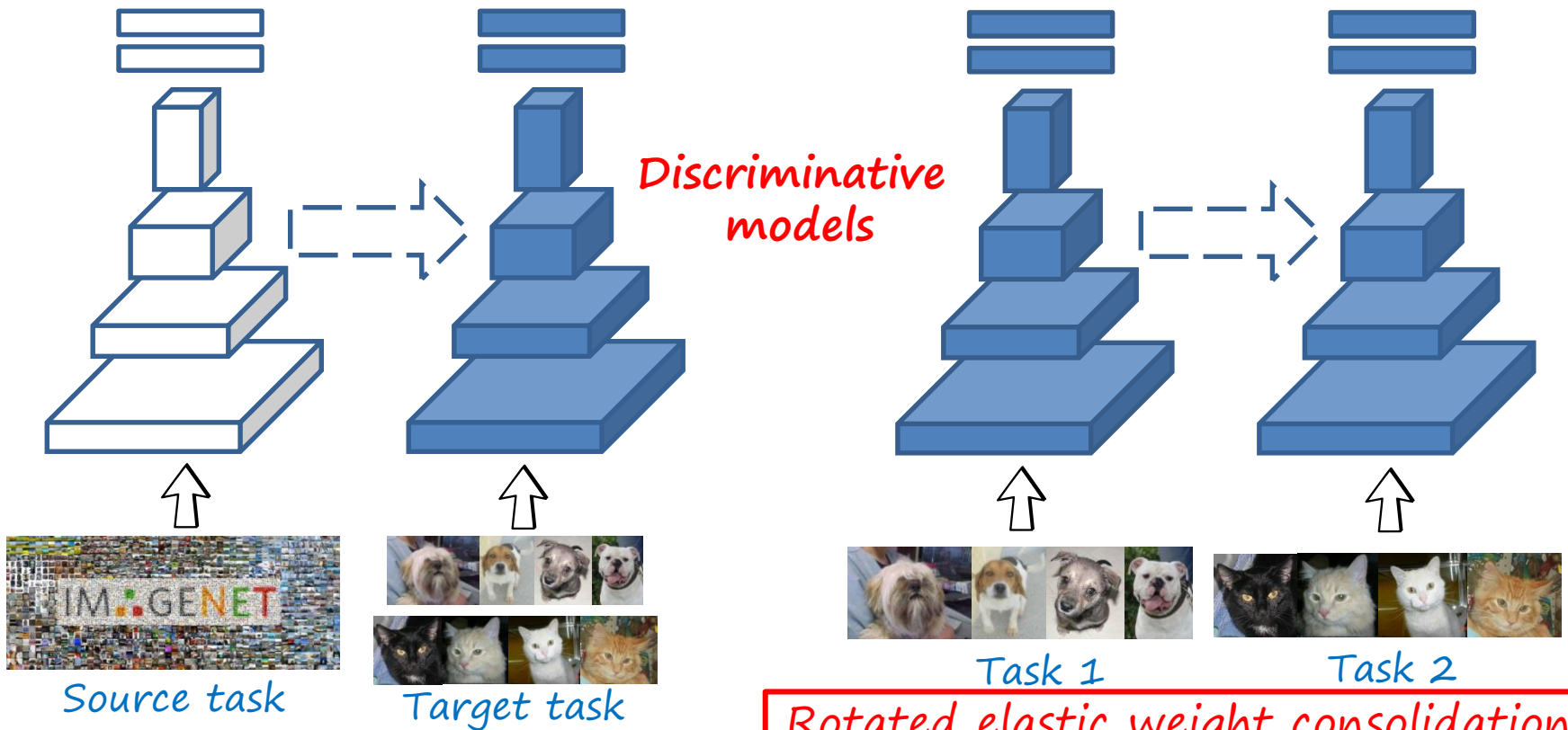
Forgets source task, i.e. catastrophic forgetting (who cares?)

Forgets task 1 (big deal!!)

Transfer+adaptation

Lifelong learning

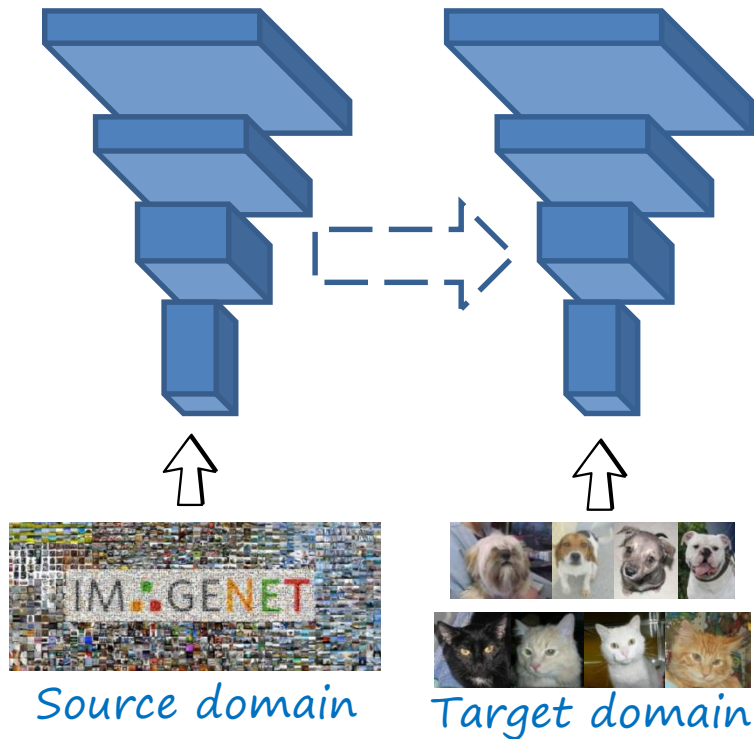
Discriminative models



Rotated elastic weight consolidation (ICPR 2018)

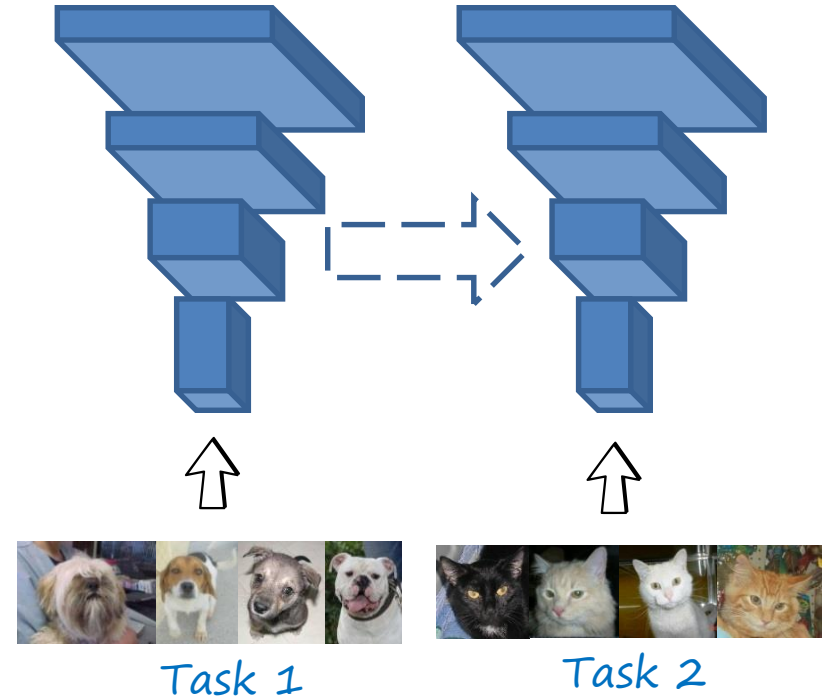
# Transfer learning and lifelong learning (now with GANs for image generation)

Transfer+adaptation (generative)



Transferring GANs  
(ECCV 2018)

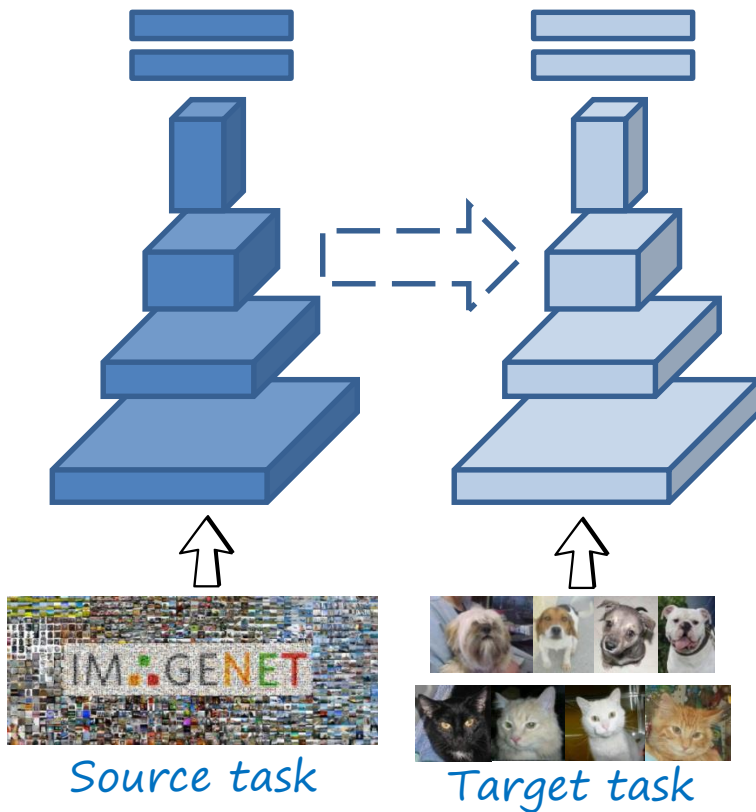
Lifelong learning (generative)



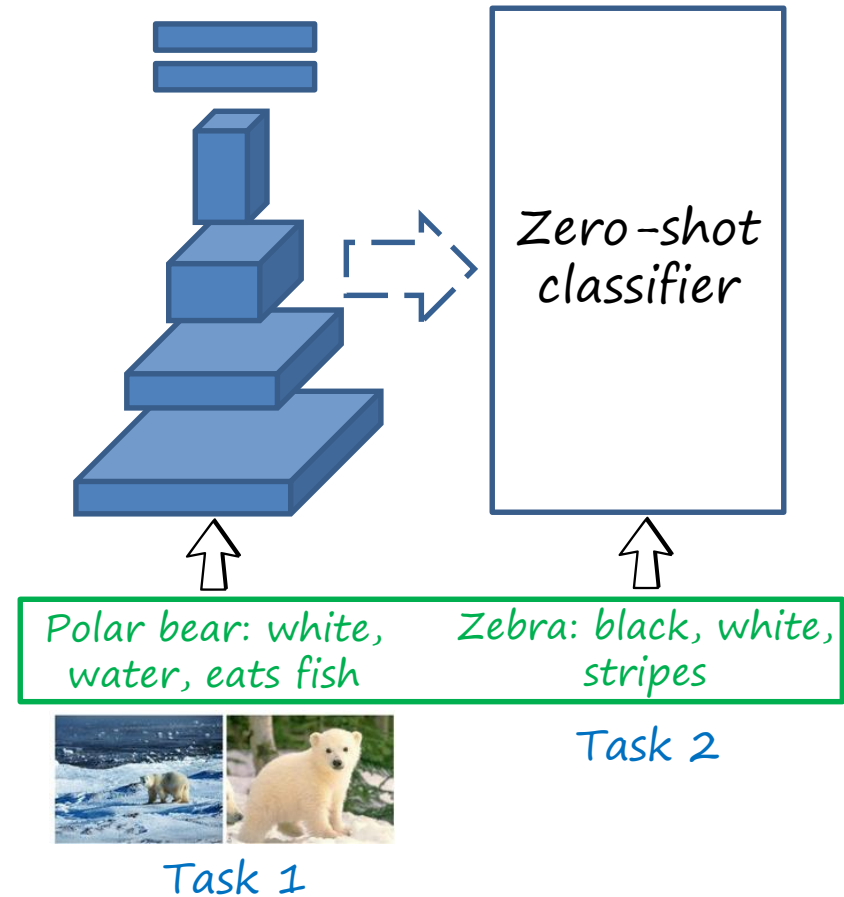
Memory Replay GANs  
(NIPS 2018)

# Transfer learning and zero-shot learning

Transfer+adaptation



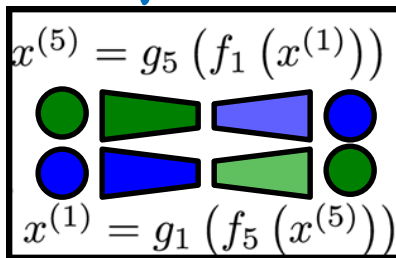
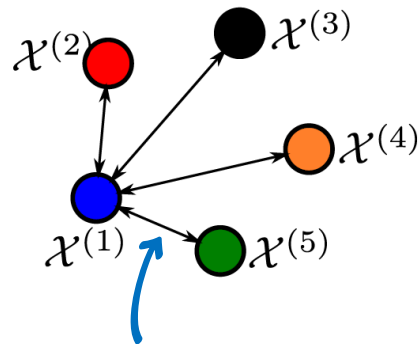
Zero-shot learning



# Zero-pair image-to-image translations

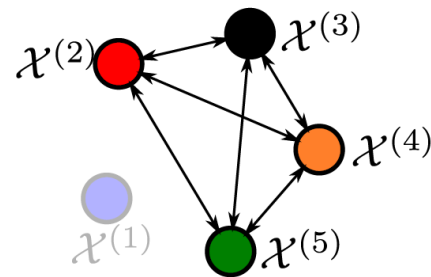
Only these translations  
are trained (seen)

Train



Evaluate on these unseen  
translations (no training pairs)

Test



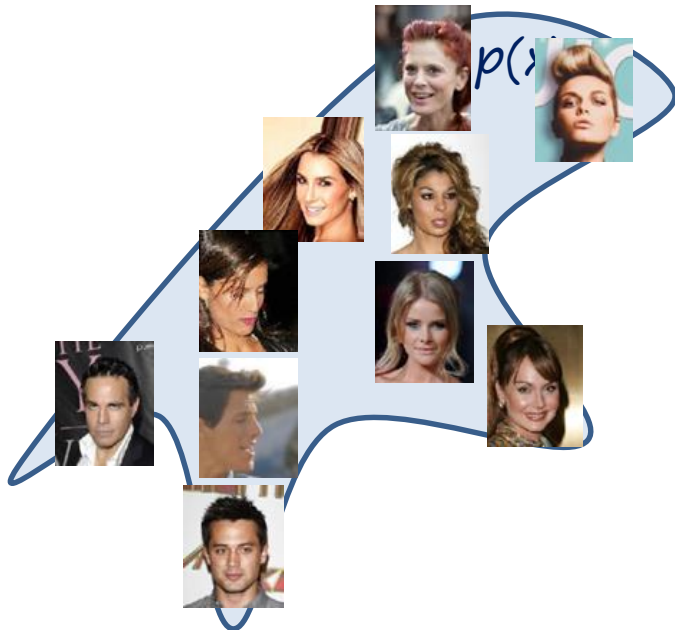
Mix and match networks  
(CVPR 2018)

# Outline

- Introduction
- **Transferring GANs (ECCV 2018)**
- Rotated elastic weight consolidation (ICPR 2018)
- Memory Replay GANs (NIPS 2018)
- Mix and match networks (CVPR 2018)

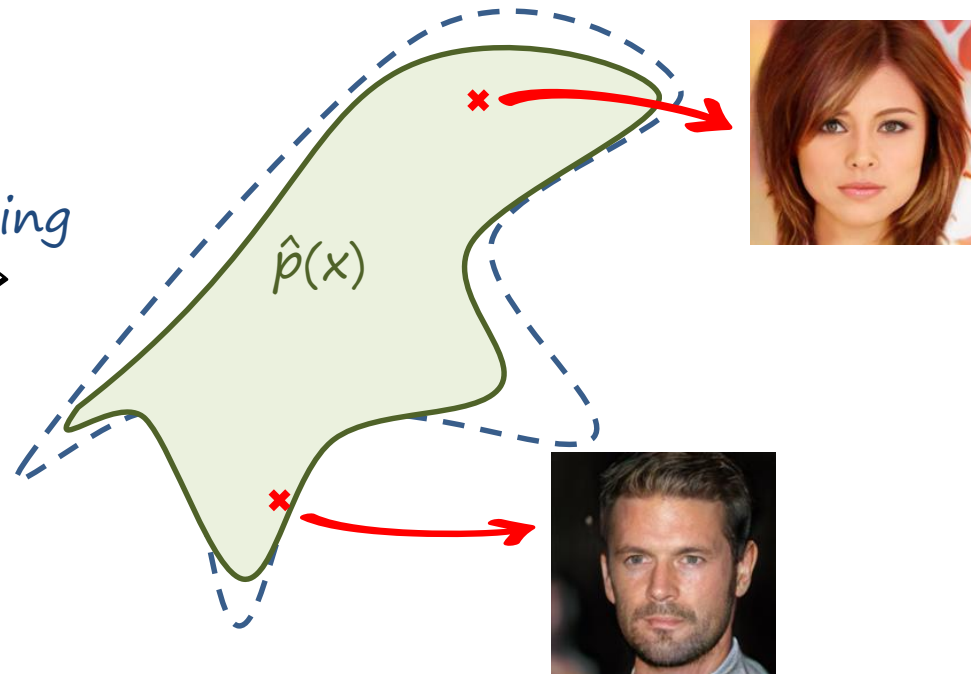
# Generative models: networks that imagine

Training data  
(e.g.  $64 \times 64 \times 3 \approx 12\text{K}$  dims)



Learning  
→

Sampling

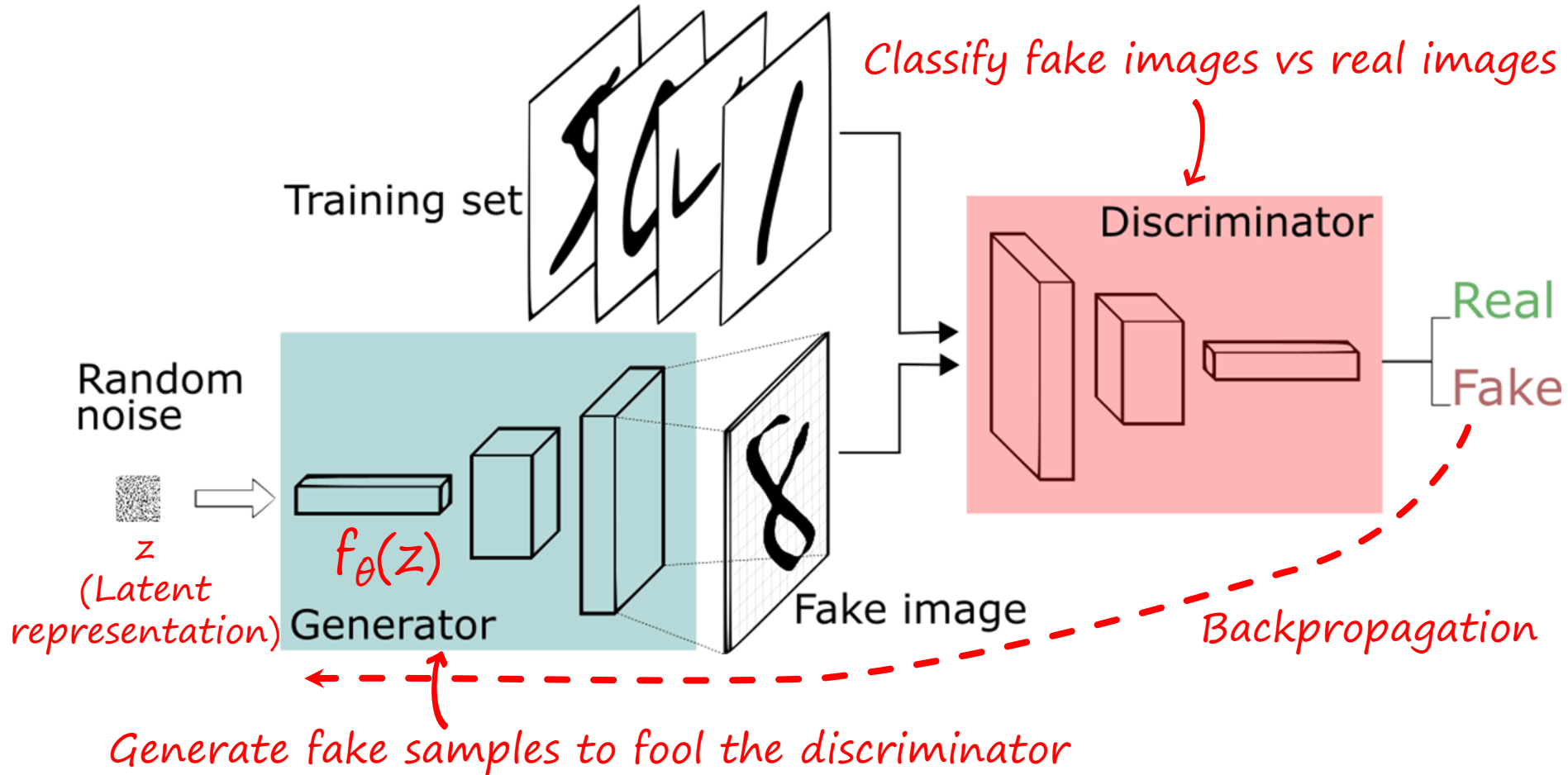


Different approaches

- Density estimation
- Variational autoencoders
- Autoregressive models

- **Generative adversarial networks (GANs)**

# Generative Adversarial Networks (GANs)



Goodfellow et al., "Generative Adversarial Networks", NIPS 2014

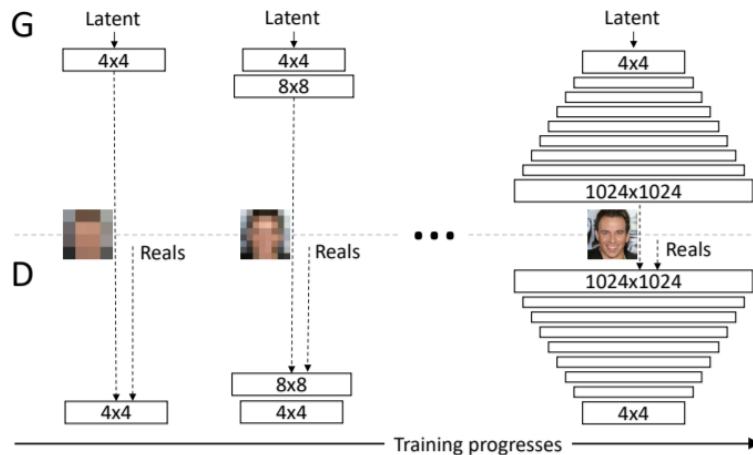
Figure from <https://deeplearning4j.org/generative-adversarial-network>

# Generative Adversarial Networks

## Wasserstein GAN (WGAN-GP)



## Progressive growing of GANs

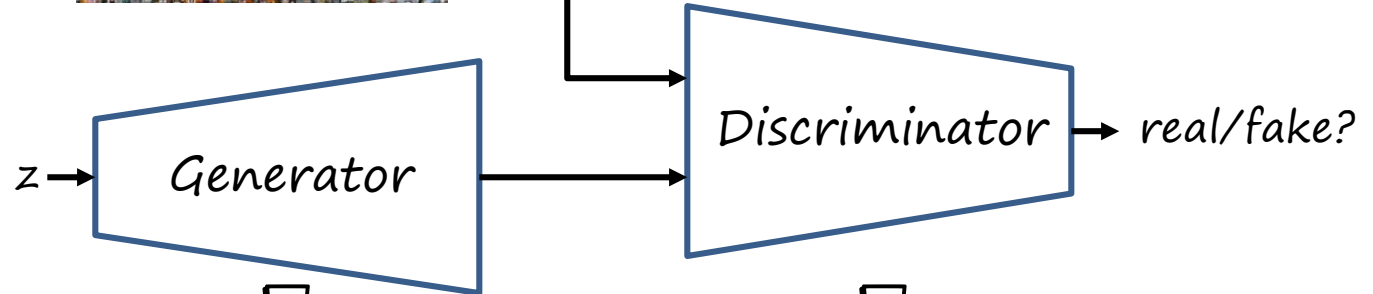


and many more...

# Transferring GAN representations

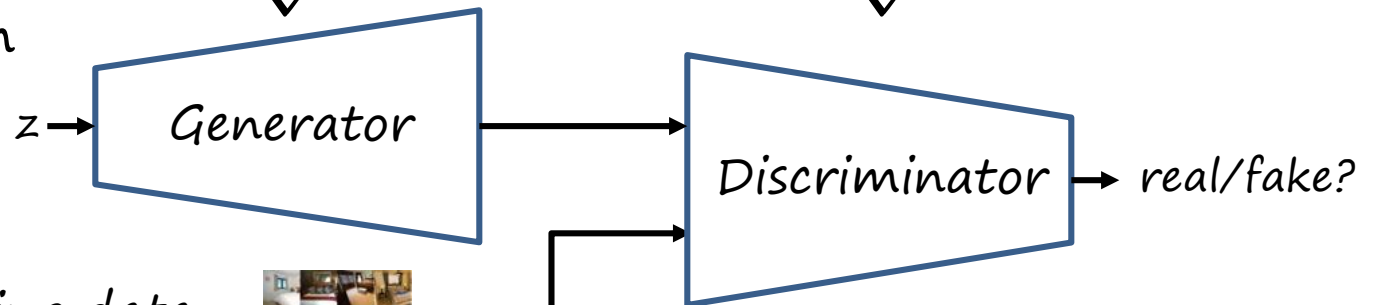
Source domain

Training data  
(ImageNet)



Target domain

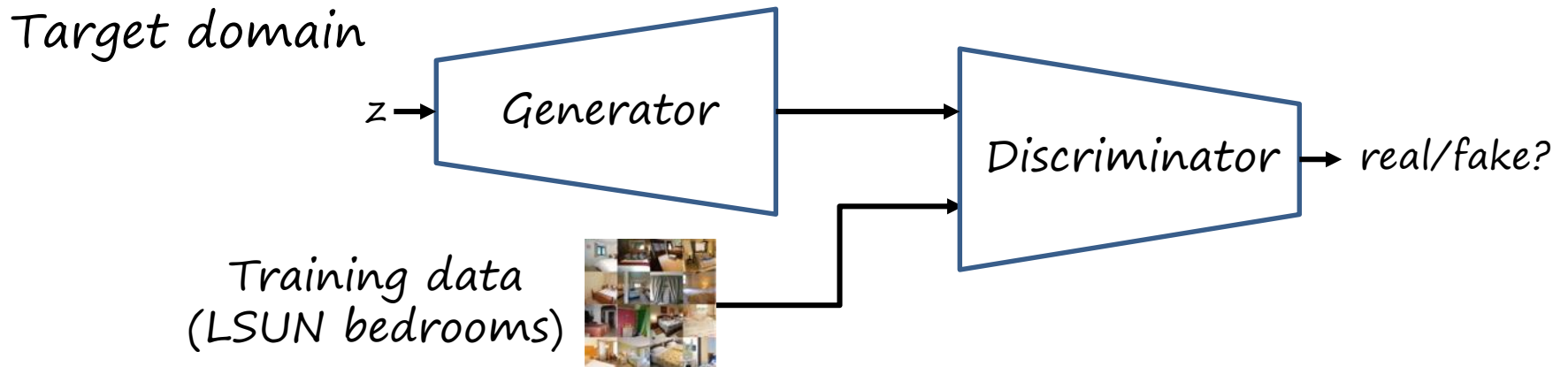
Training data  
(LSUN bedrooms)



# Transfer configuration

Training from scratch

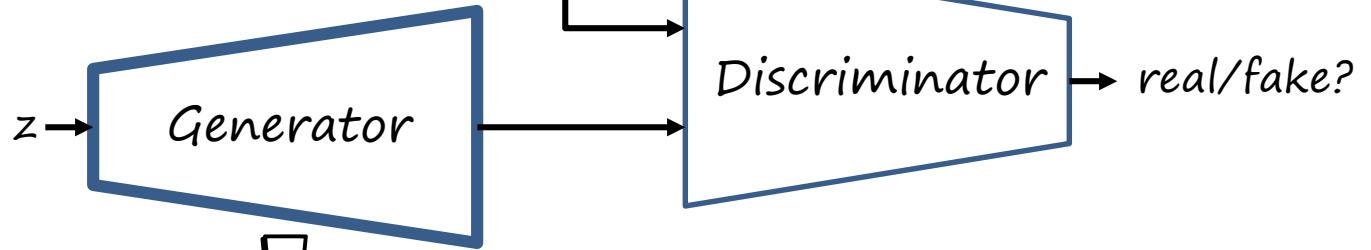
- Discr: from scratch
- Gen: from scratch



# Transfer configuration

Source domain

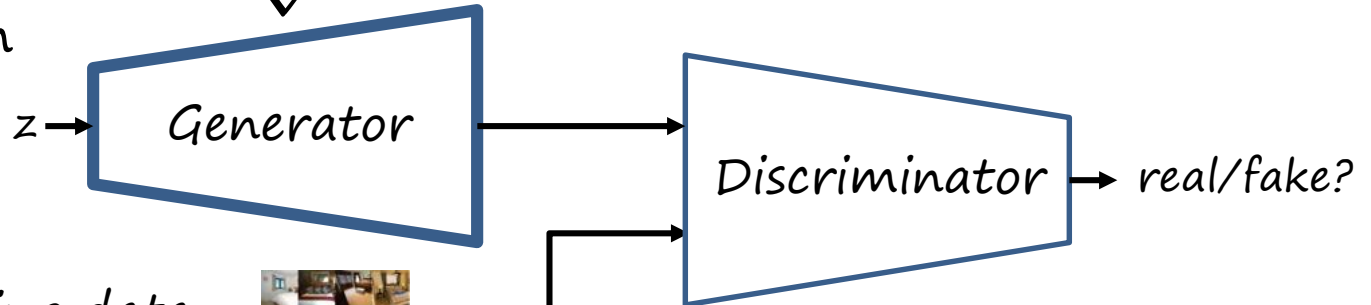
Training data  
(ImageNet)



Transfer only discr.  
- Discr: pretrained  
- Gen: from scratch

Target domain

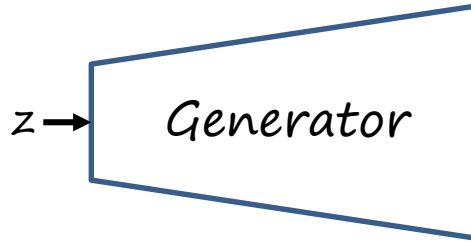
Training data  
(LSUN bedrooms)



# Transfer configuration

Source domain

Training data  
(ImageNet)



real/fake?

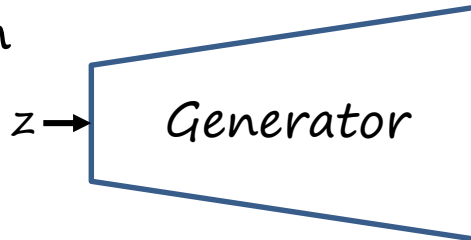
Transfer only gen.

- Discr: from scratch
- Gen: pretrained



Target domain

Training data  
(LSUN bedrooms)



real/fake?

# Transfer configuration

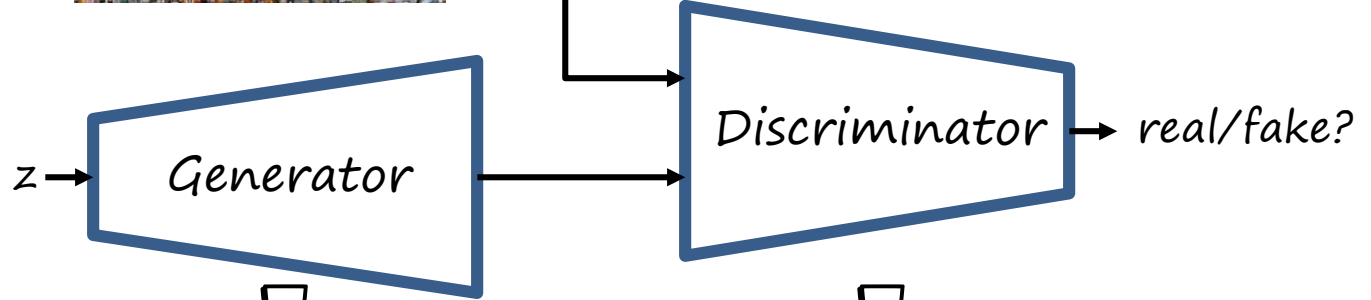
Source domain

Training data  
(ImageNet)



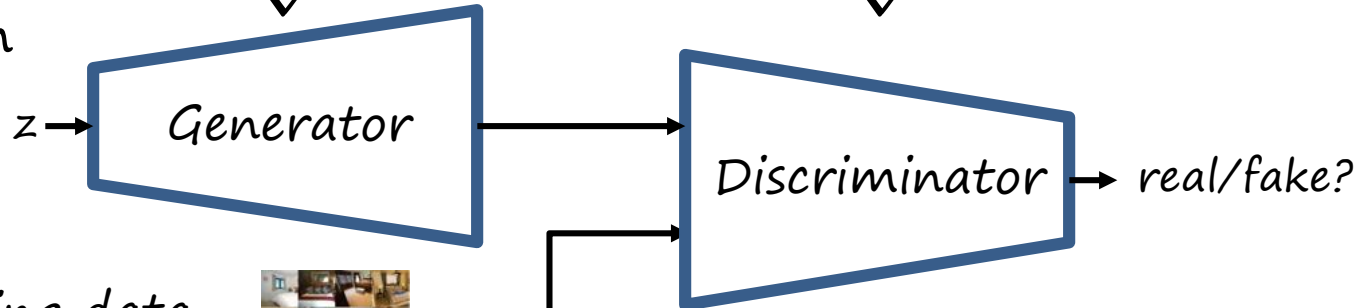
Transfer both

- Discr: pretrained
- Gen: pretrained



Target domain

Training data  
(LSUN bedrooms)



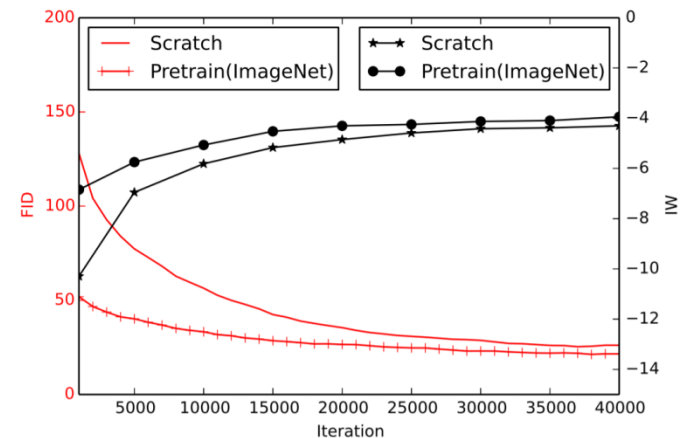
# Transfer configuration

- What should I transfer?
  - Experiment: ImageNet to Bedrooms (100K images)

Generator		Scratch		Pretrained	
Discriminator		Scratch	Pretrained	Scratch	Pretrained
<i>Lower better</i>	FID ( $\mathcal{X}_{data}^{tgt}, \mathcal{X}_{gen}^{tgt}$ )	32.87	30.57	56.16	<b>24.35</b>
<i>Higher better</i>	IW ( $\mathcal{X}_{val}^{tgt}, \mathcal{X}_{gen}^{tgt}$ )	-4.27	-4.02	-6.35	<b>-3.88</b>

*None   Only discr.   Only gen.   Both*

- Training is faster and images have better quality
  - Especially when data is limited



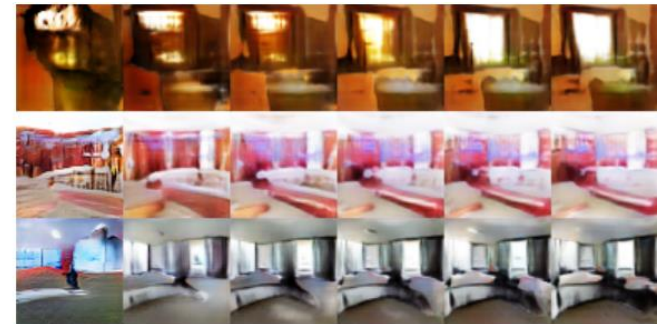
# Learning with limited data

Target dataset:  
Bedroom

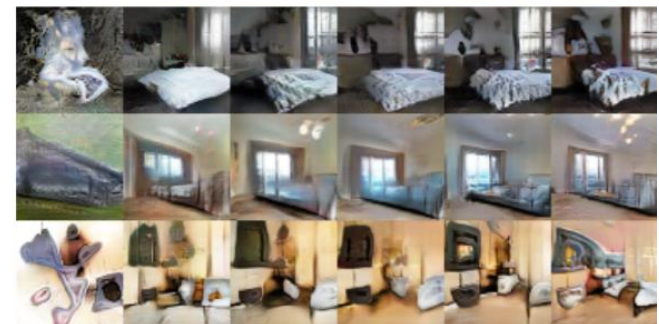
1000 images



Pretrained (ImageNet)



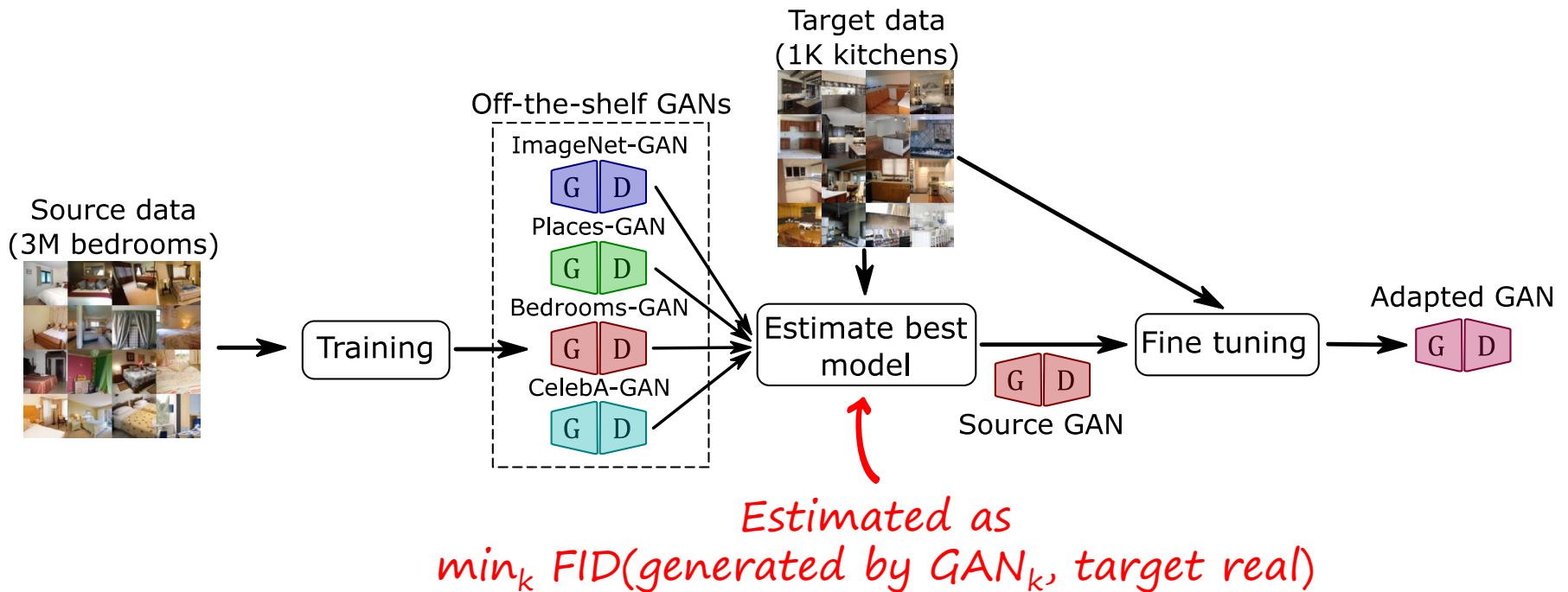
10000 images



100000 images

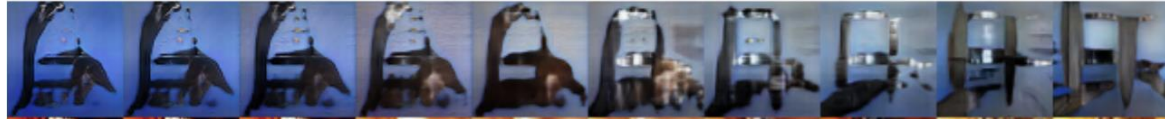


# Pretrain model selection



# Good/bad source datasets

Source



Target

Places

(205 classes  
~10K img/class)



Kitchen  
(50K)



Bedrooms  
(1 class  
3M images)



Kitchen  
(50K)



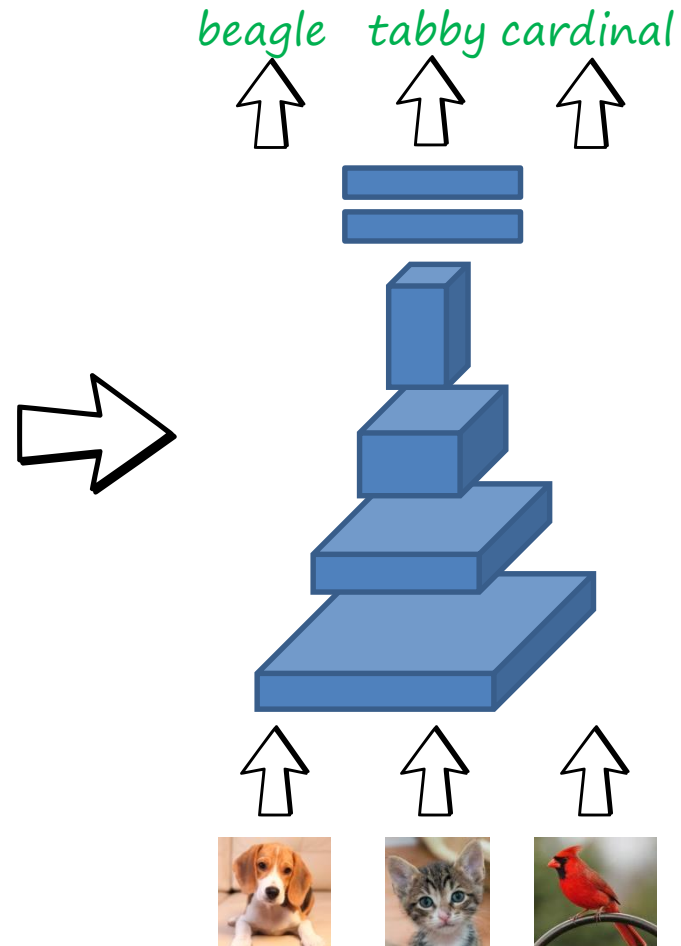
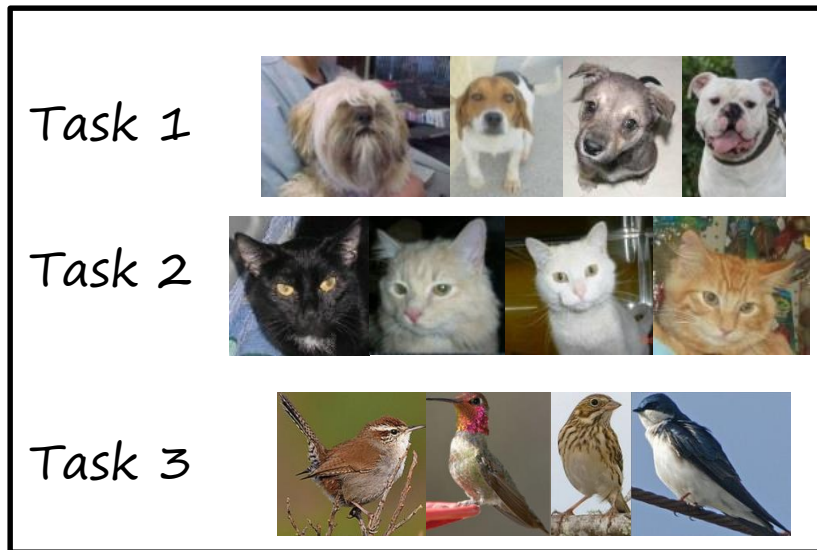
*Generative: very dense, diversity not so important (e.g. LSUN Bedrooms)*

*Discriminative: very diverse, medium density (e.g. ImageNet, Places)*

# Outline

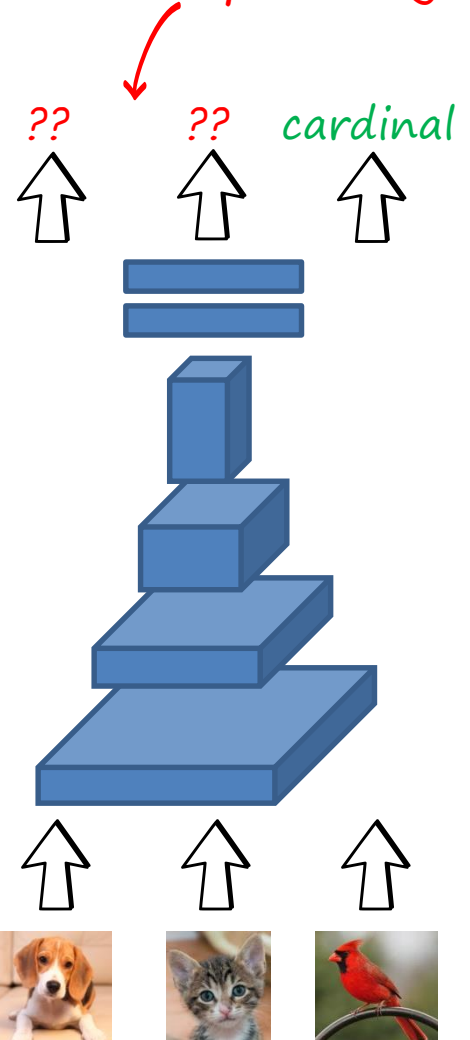
- Introduction
- Transferring GANs (ECCV 2018)
- **Rotated elastic weight consolidation (ICPR 2018)**
- Memory Replay GANs (NIPS 2018)
- Mix and match networks (CVPR 2018)

# When are neural networks good?

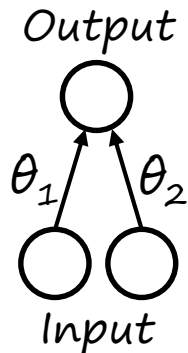


# Sequential learning

*Catastrophic forgetting*

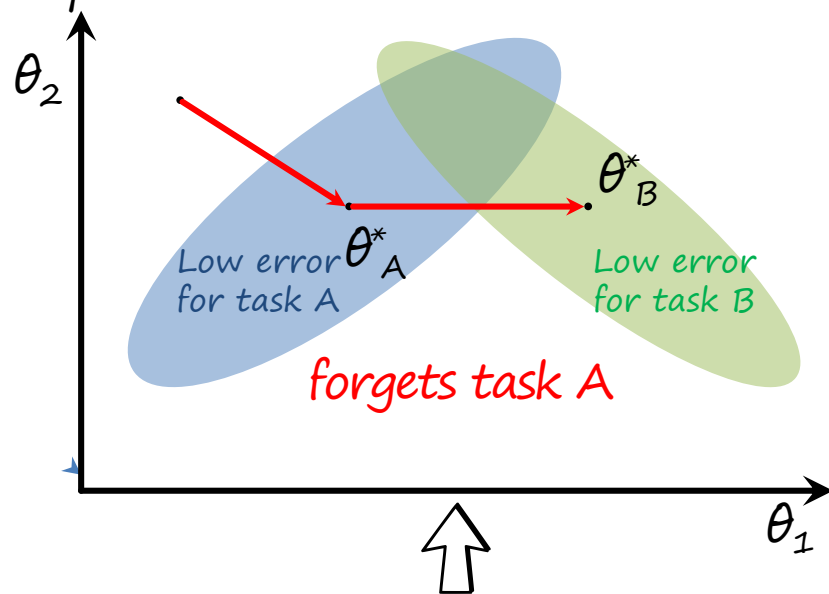


# Catastrophic interference and forgetting



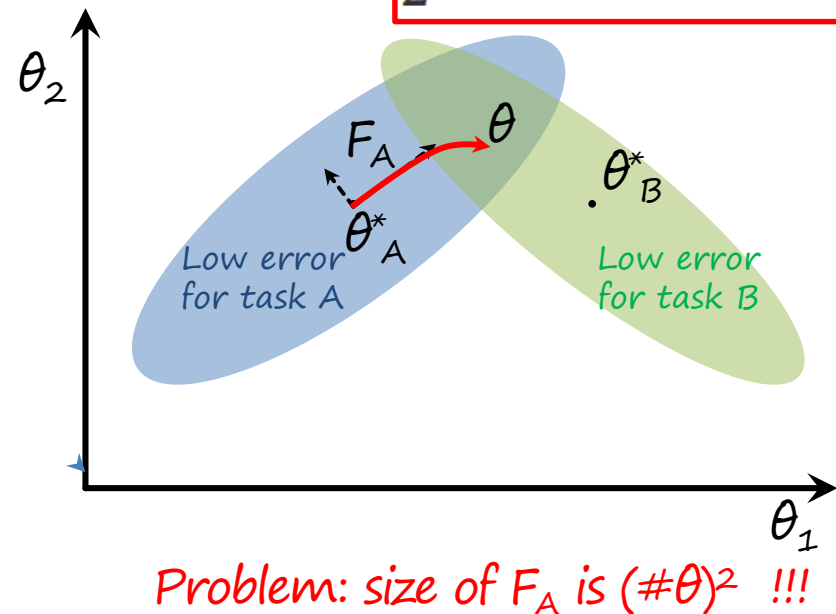
Training/fine tuning

$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta)$$



Elastic weight consolidation (EWC)

$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \frac{\lambda}{2} (\theta - \theta_A^*)^\top F_A (\theta - \theta_A^*)$$



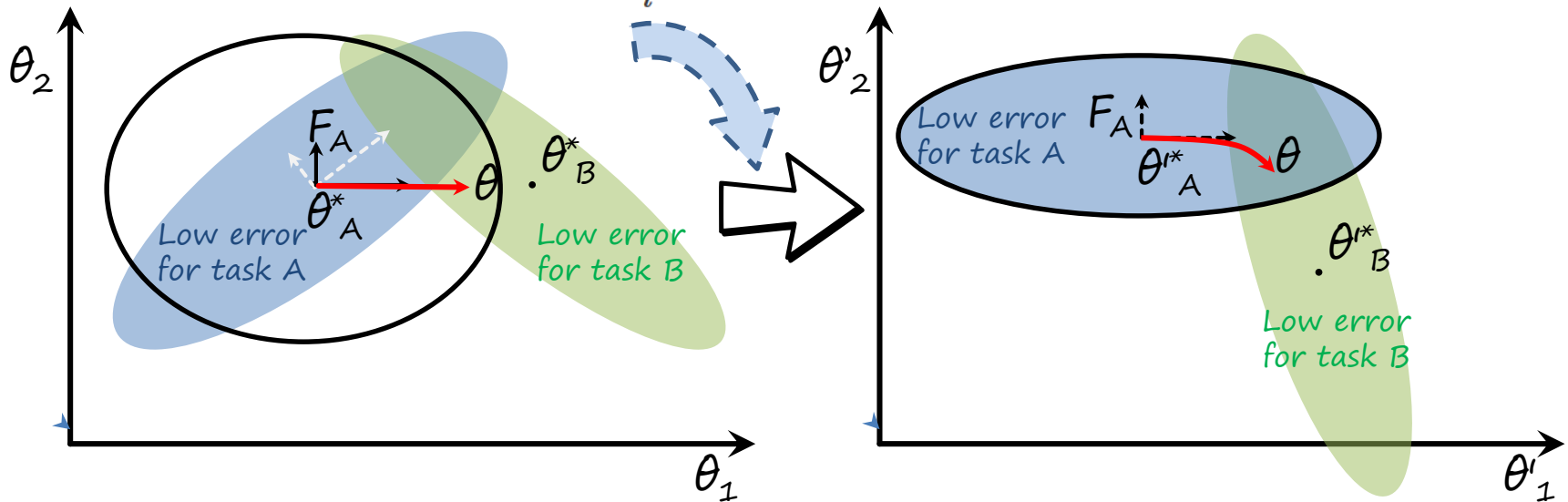
# Rotated elastic weight consolidation

Size of the diagonal of  $F_A$  is  $\#\theta$

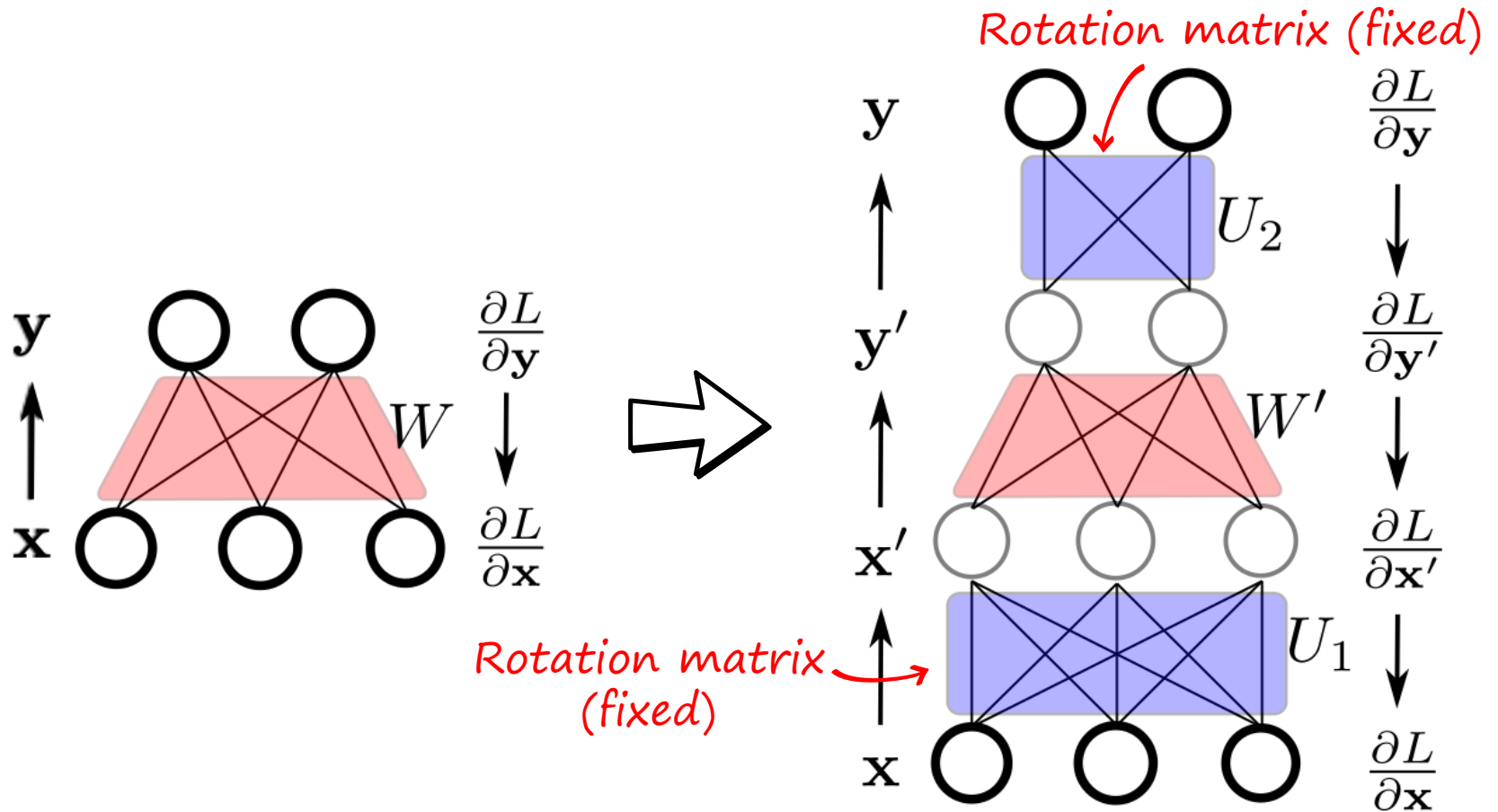
Elastic weight consolidation in practice  
(with diagonal approx. of  $F_A$ )

Rotated elastic weight consolidation (R-EWC)

$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \frac{\lambda}{2} \sum_i (F_A)_i (\theta_i - (\theta_A^*)_i)^2$$



# Rotating fully connected layers



# Computing the rotations

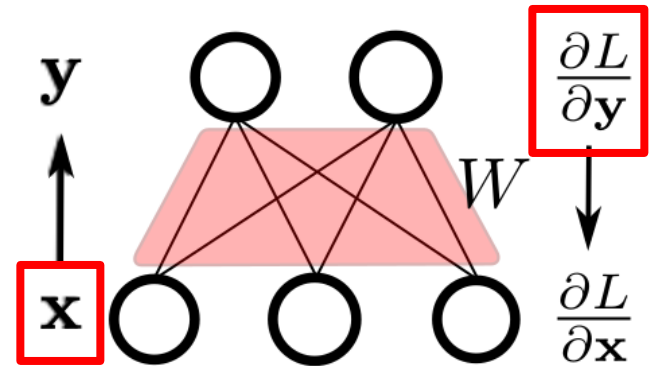
$$F_W = \mathbb{E}_{p \sim \pi} \left[ \left( \frac{\partial L}{\partial \mathbf{y}} \right) \mathbf{x} \mathbf{x}^\top \left( \frac{\partial L}{\partial \mathbf{y}} \right)^\top \right]$$

*Assuming  $\mathbf{x}$  and  $\delta L / \delta \mathbf{y}$  independent*

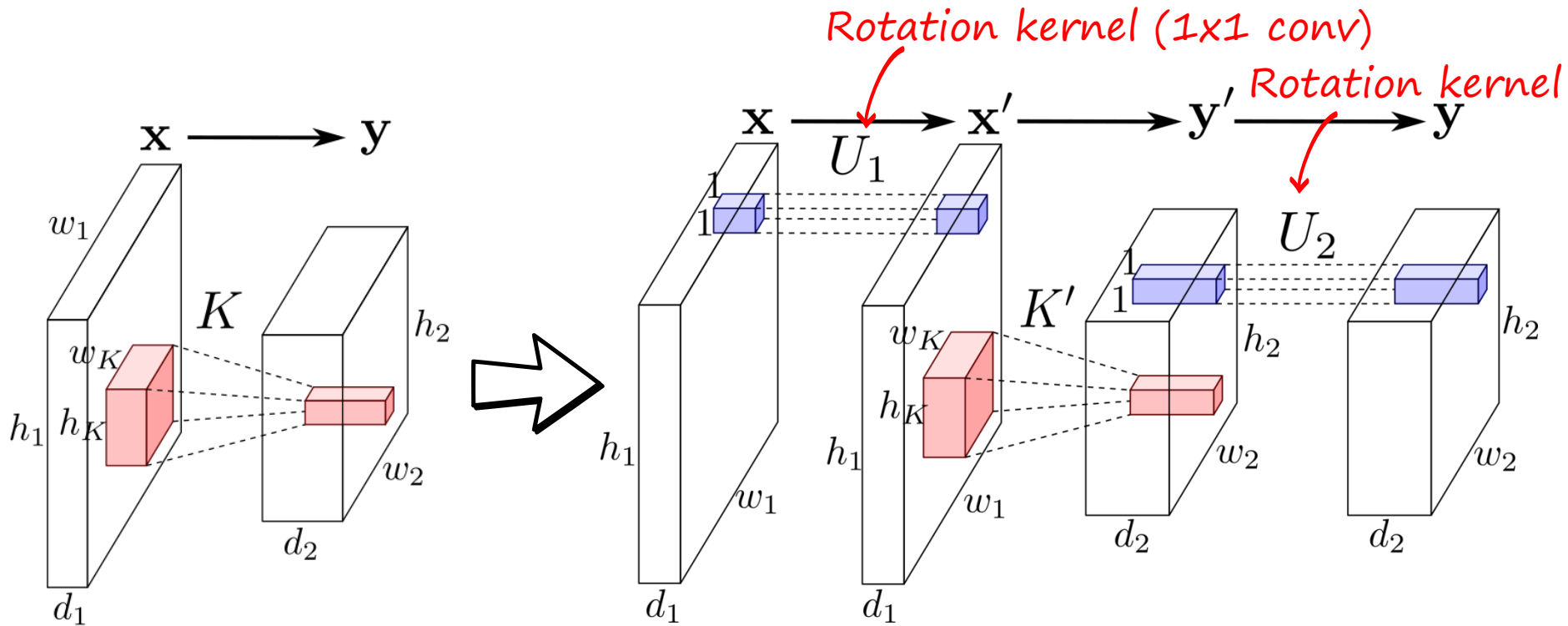
$$F_W \approx \mathbb{E}_{\substack{x \sim \pi \\ y \sim p}} \left[ \left( \frac{\partial L}{\partial \mathbf{y}} \right) \left( \frac{\partial L}{\partial \mathbf{y}} \right)^\top \right] \mathbb{E}_{x \sim \pi} [\mathbf{x} \mathbf{x}^\top]$$

$$\begin{aligned} \mathbb{E}_{x \sim \pi} [\mathbf{x} \mathbf{x}^\top] &= U_1 S_1 V_1^\top \\ \mathbb{E}_{\substack{x \sim \pi \\ y \sim p}} \left[ \left( \frac{\partial L}{\partial \mathbf{y}} \right) \left( \frac{\partial L}{\partial \mathbf{y}} \right)^\top \right] &= U_2 S_2 V_2^\top \end{aligned}$$

*Using SVD*

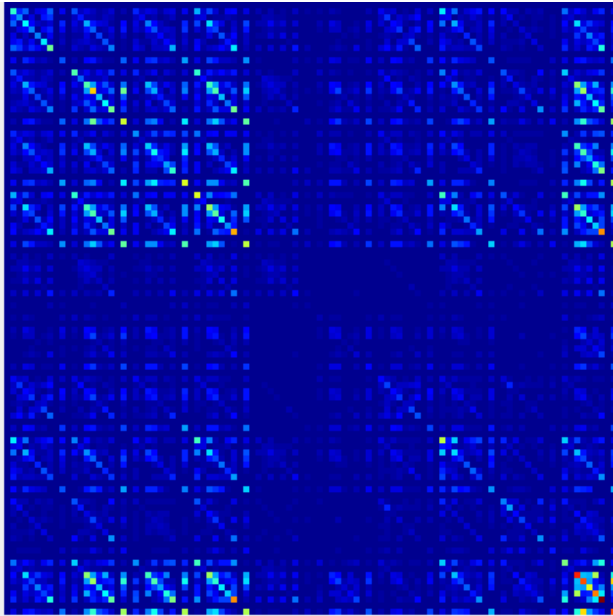


# Rotating convolutional layers



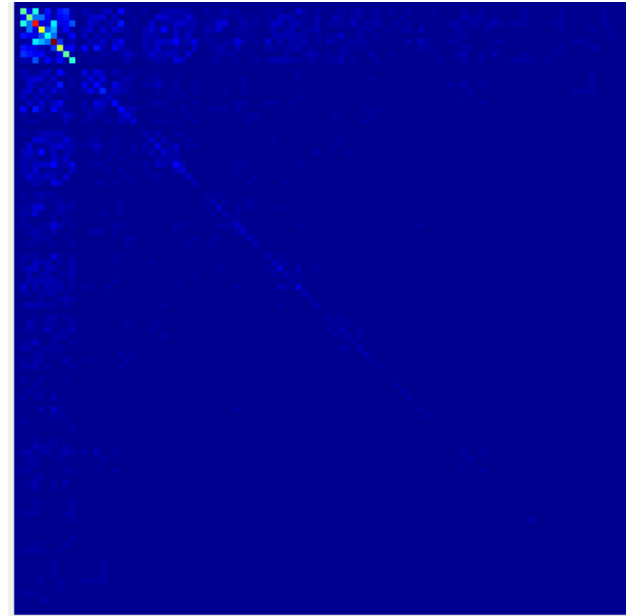
# Fisher Information matrix

*No rotation (i.e. EWC)*



*Energy in the diagonal: 40%*

*After rotation (i.e. R-EWC)*



*Energy in the diagonal: 74%*

# Experimental results (2 tasks)

- MNIST dataset. Two tasks: 0-4 and 5-9

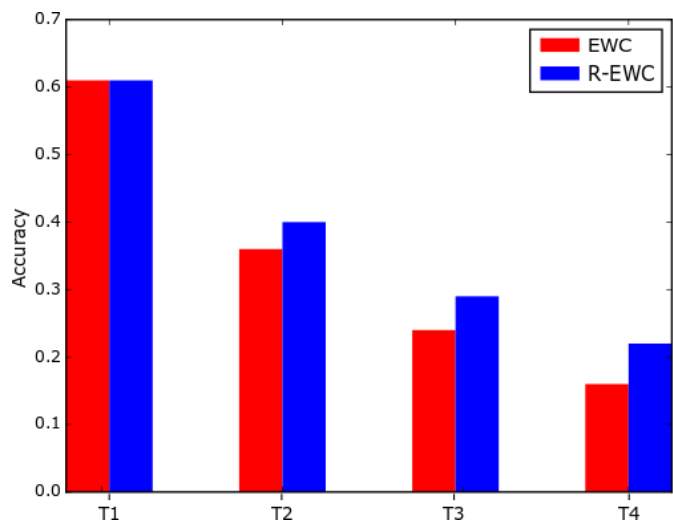
	$\lambda = 1$		$\lambda = 10$		$\lambda = 100$		$\lambda = 1000$		$\lambda = 10000$	
	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2	Task 1	Task 2
FT	6.1	97.6	6.1	97.6	6.1	97.6	6.1	97.6	6.1	97.6
EWC [5]	66.8	90.9	75.3	95.6	<b>85.8</b>	<b>92.8</b>	78.4	93.7	81.0	88.8
R-EWC - conv only	62.7	89.2	67.5	96.1	80.4	91.4	<b>84.7</b>	<b>93.1</b>	75.5	93.7
R-EWC - fc only	78.9	95.3	79.0	95.8	87.4	93.5	93.0	82.3	<b>94.3</b>	<b>88.0</b>
R-EWC - all	77.2	96.7	<b>91.7</b>	<b>91.2</b>	86.9	95.9	96.3	81.1	92.1	86.0
R-EWC - all no last	71.5	91.8	84.9	97.0	<b>91.6</b>	<b>94.5</b>	94.6	88.4	97.9	79.4

- Several datasets

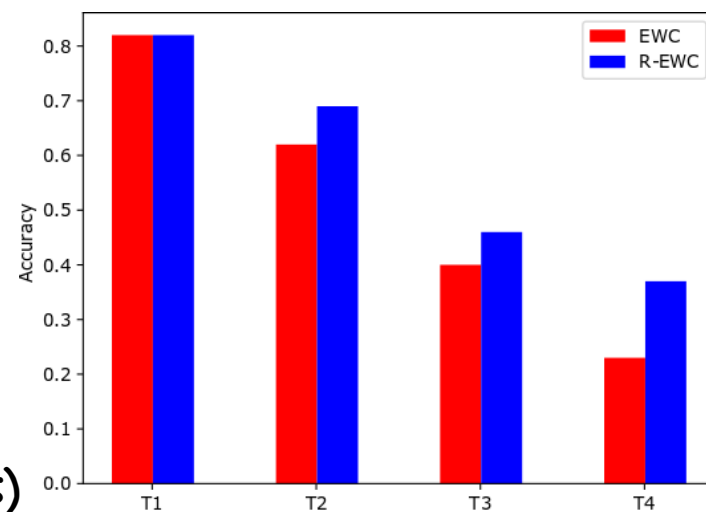
	EWC [5] (T1 / T2)	R-EWC (T1 / T2)
MNIST	89.3 (85.8 / 92.8)	<b>93.1</b> (91.6 / 94.5)
CIFAR-100	37.5 (23.5 / 51.5)	<b>42.5</b> (30.2 / 54.7)
CUB-200 Birds	45.3 (42.3 / 48.6)	<b>48.4</b> (53.3 / 45.2)
Stanford-40 Actions	50.4 (44.3 / 58.4)	<b>52.5</b> (52.3 / 52.6)

# Experimental results (4 tasks)

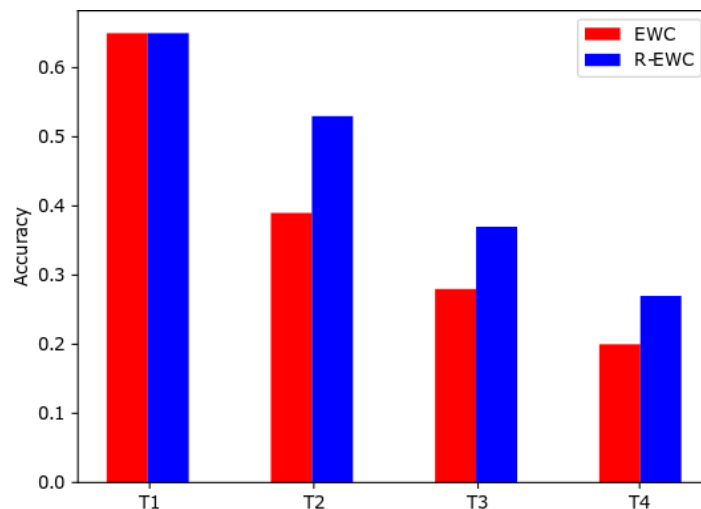
*CIFAR 100*



*Stanford Actions*



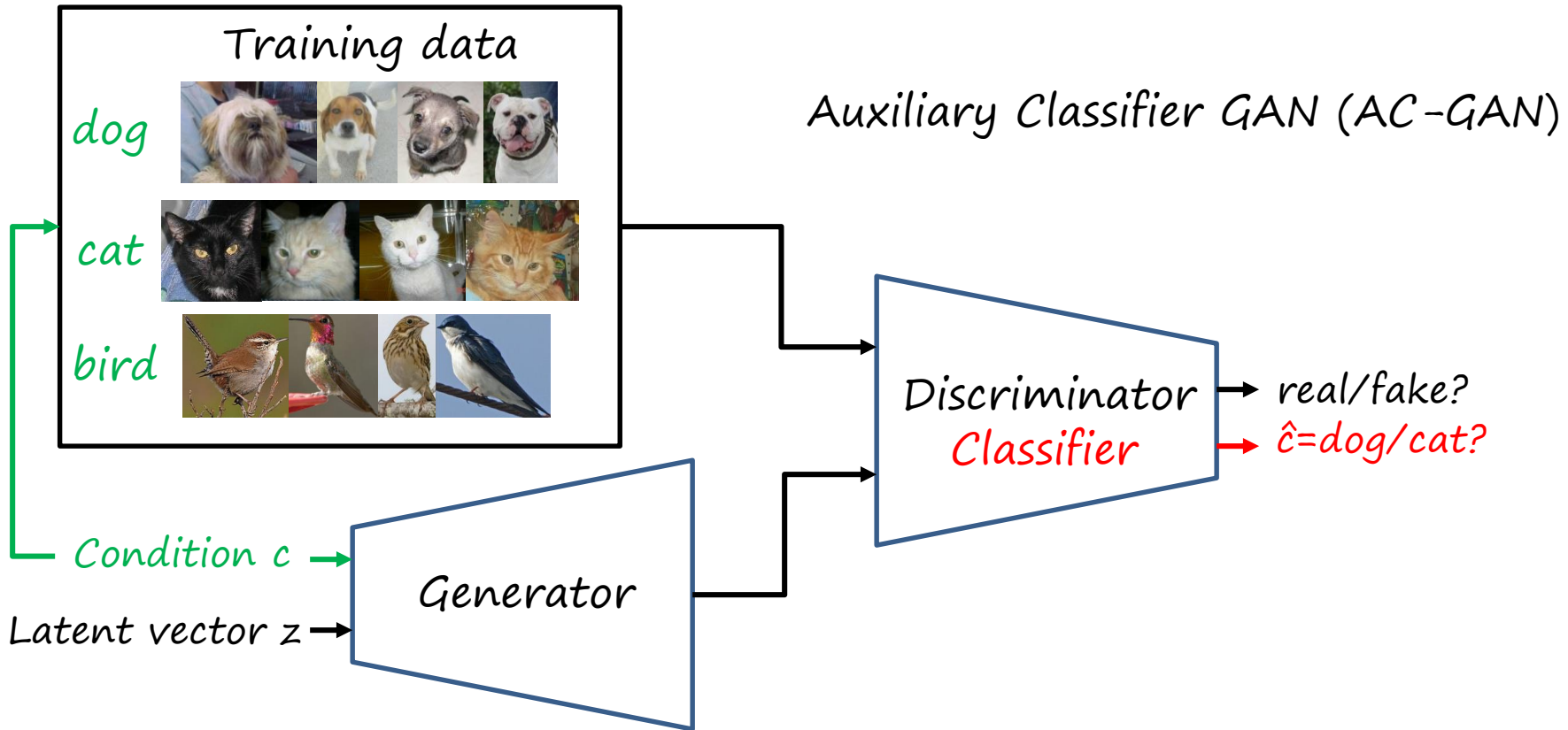
*CUB 200 (birds)*



# Outline

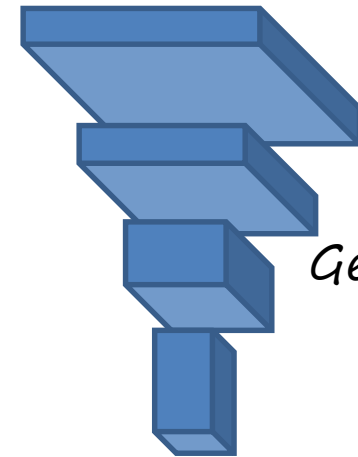
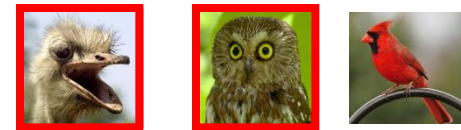
- Introduction
- Transferring GANs (ECCV 2018)
- Rotated elastic weight consolidation (ICPR 2018)
- **Memory Replay GANs (NIPS 2018)**
- Mix and match networks (CVPR 2018)

# Joint training (non-sequential)

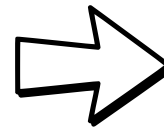


# Sequential learning for image generation

*Catastrophic forgetting*



Generator



$c=\text{dog}$     $c=\text{cat}$     $c=\text{bird}$   
 $z=0.643$     $z=0.453$     $z=0.132$



# Sequential fine tuning and forgetting

LSUN 4 categories (4 tasks)

Task 1 Task 2 Task 3 Task 4

$c=bedroom$



$c=kitchen$



$c=church$



$c=tower$



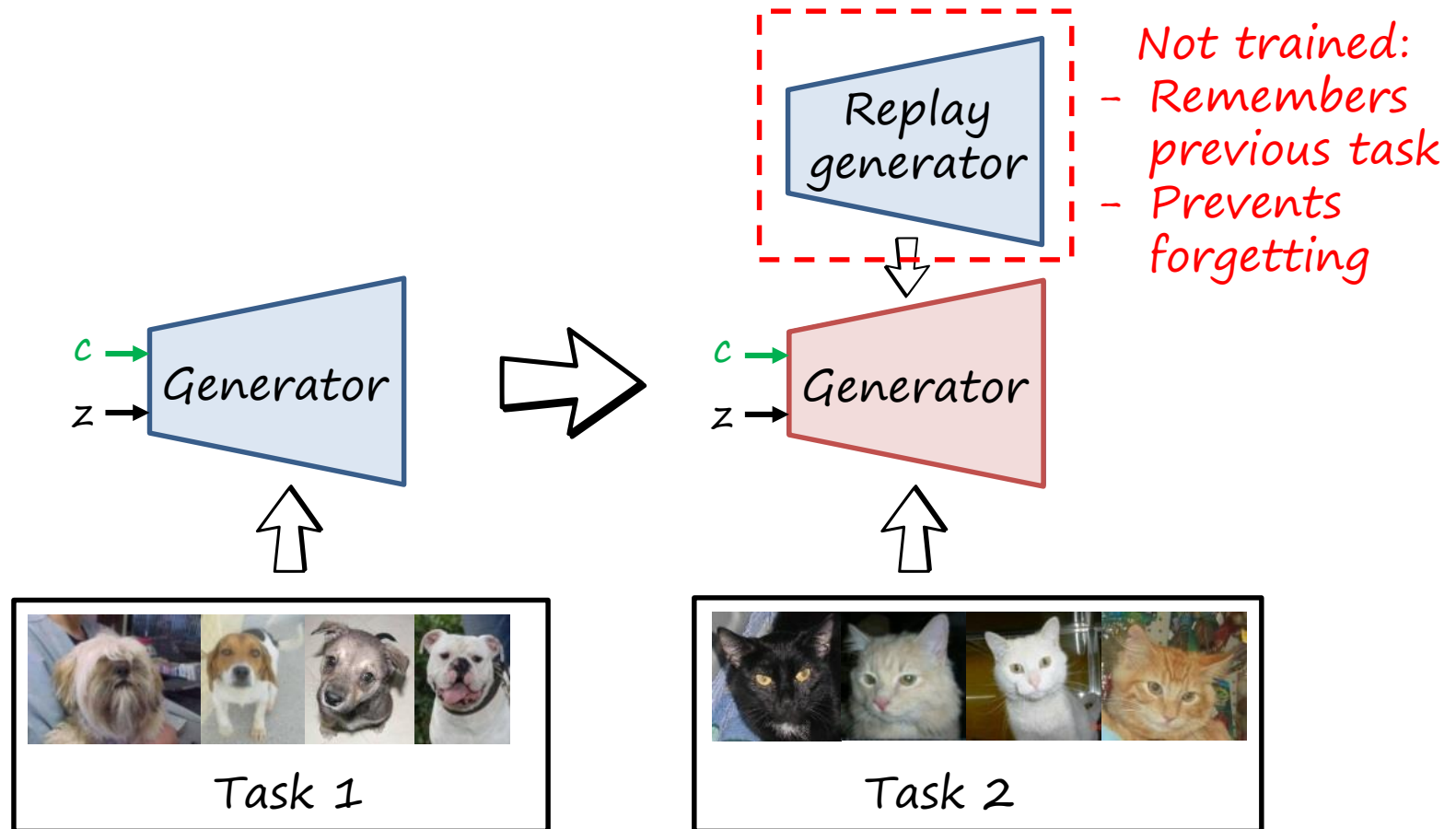
MNIST 10 categories (10 tasks)

$Z_1 Z_2 Z_3 Z_4$

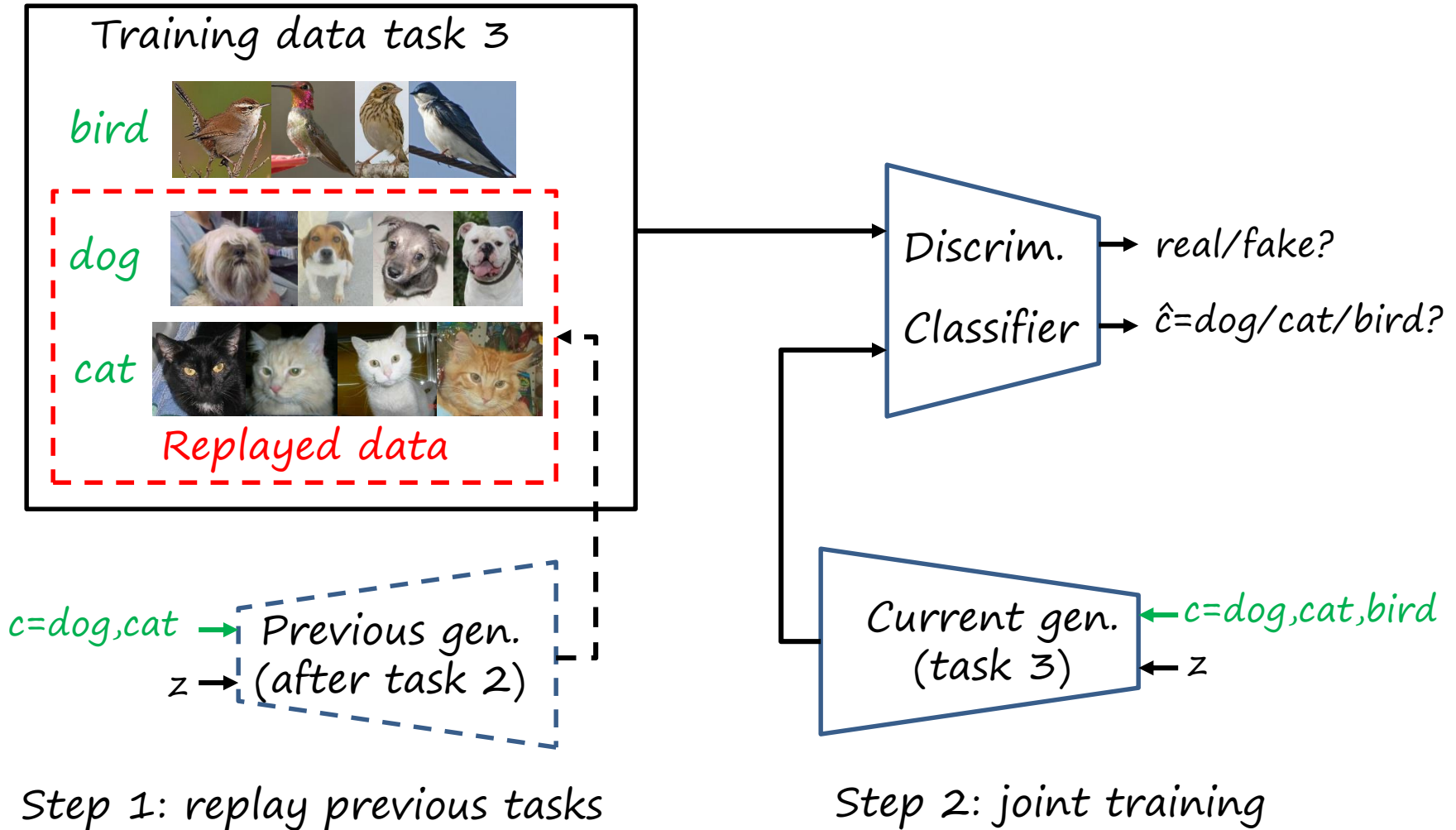
$c=0$	9	9	9	9
$c=1$	9	9	9	9
$c=2$	9	9	9	9
$c=3$	9	9	9	9
$c=4$	9	9	9	9
$c=5$	9	9	9	9
$c=6$	9	9	9	9
$c=7$	9	9	9	9
$c=8$	9	9	9	9
$c=9$	9	9	9	9

After task 10

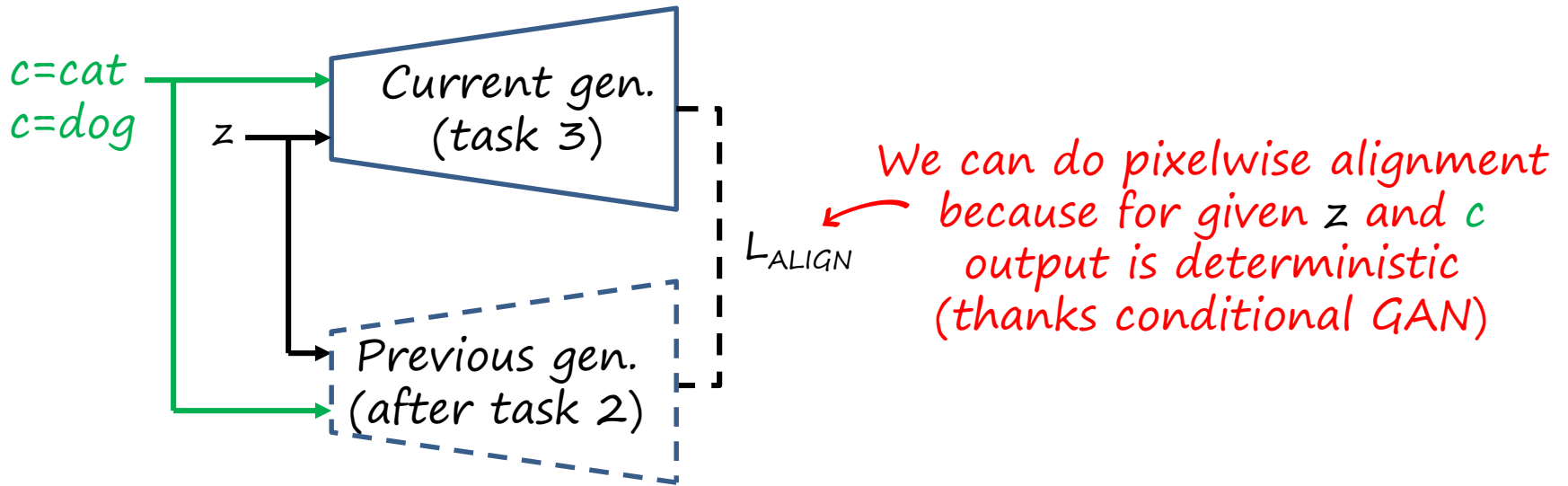
# Memory Replay GANs (MeRGANs)



# MeRGAN-JTR: joint training w/ replay

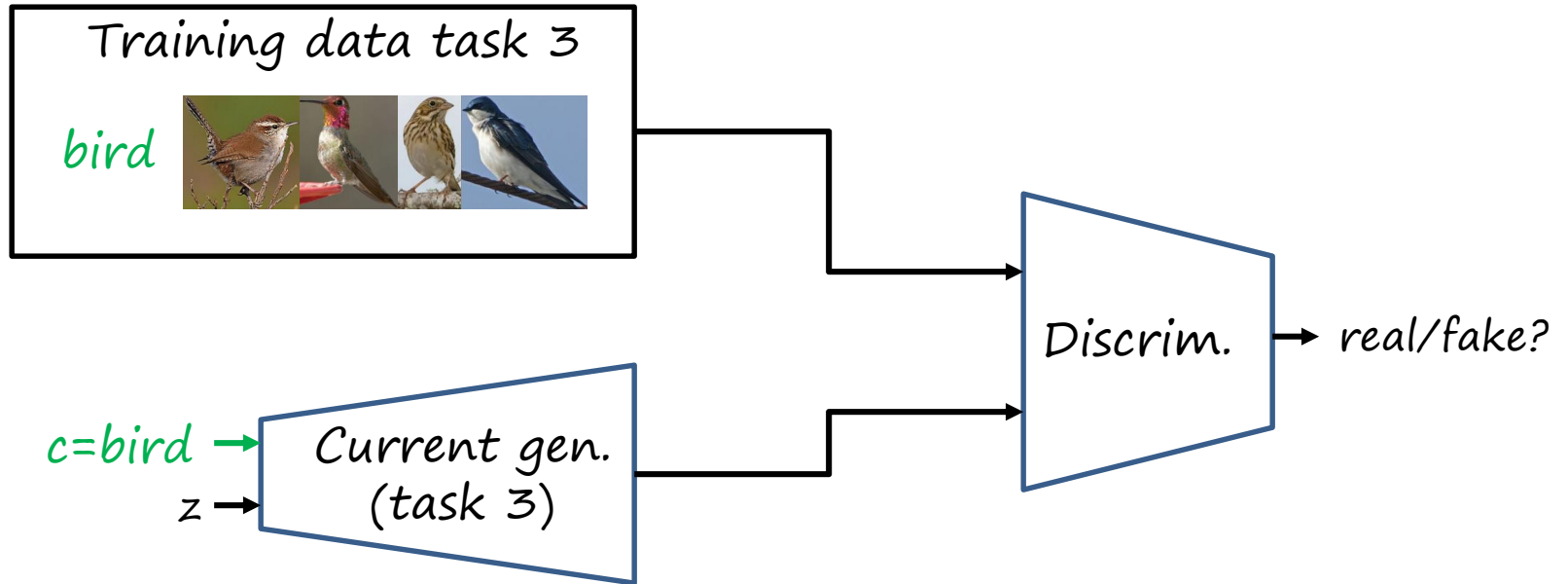


# MeRGAN-RA: replay alignment



Step 1: replay previous tasks and align

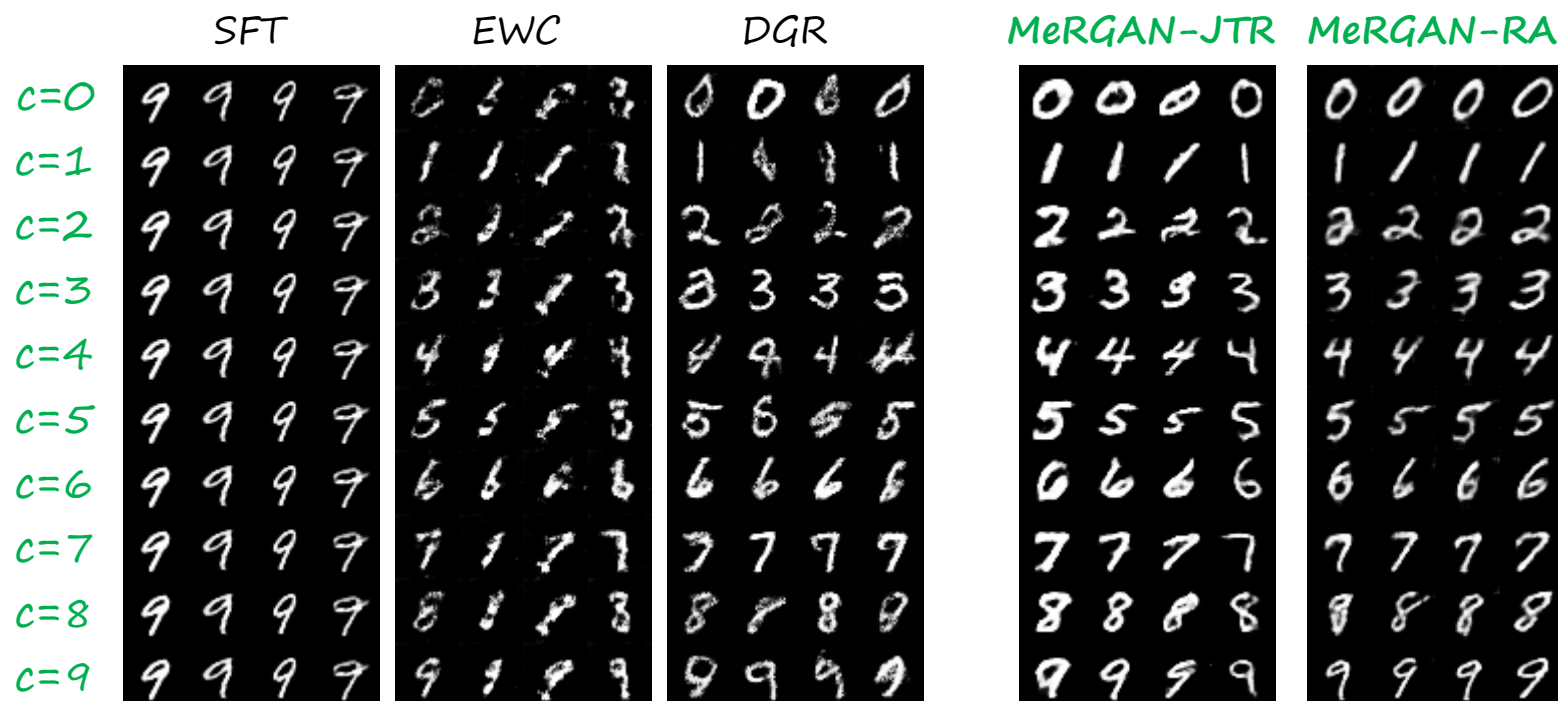
# MeRGAN-RA: replay alignment



Step 2: learning new task

# Digit generation (MNIST)

10 tasks (10 categories) 3 layers generator



SFT: Sequential fine tuning EWC: GAN with EWC [arxiv17] DGR: Deep generative replay [NIPS17]

# Scene generation (LSUN)

4 tasks (4 categories) 18-layer ResNet generator

Sequential fine tuning  
Task 1 Task 2 Task 3 Task 4



*Different bedrooms!*

MeRGAN-JTR  
Task 1 Task 2 Task 3 Task 4



*Remembers the category*

*Same bedroom!*

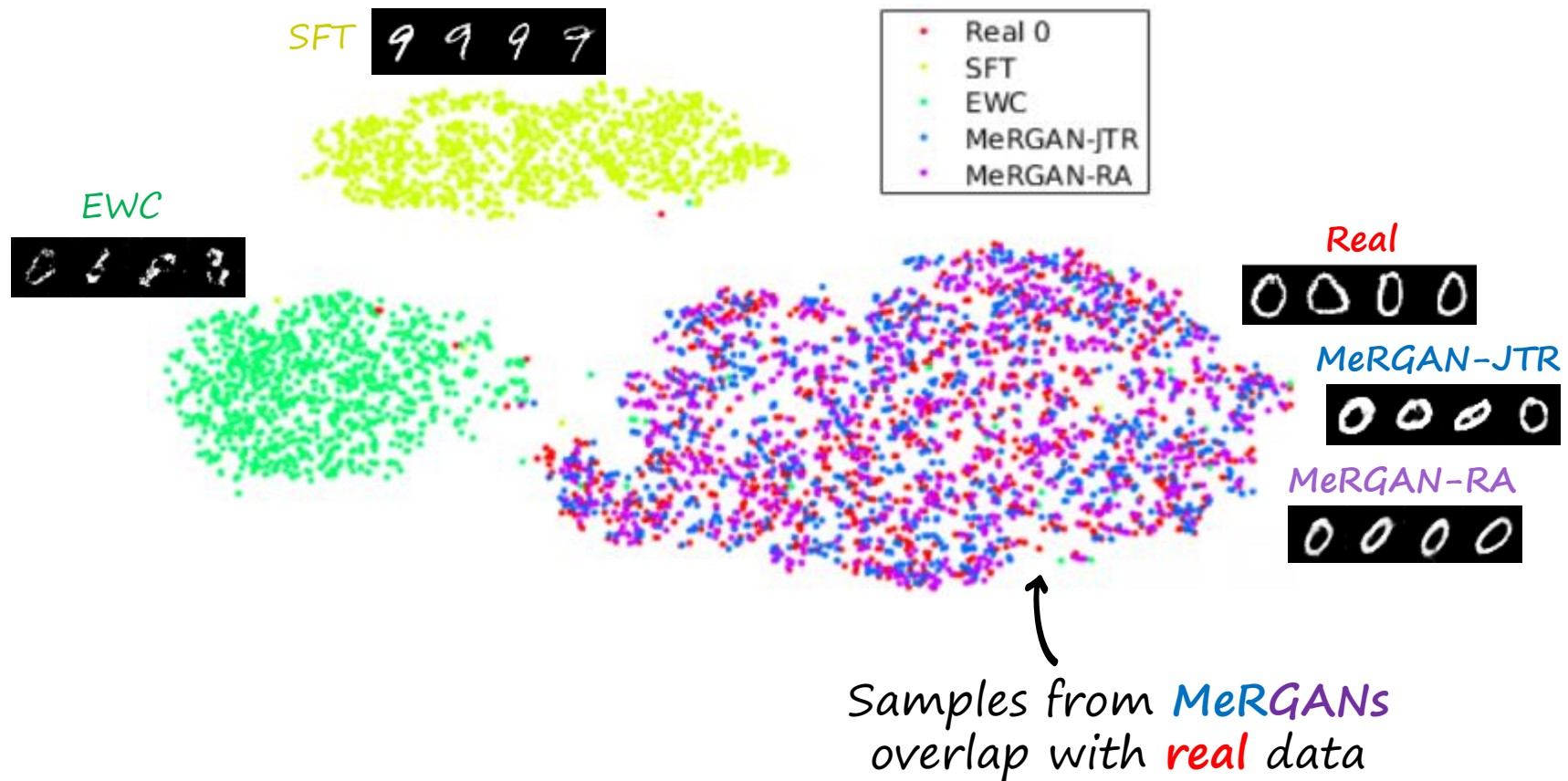
MeRGAN-RA  
Task 1 Task 2 Task 3 Task 4



*Remembers the instance*

# t-SNE visualizations (MNIST)

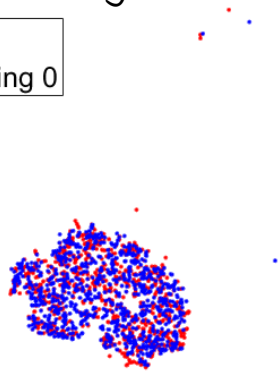
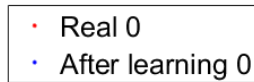
Generating digit 0 (i.e. first task) after learning 10 tasks



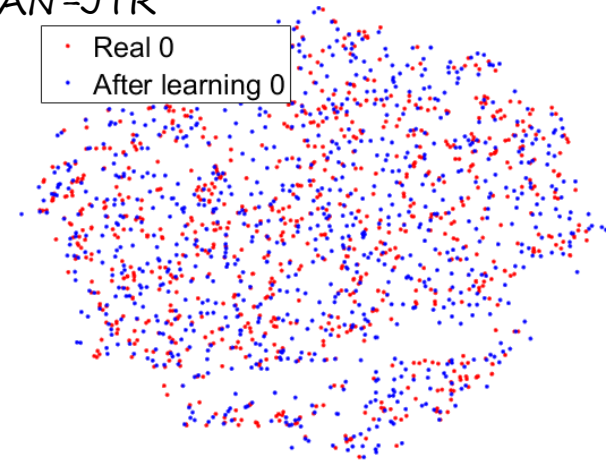
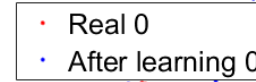
# Learning and forgetting in t-SNE (MNIST)

Generating digit 0 (i.e. task 1)

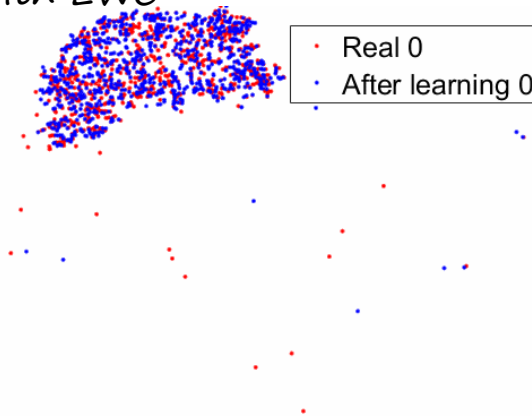
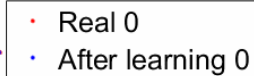
*Sequential fine tuning*



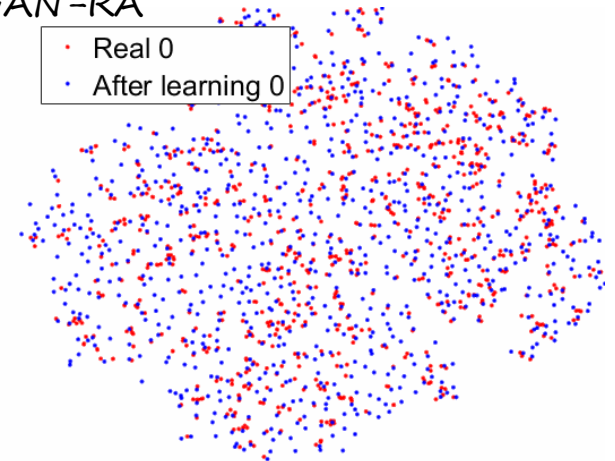
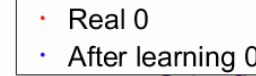
*MeRGAN-JTR*



*GAN with EWC*



*MeRGAN-RA*



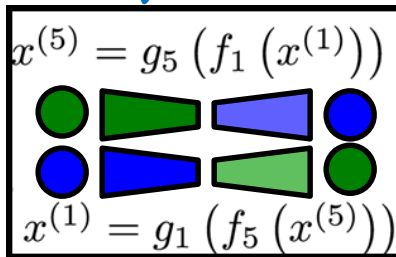
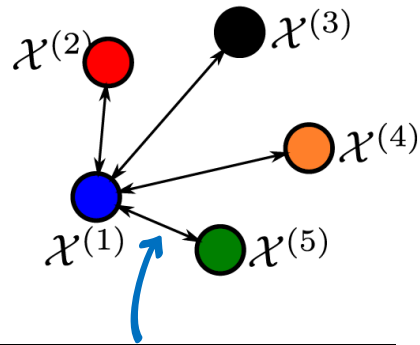
# Outline

- Introduction
- Transferring GANs (ECCV 2018)
- Rotated elastic weight consolidation (ICPR 2018)
- Memory Replay GANs (NIPS 2018)
- **Mix and match networks (CVPR 2018)**

# Unseen image-to-image translations

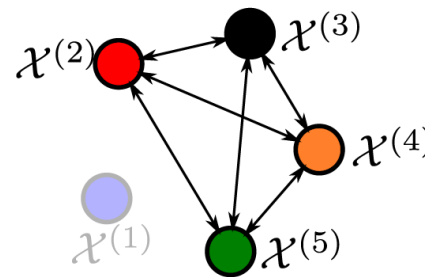
Only these translations  
are trained (seen)

Train

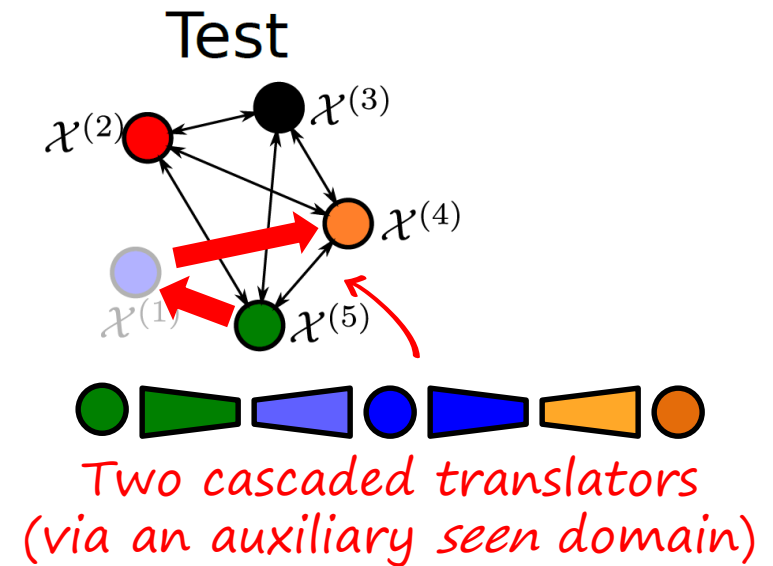
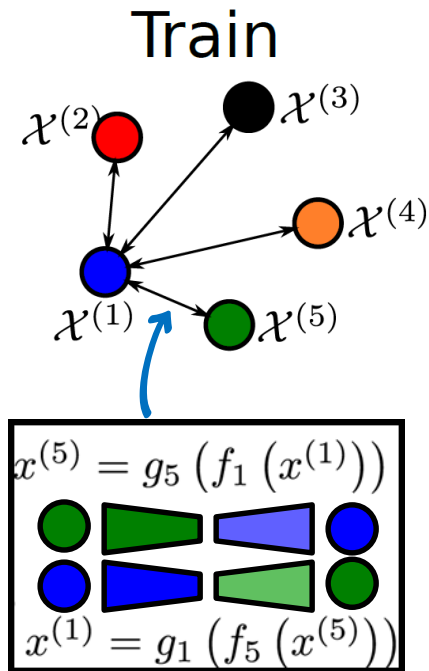


Evaluate on these unseen  
translations

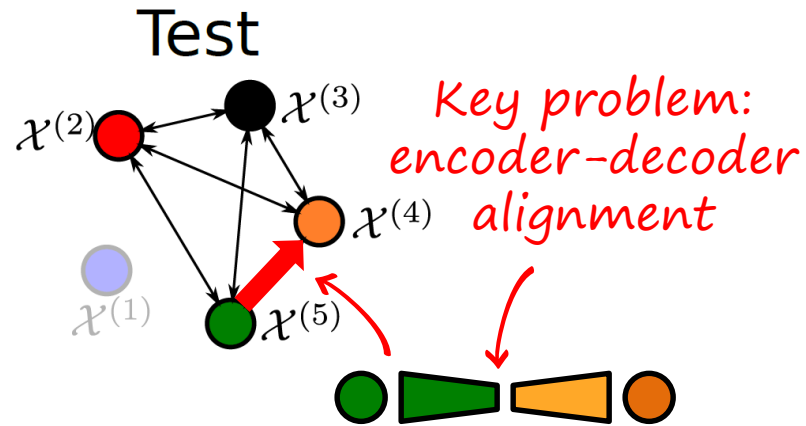
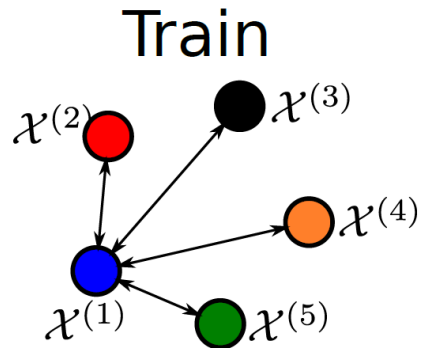
Test



# Cascading image-to-image translators



# Mix and match networks



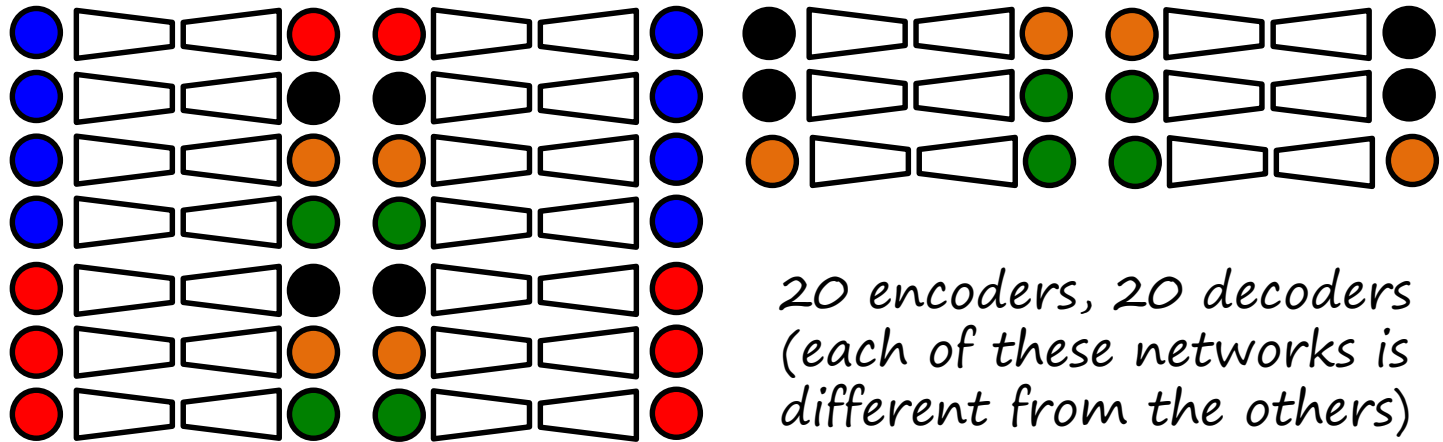
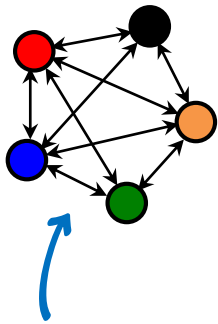
Mix&match encoder-decoders  
(they haven't seen each other  
during training)

# Application: many-to-many translations

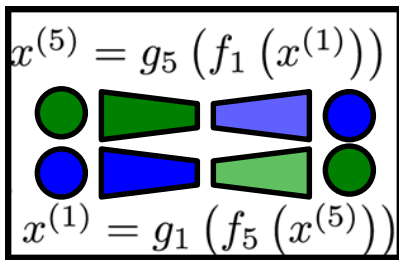
We could train all possible translators

Problems:

- No sharing
- Poor scalability: number of networks  $O(N^2)$



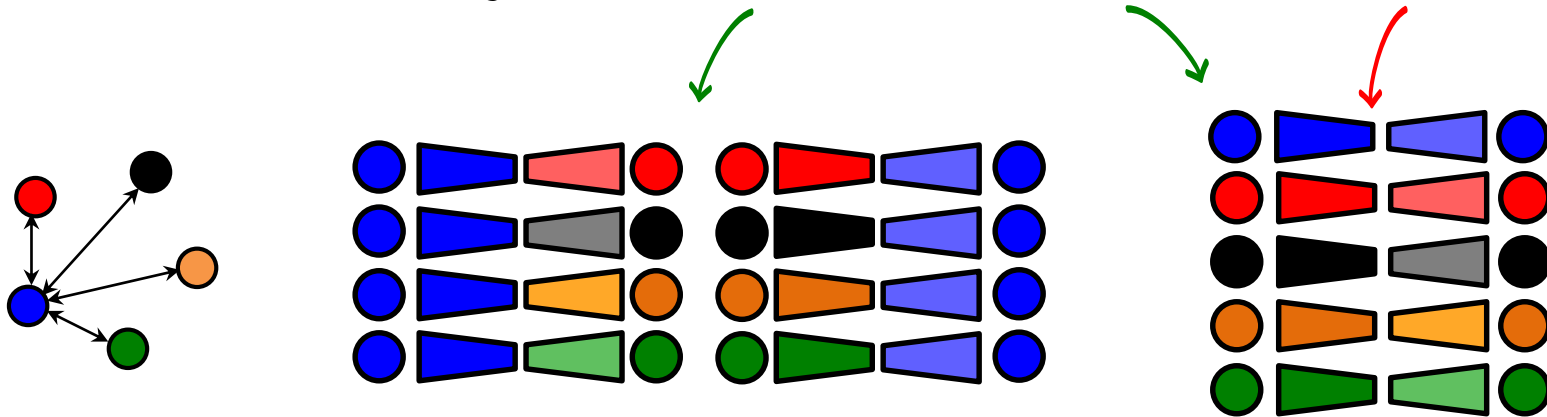
20 encoders, 20 decoders  
(each of these networks is different from the others)



# Mix and match networks

Unseen encoder-decoder alignment

- Latent representation should be **domain-independent**
- Achieved using **shared encoder/decoders** and **autoencoders**



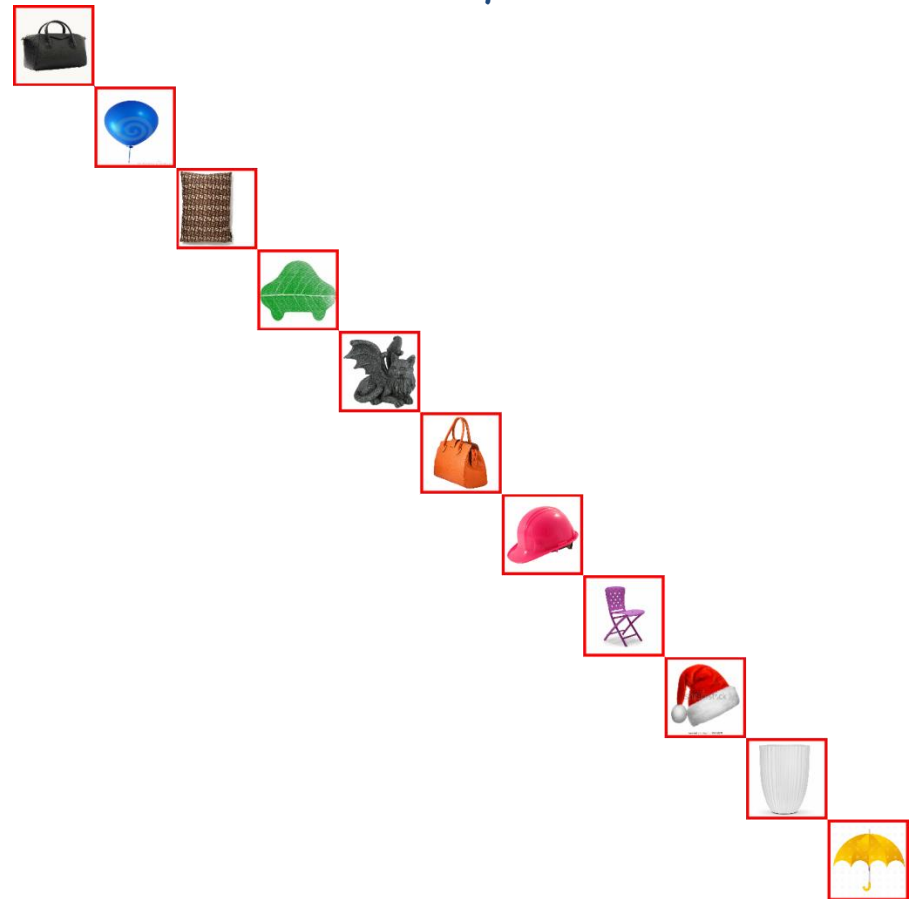
5 encoders, 5 decoders

- Scalable: number of networks  $O(N)$

# Example: scalable recolorization

*Unpaired translation  
11 colors (i.e. 11 domains)*

*Input*





# Example: scalable recolorization

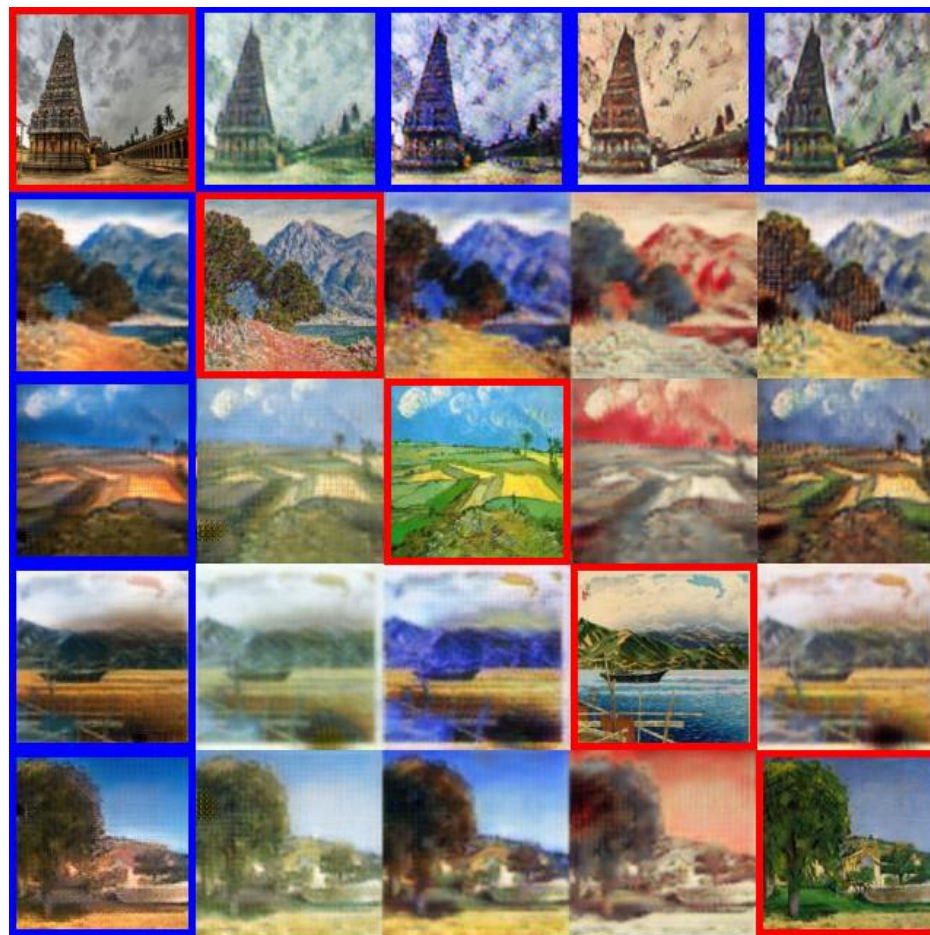
Unpaired translation  
11 colors (i.e. 11 domains)

Requires training 11  
encoders and 11 decoders

CycleGANs for all  
combinations would require  
55 encoders and 55 decoders



# Example: scalable style transfer



*Unpaired translation  
Five domains  
(photo, Monet, van Gogh,  
Ukiyo-e, Cezanne)*

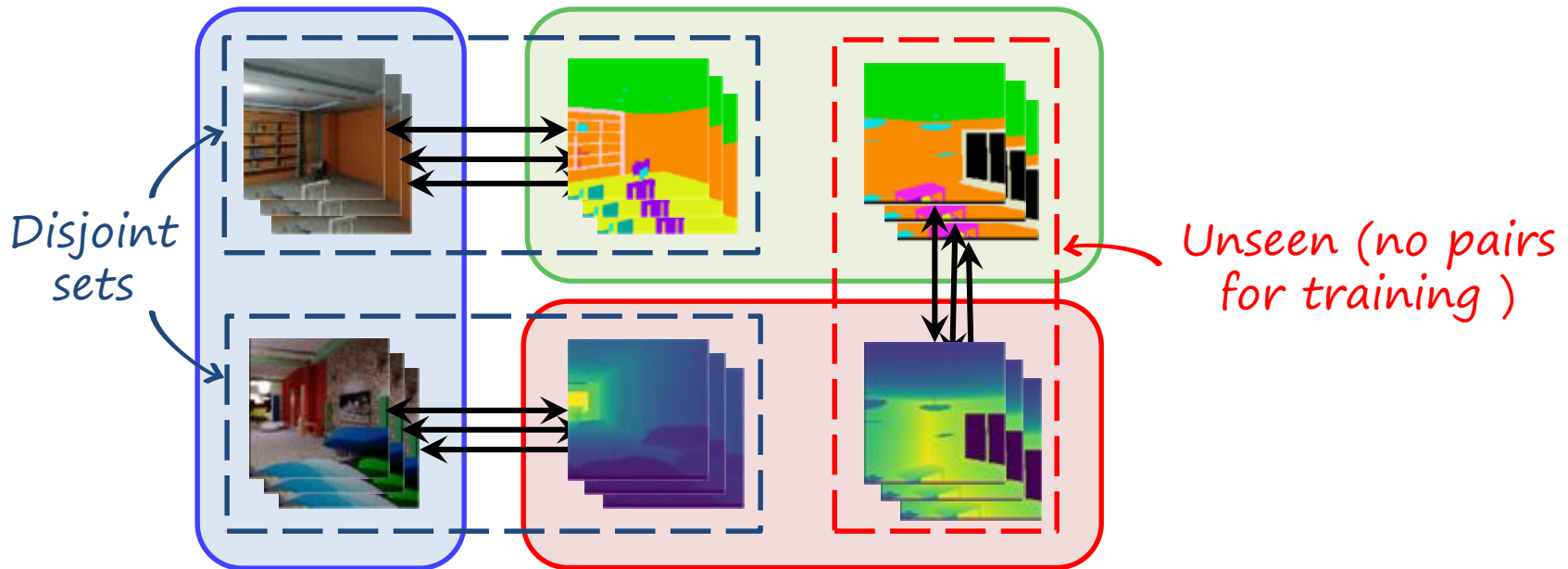
*(5 encoders and 5  
decoders)*

# Application: zero-pair translation

Cross-modal translation setting (RGB, segmentation and depth)

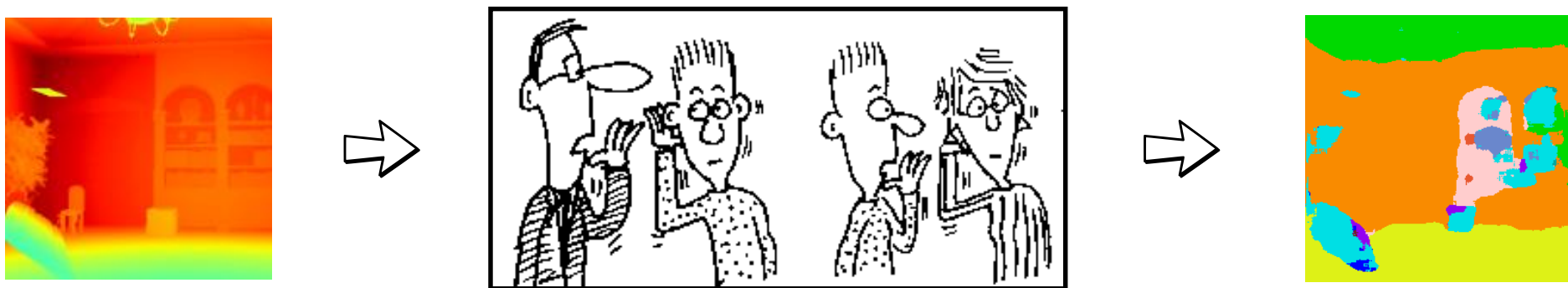
Paired data available for (RGB, segm.) and (RGB, depth)

Evaluate on the unseen zero-pair translations (depth, segm.)

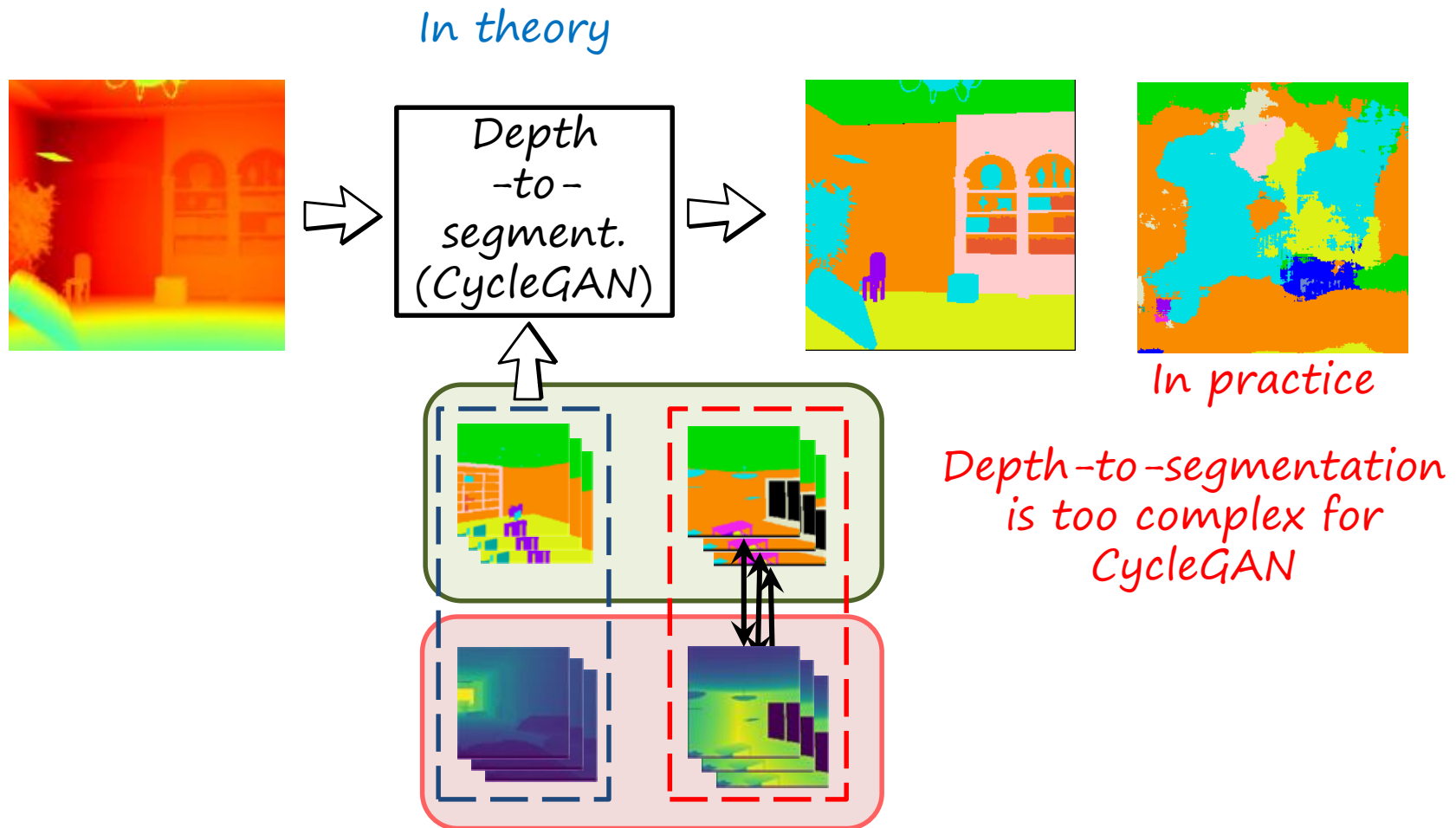


# Zero-pair translation with two cascaded pix2pix (paired translations)

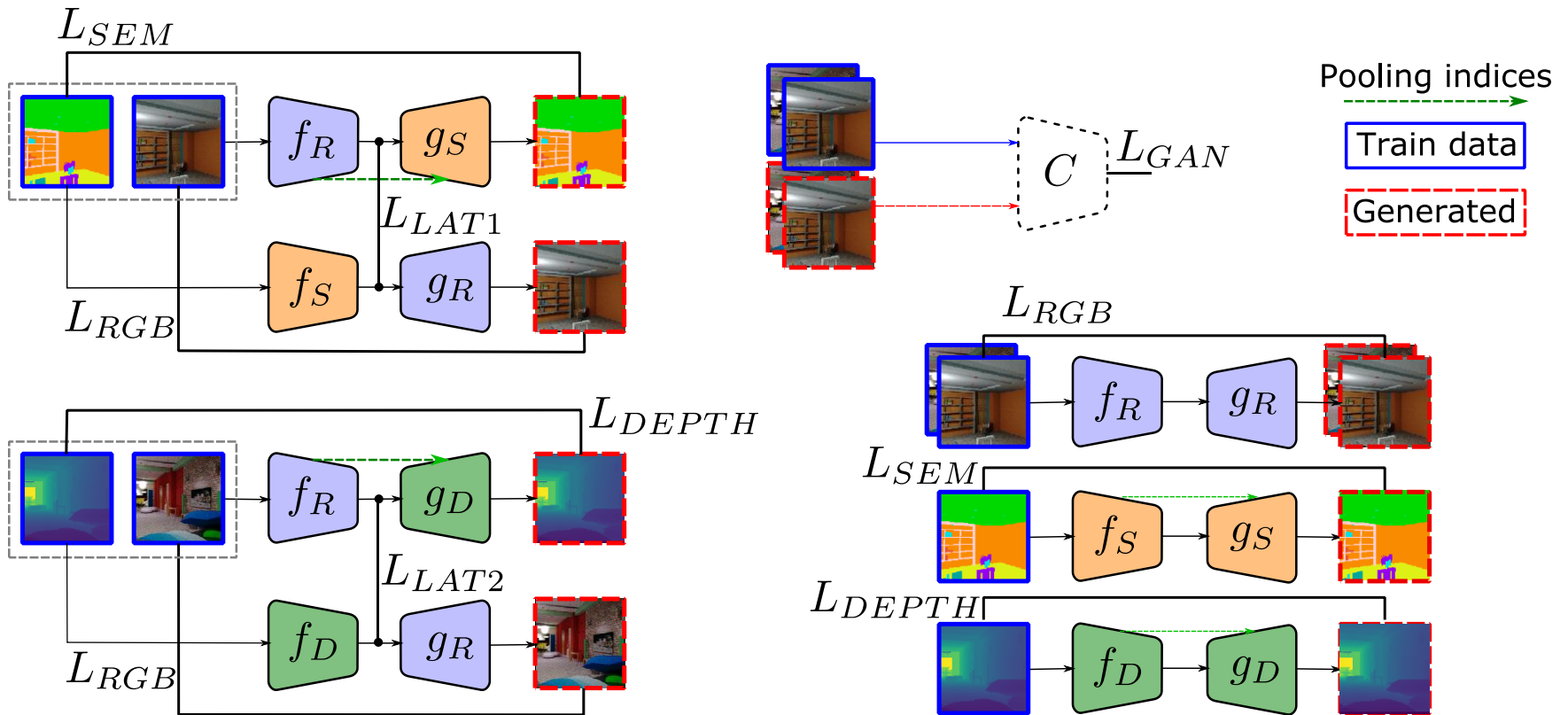
*In theory*



# Zero-pair translation with CycleGAN (unpaired translation)

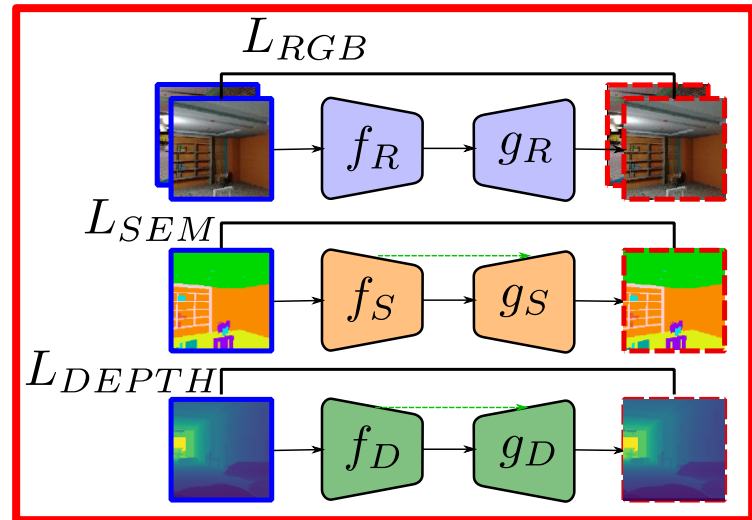
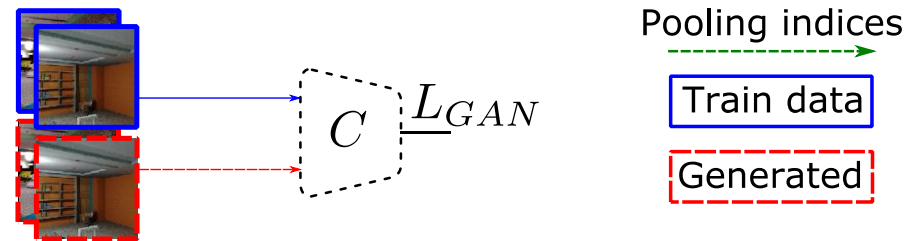
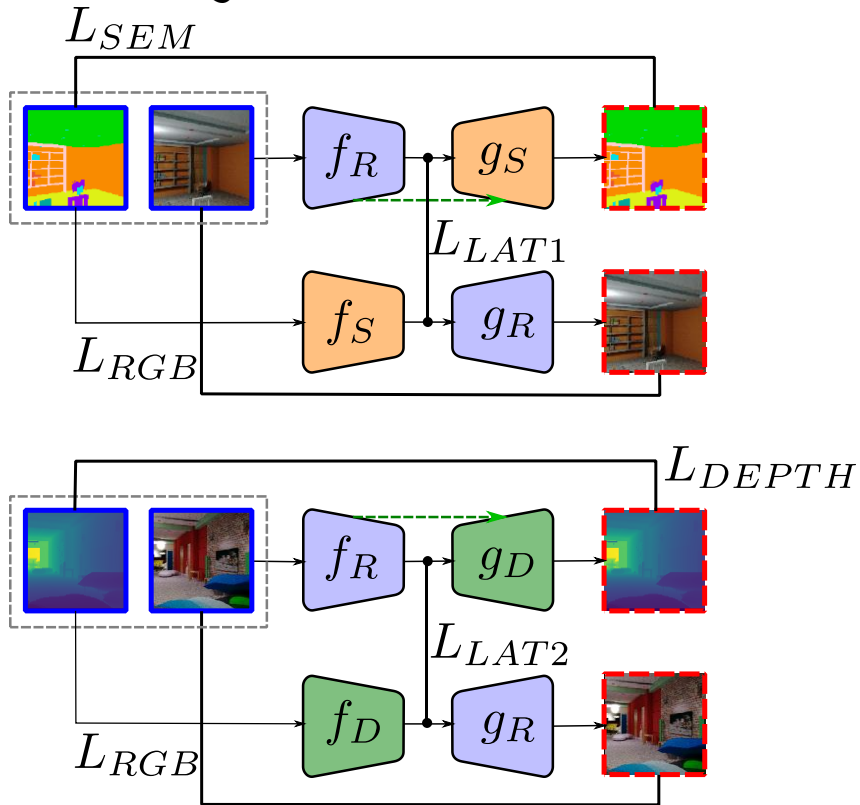


# Zero-pair translation with mix and match networks



# Zero-pair translation with mix and match networks

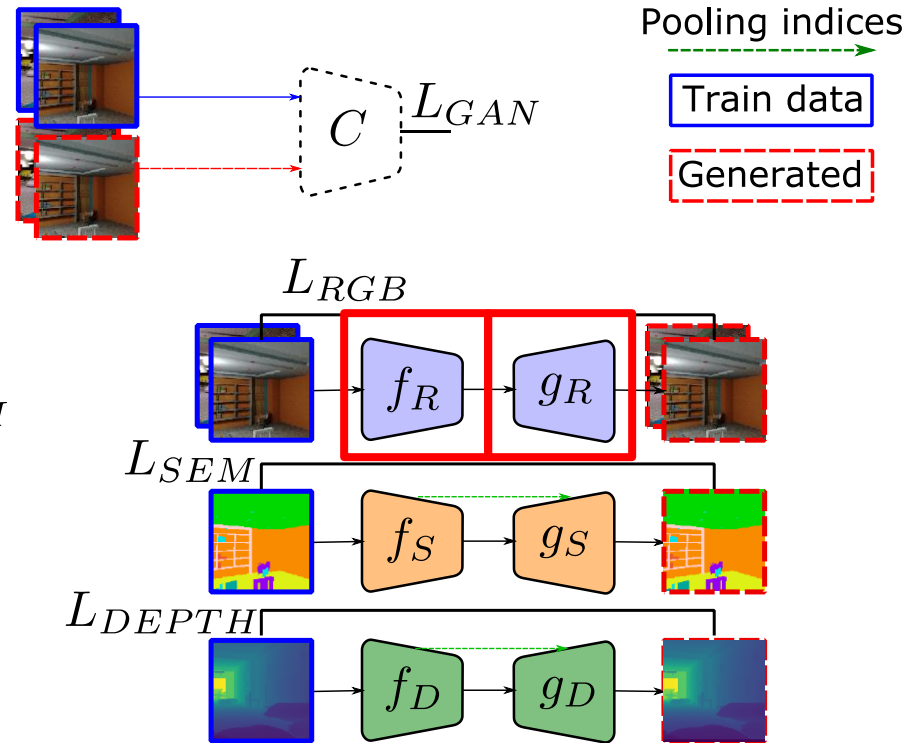
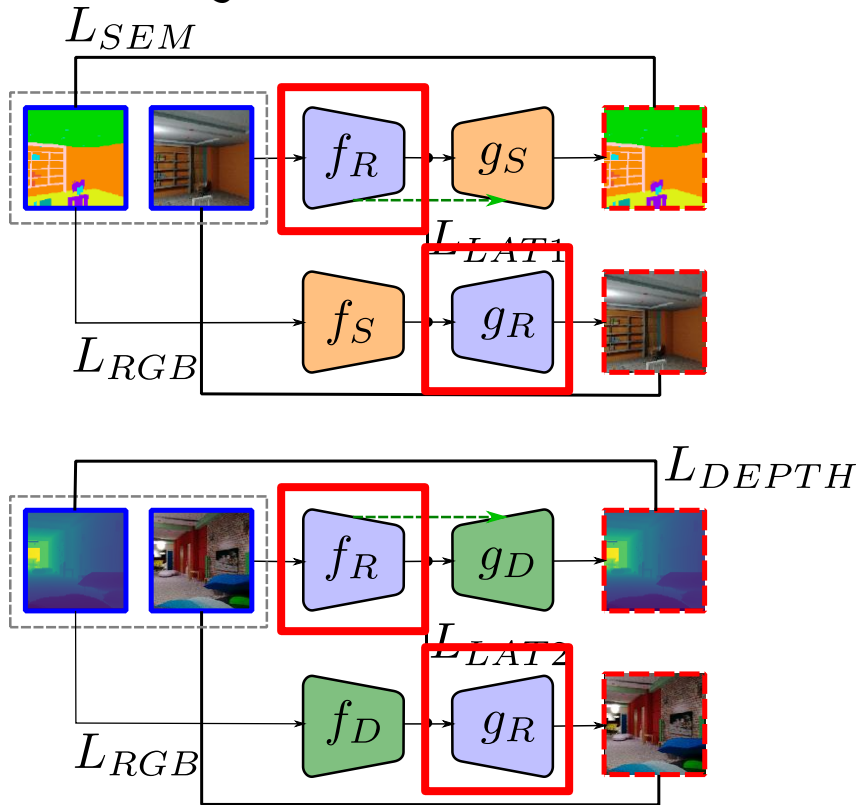
Training for encoder-decoder alignment: *Autoencoders*



# Zero-pair translation with mix and match networks

Training for encoder-decoder alignment:

Autoencoders  
*Shared encoder/decoders*

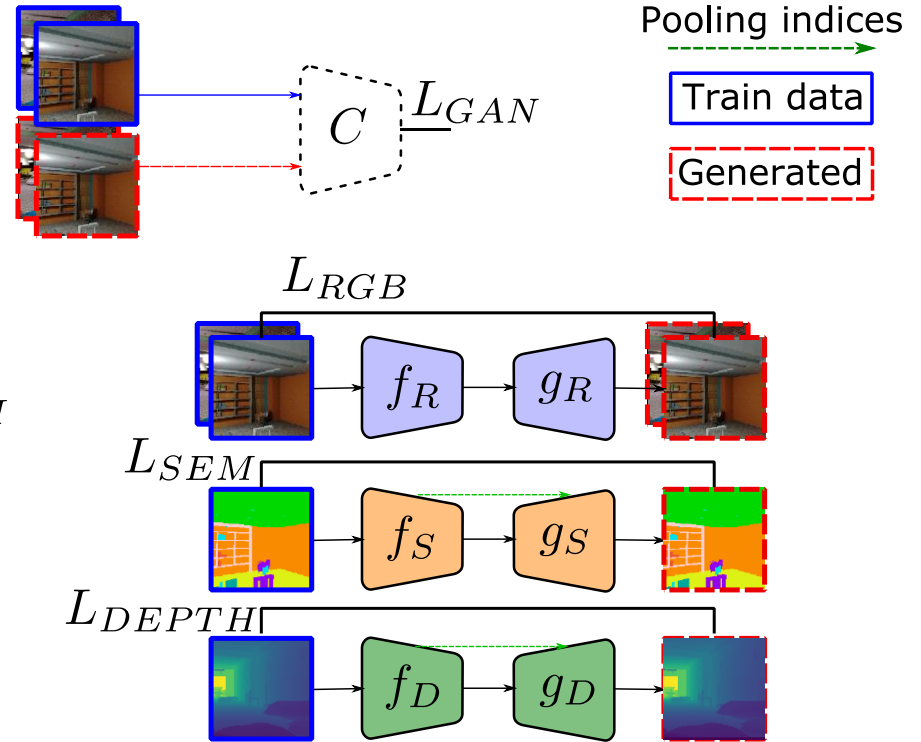
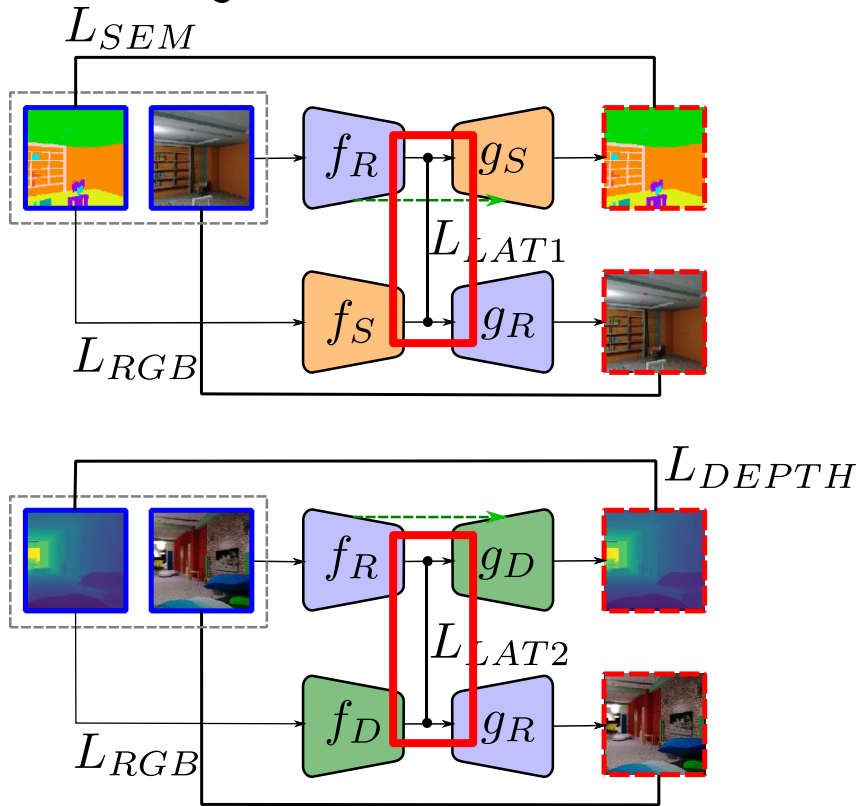


# Zero-pair translation with mix and match networks

Training for encoder-decoder alignment:

Autoencoders  
Shared encoder/decoders

Latent losses

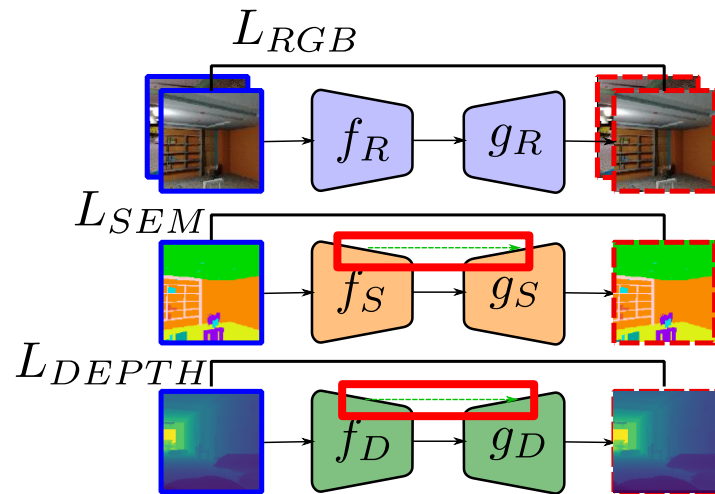
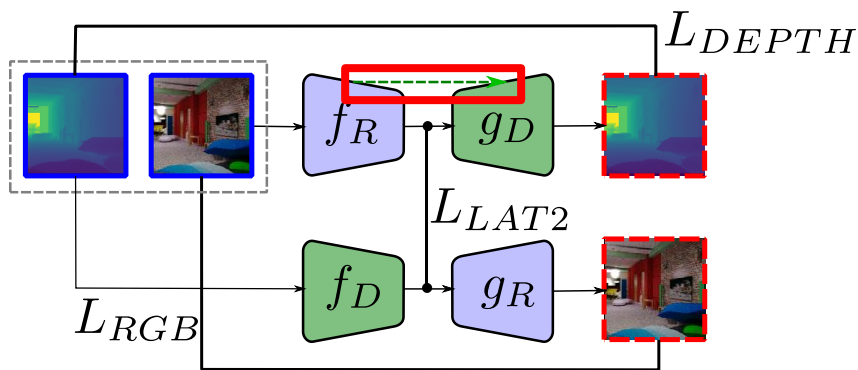
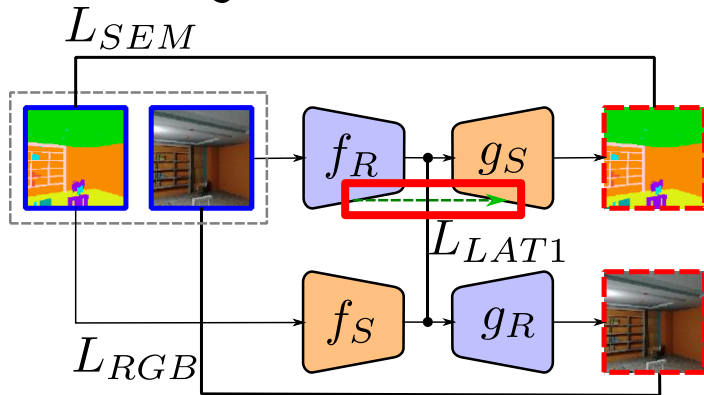


# Zero-pair translation with mix and match networks

Training for encoder-decoder alignment:

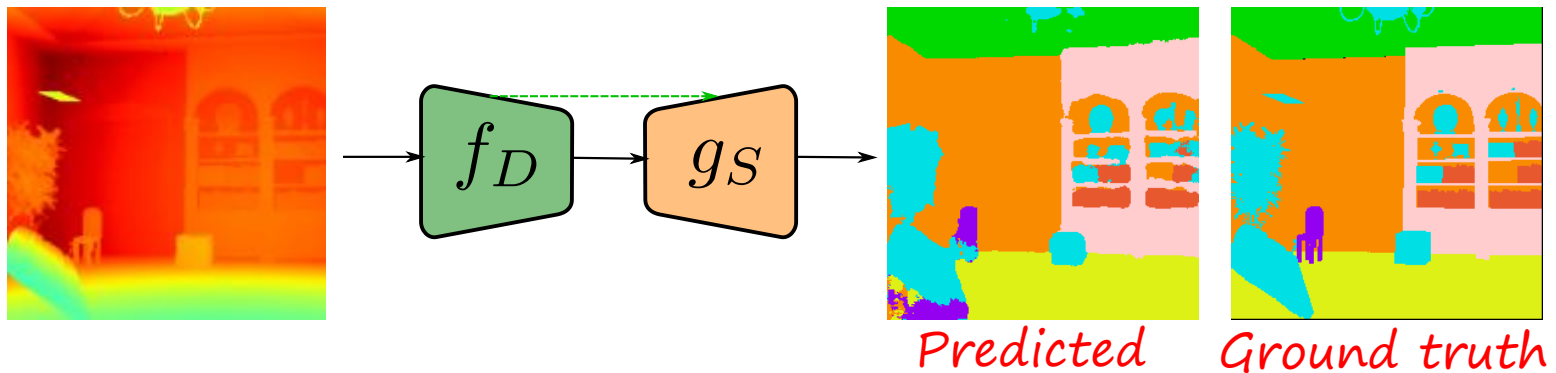
Autoencoders  
Shared encoder/decoders

Latent losses  
**Robust side information (pooling indices)**



# Zero-pair translation with mix and match networks

*Test on zero-pair translation depth-to-segmentation*



# Comparison: depth-to-segmentation

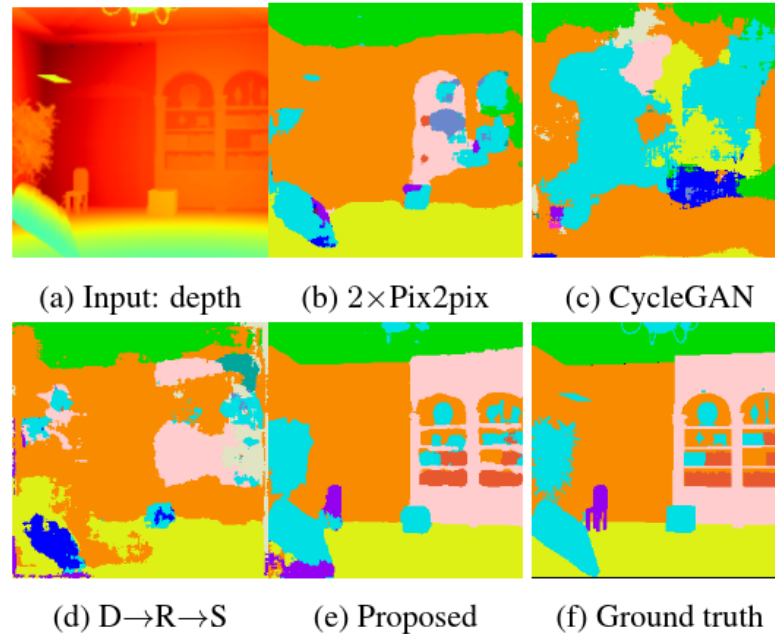
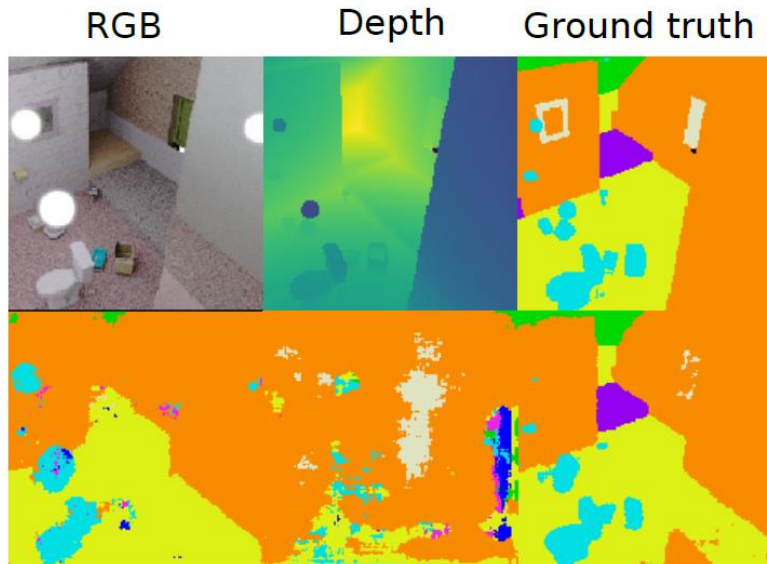


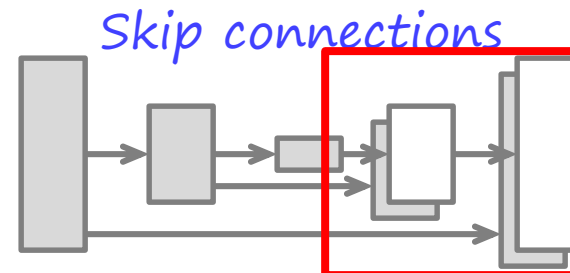
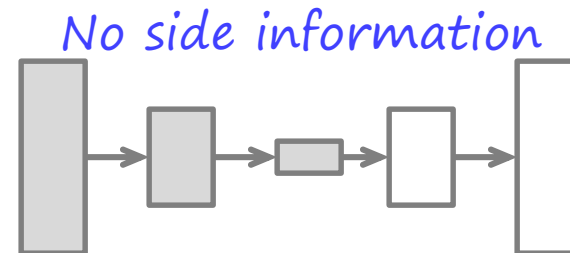
Figure 1: Zero-pair depth→segmentation, trained on (depth,RGB) and (RGB,segmentation).

# Side information in mix and match networks

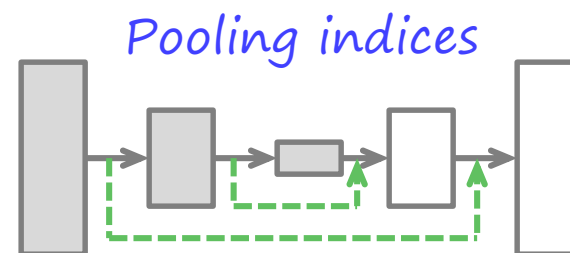


No side information      Skip connections      Pooling indices

Side information	Pretrained	mIoU	Global
-	N	32.2%	63.5%
Skip connections	N	14.1%	52.6%
Pooling indices	N	45.6%	73.4%
Pooling indices	Y	49.5%	80.0%



*Decoder conditioned on seen encoder(s)*



*Seems more invariant and robust*

# Quantitative evaluation

Method	Conn.	$L_{SEM}$	Bed	Book	Ceiling	Chair	Floor	Furniture	Object	Picture	Sofa	Table	TV	Wall	Window	mIoU	Global
<b>Baselines</b>																	
CycleGAN [34]	SC	CE	2.79	0.00	16.9	6.81	4.48	0.92	7.43	0.57	9.48	0.92	0.31	17.4	15.1	6.34	14.2
2×pix2pix [10]	SC	CE	34.6	1.88	70.9	20.9	63.6	17.6	14.1	0.03	38.4	10.0	4.33	67.7	20.5	25.4	57.6
M&MNet $D \rightarrow R \rightarrow S$	PI	CE	0.02	0.00	8.76	0.10	2.91	2.06	1.65	0.19	0.02	0.28	0.02	58.2	3.3	5.96	32.3
M&MNet $D \rightarrow R \rightarrow S$	SC	CE	25.4	0.26	82.7	0.44	56.6	6.30	23.6	5.42	0.54	21.9	10.0	68.6	19.6	24.7	59.7
<b>Zero-pair</b>																	
M&MNet $D \rightarrow S$	PI	CE	50.8	18.9	89.8	31.6	88.7	48.3	44.9	62.1	17.8	49.9	51.9	86.2	79.2	55.4	80.4
<b>Multi-modal</b>																	
M&MNet $(R, D) \rightarrow S$	PI	CE	49.9	25.5	88.2	31.8	86.8	56.0	45.4	70.5	17.4	46.2	57.3	87.9	79.8	57.1	81.2

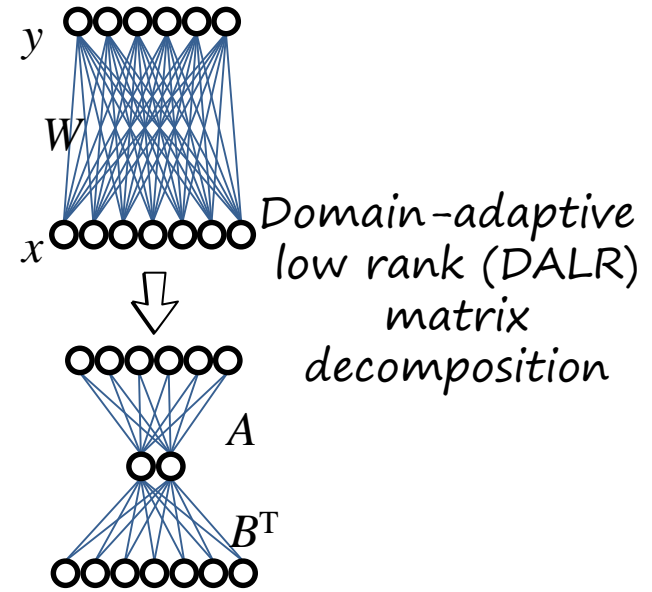
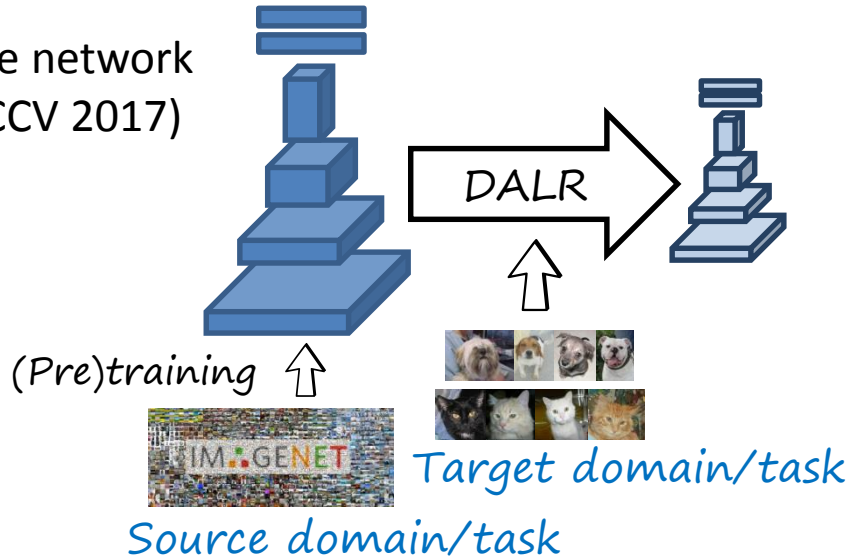
Table 3: Zero-pair depth-to-semantic segmentation. **SC**: skip connections, **PI**: pooling indexes, **CE**: cross-entropy

# Outline

- Introduction
- Transferring GANs (ECCV 2018)
- Rotated elastic weight consolidation (ICPR 2018)
- Memory Replay GANs (NIPS 2018)
- Mix and match networks (CVPR 2018)
- **Other works**

# Other works

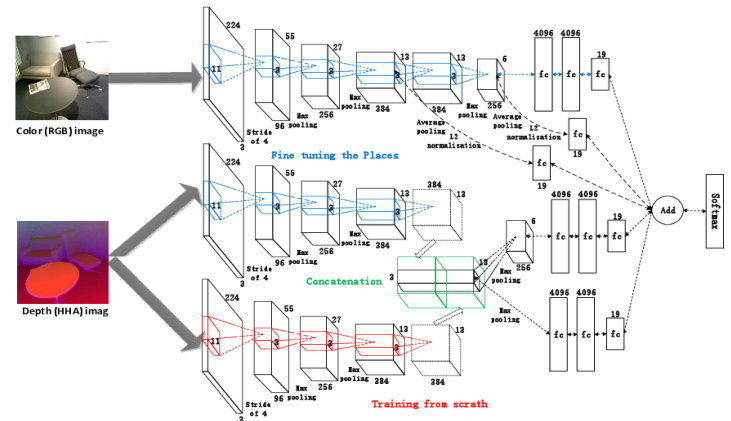
Domain-adaptive network compression (ICCV 2017)



Contextual food recognition and analysis (TMM17, TMM18)



RGB-D deep representations (AAAI17, IJCAI17, TIP18)



# THANK YOU!

[lherranz@cvc.uab.es](mailto:lherranz@cvc.uab.es)

[www.lherranz.org](http://www.lherranz.org)

