

**Spring 2021**

# ADVANCED TOPICS IN COMPUTER VISION

---

**Atlas Wang**

Assistant Professor, The University of Texas at Austin

# Why Synthetic Training

- Collections of real data are costly
  - Massive real image
  - Classification / Segmentation / Detection
- Synthetic data are relatively cheap to generate



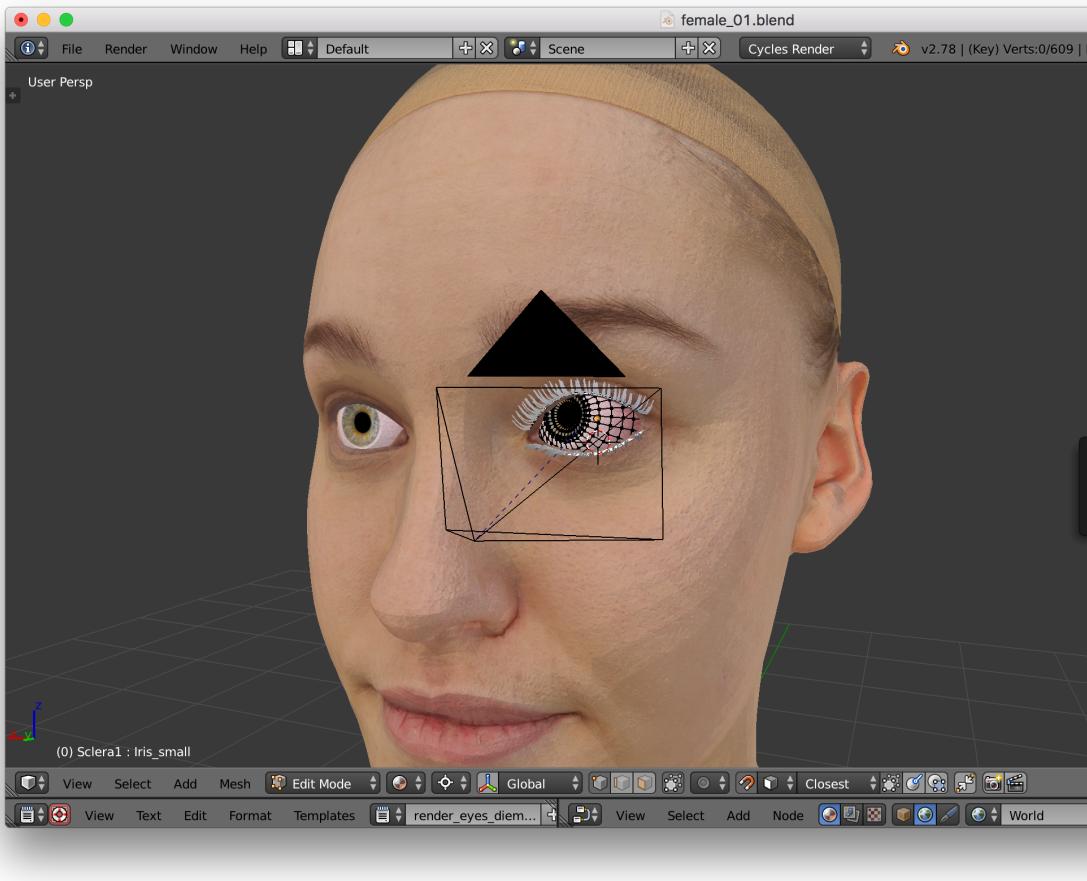
Cityscapes (3K annotations)



GTA5 (24,966 annotations)

# Why Synthetic Training

- In some cases, synthetic data is all you have...
- EyeGaze / Depth / Flow / 3D Mesh reconstruction / Robotics



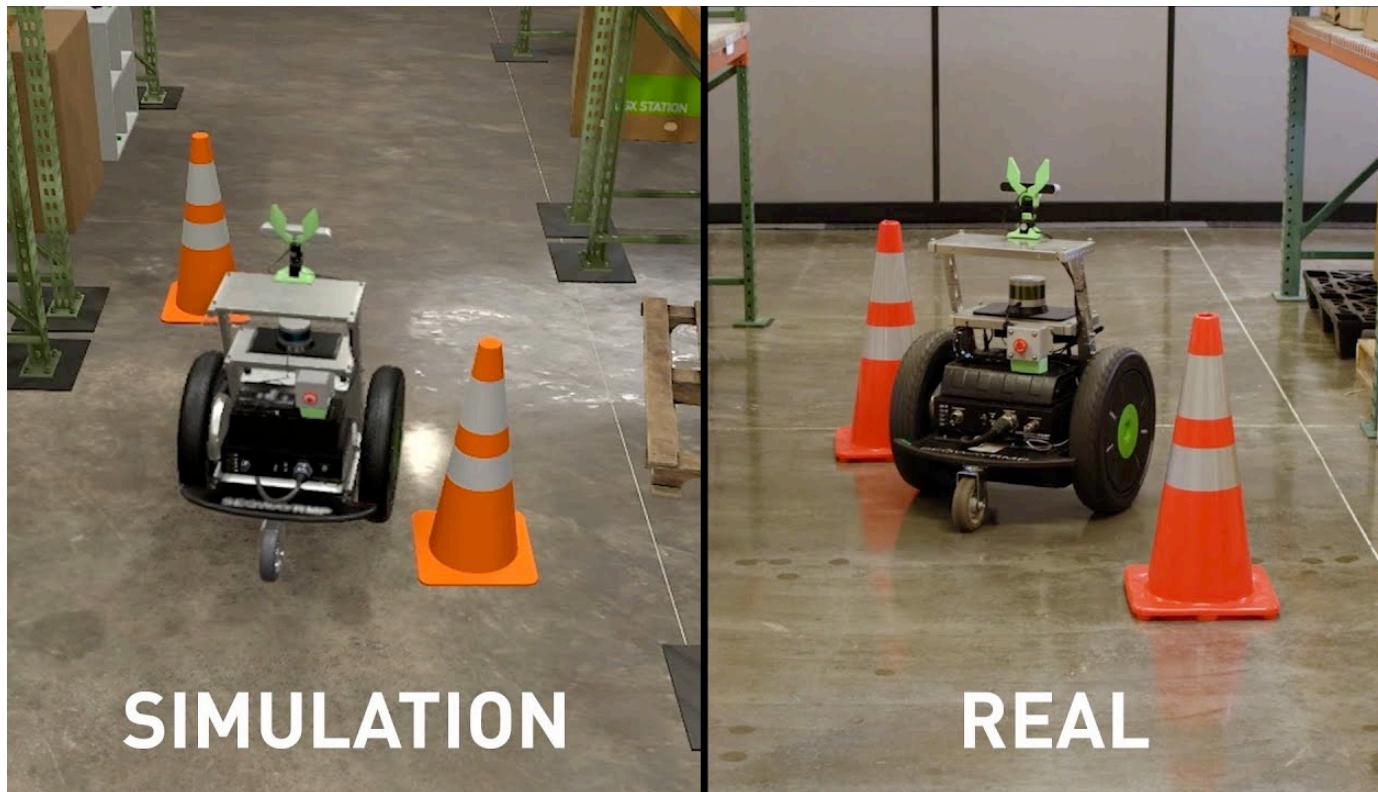
Wood et al. ICCV 2015



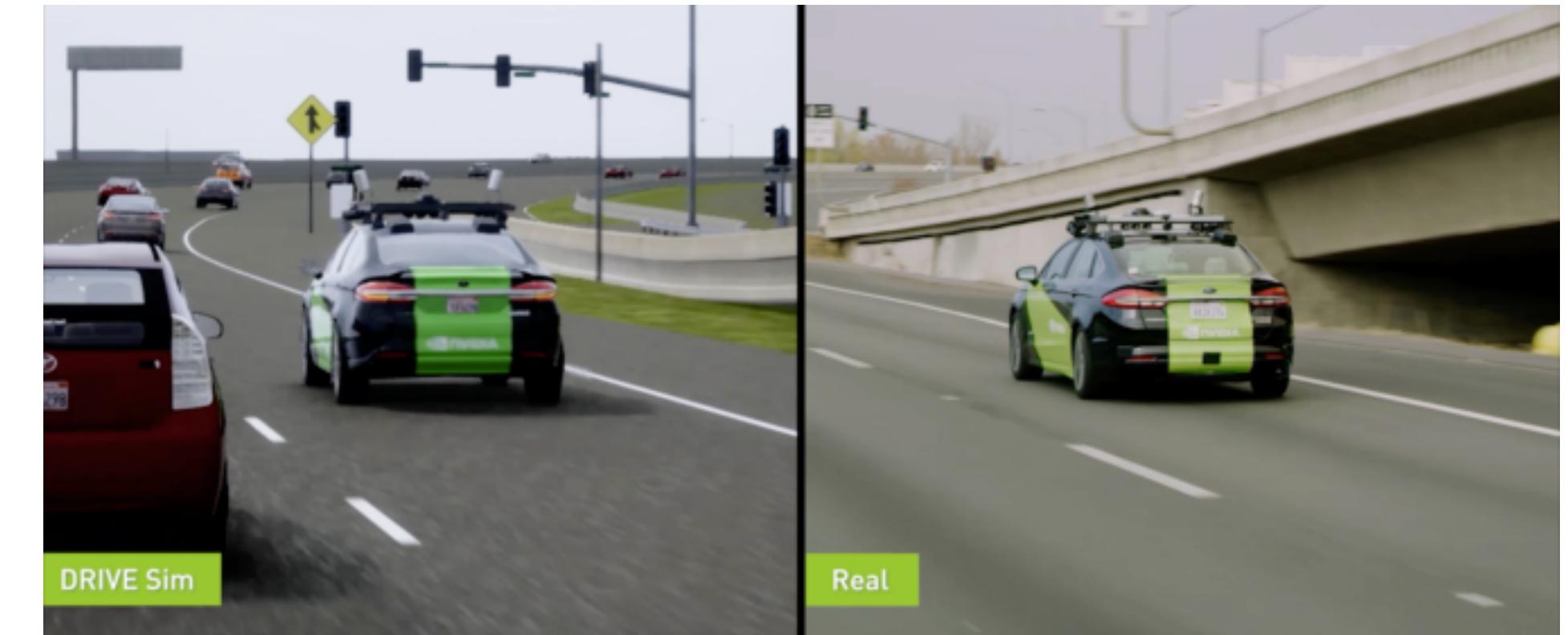
Habitat (Facebook)

# Synthetic Simulation Empowers Some Most Important Applications

- Autonomous Driving: Omniverse, ISAAC, DRIVE Sim, etc.



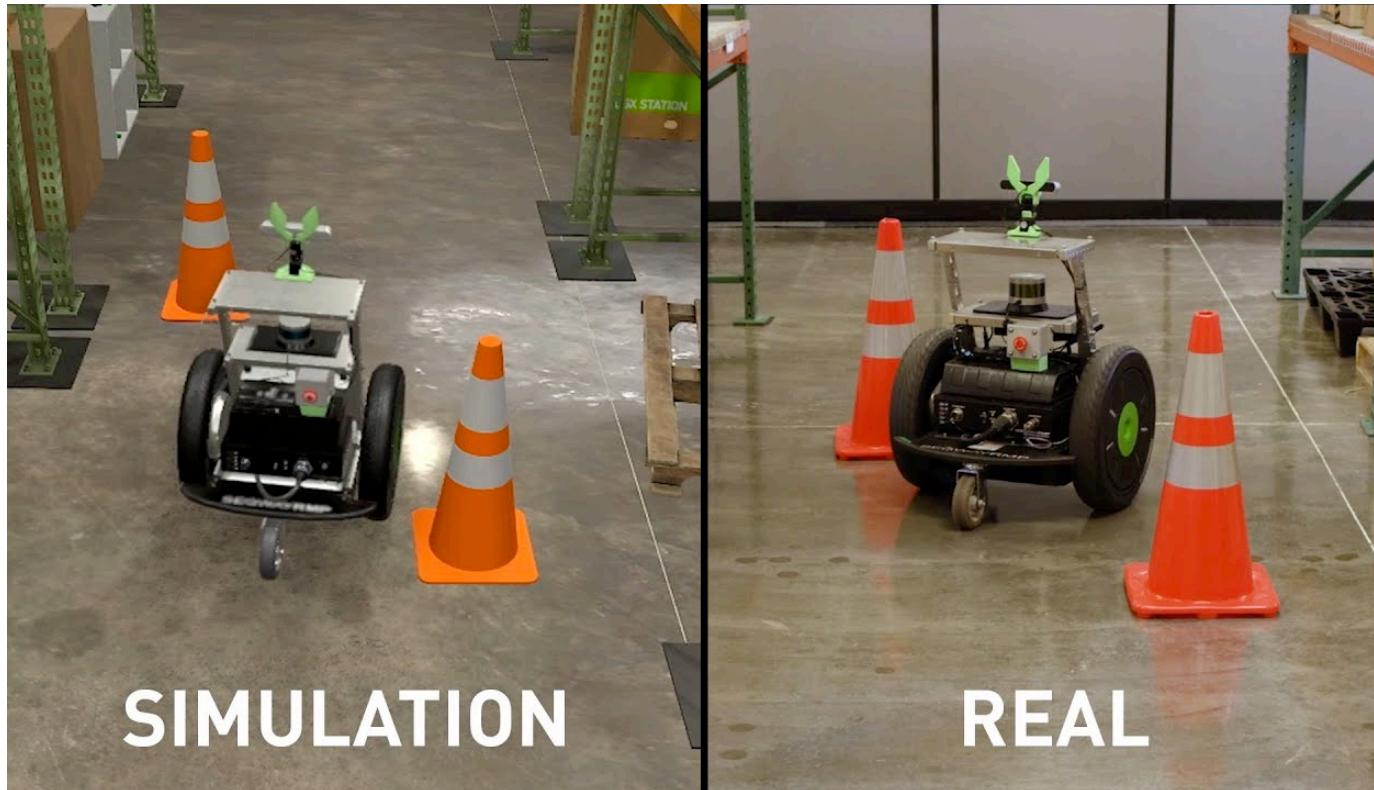
ISAAC platform



DRIVE Sim

# Synthetic Simulation Empowers Some Most Important Applications

- Autonomous Driving: Omniverse, ISAAC, DRIVE Sim, etc.



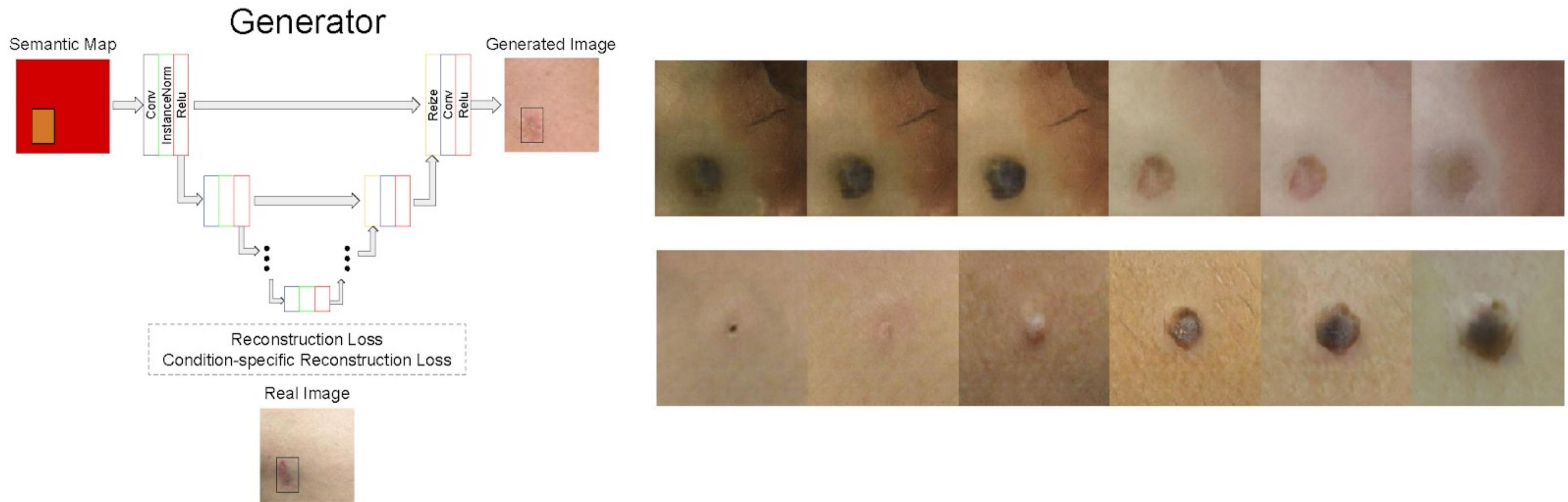
ISAAC platform



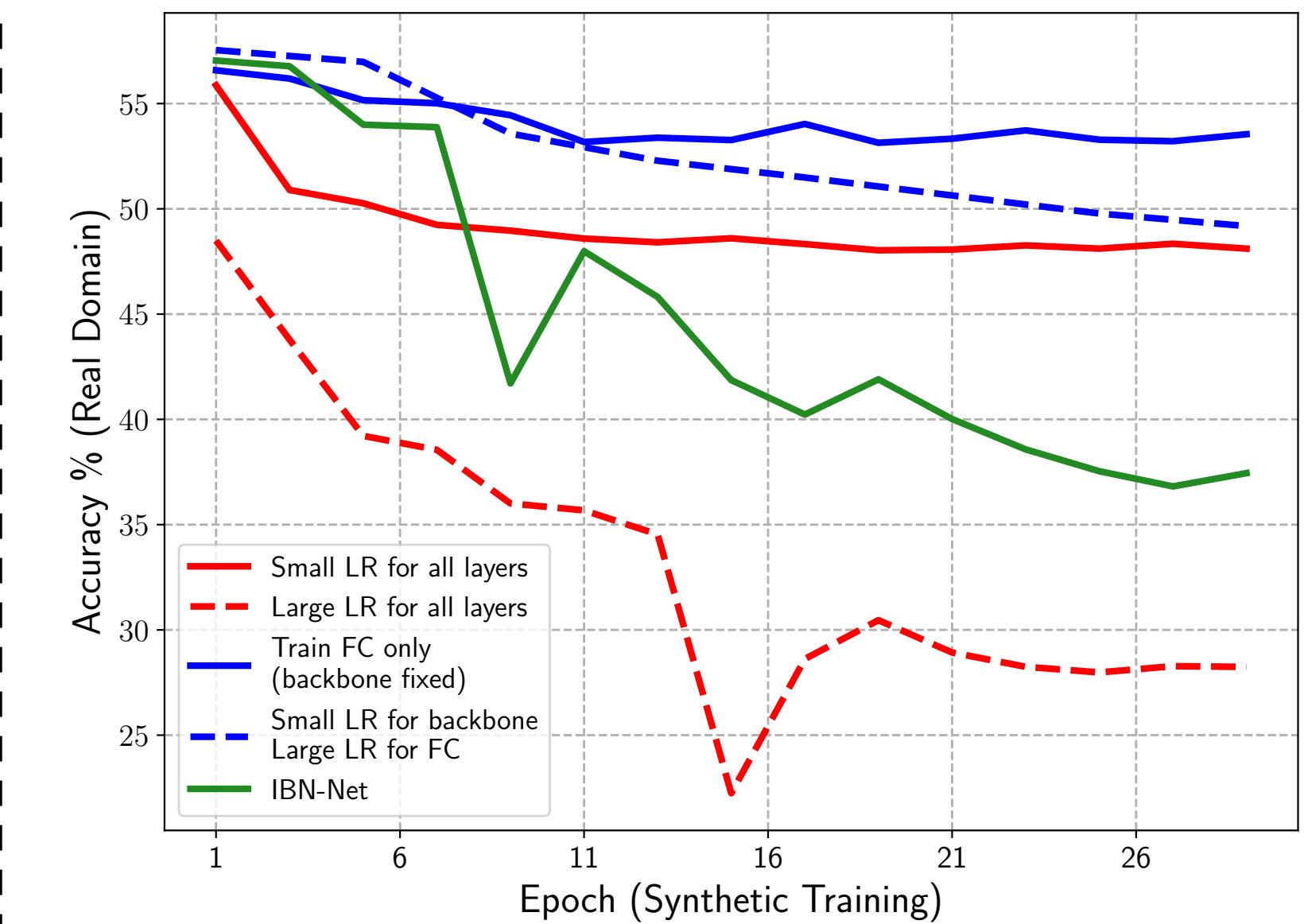
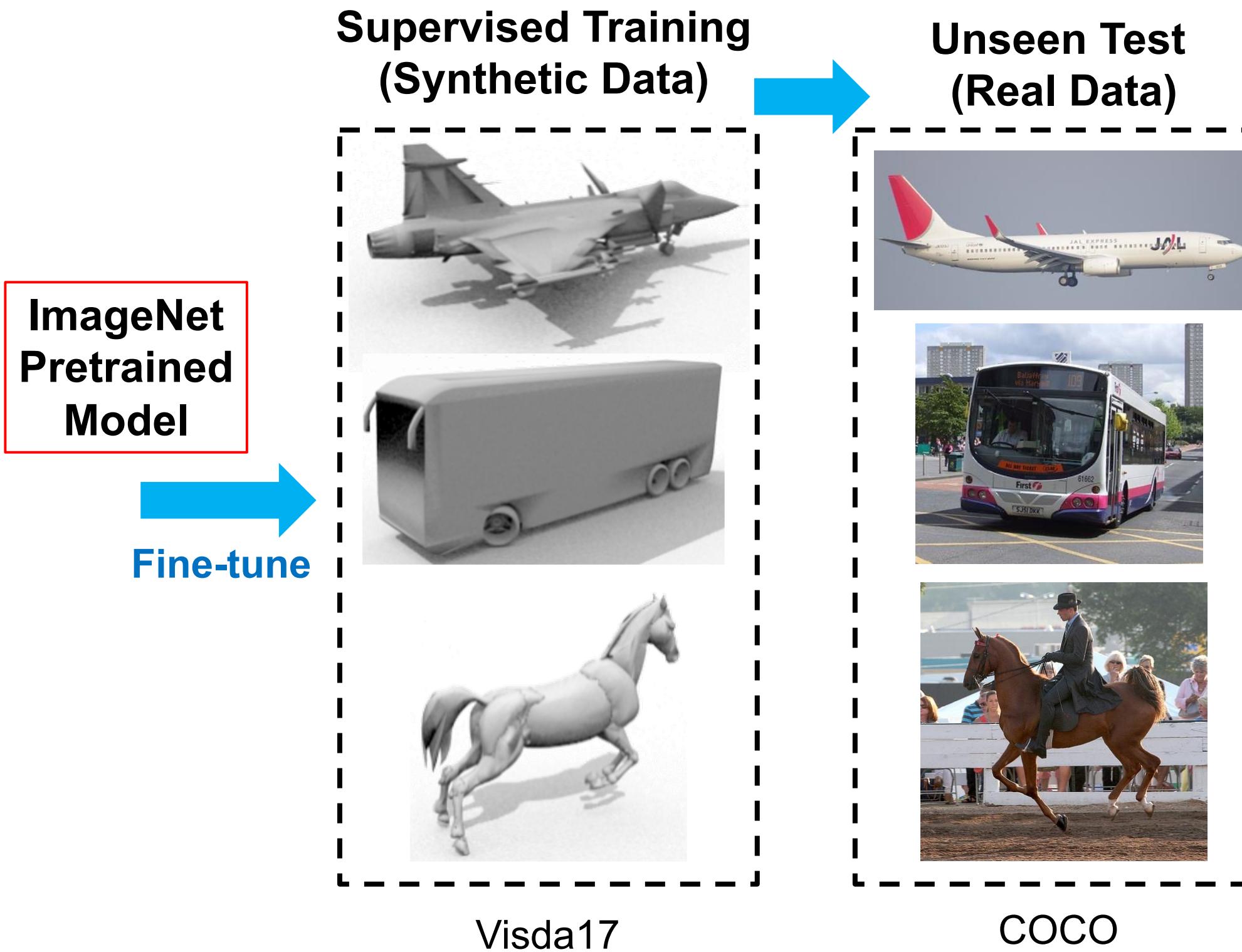
DRIVE Sim

# Synthetic Simulation Empowers Some Most Important Applications

- Medical Image Analysis: cover more corner cases, resolve privacy concerns...



# Challenging Domain Gap: Synthetic vs Real



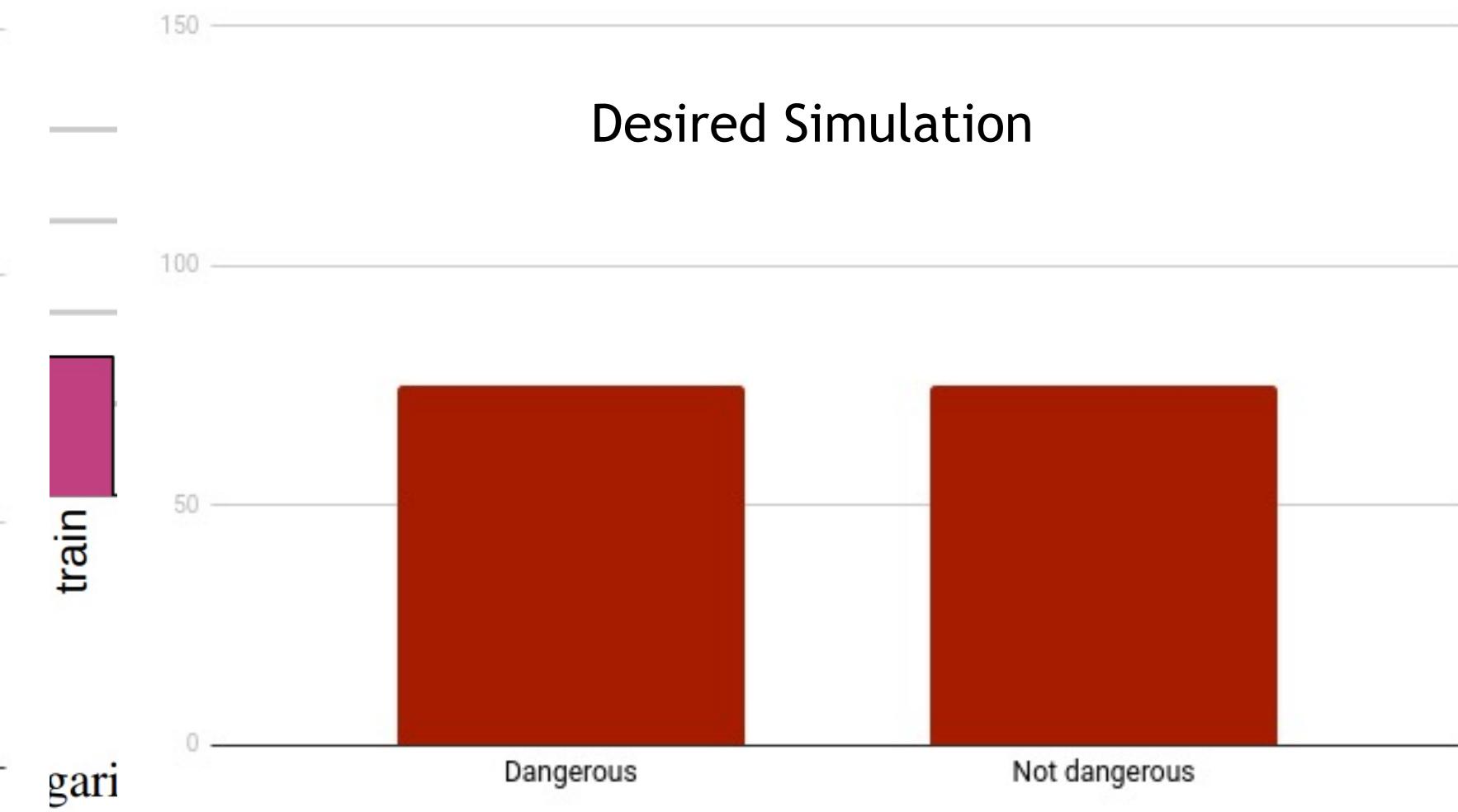
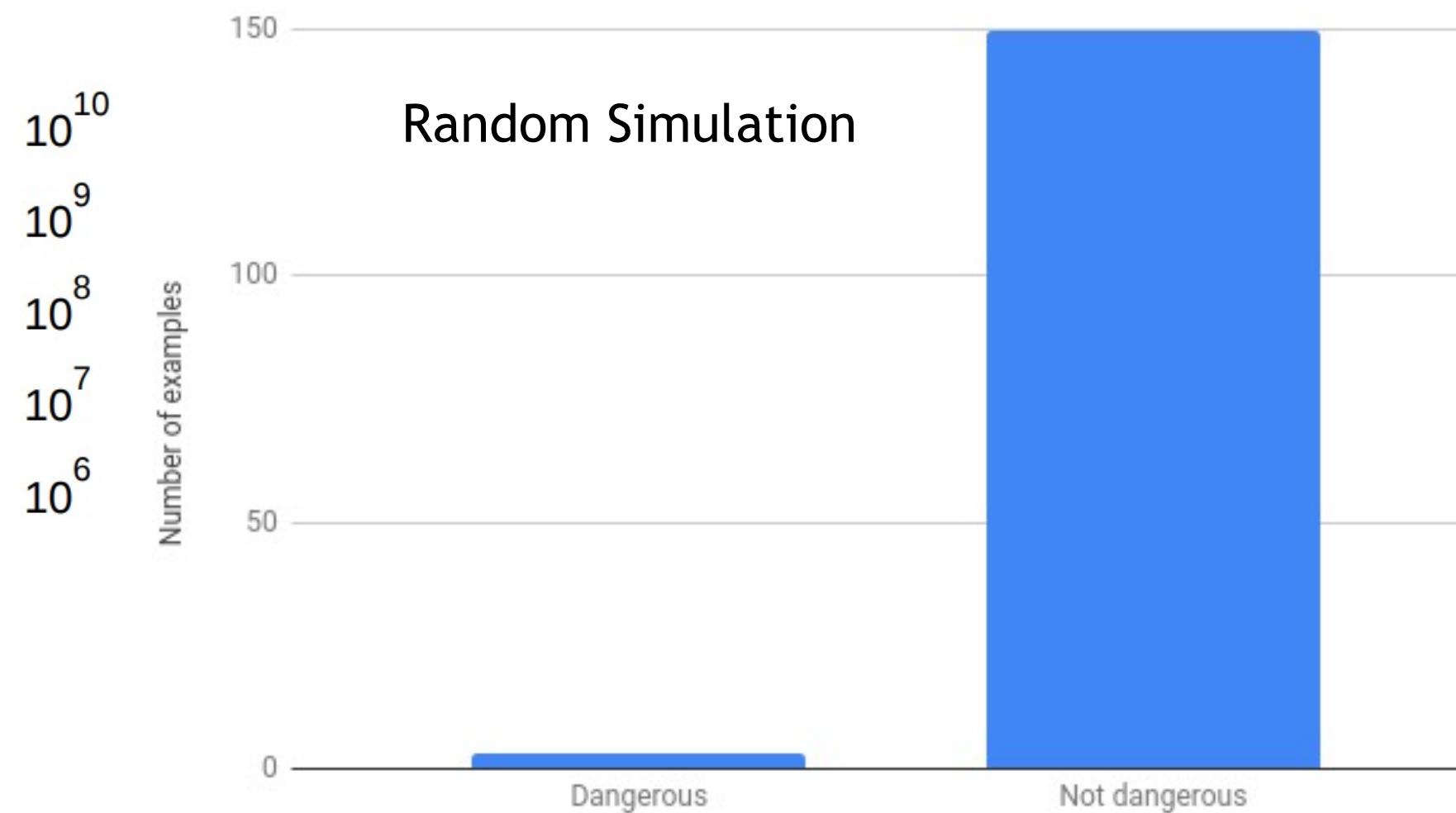
# Domain Randomization (IROS'17)

- To handle the variability in real-world data, the simulator parameters (lighting, pose, object textures, etc) are randomized in non-realistic ways to force the learning of essential diverse features.



# Can We Do Better than Random?

- Learn to simulate better data for a particular downstream task?
- Learn to simulate edge cases?



gari

# Learning to Simulate (ICLR'19)

- We want to solve the following bi-level optimization problem.

Simulation parameters

$$\psi^* = \arg \min_{\psi} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{\text{val}}} \mathcal{L}(y, h_{\theta}(\mathbf{x}; \theta^*(\psi))) \quad \xrightarrow{\text{meta-learner}}$$

$$\text{s.t. } \theta^*(\psi) = \arg \min_{\theta} \sum_{(\mathbf{x}, \mathbf{y}) \in D_{q(\mathbf{x}, \mathbf{y} | \psi)}} \mathcal{L}(\mathbf{y}, h_{\theta}(\mathbf{x}, \theta)), \quad \xrightarrow{\text{main task model } h_{\theta} \text{ trained on simulated data}}$$

Main task model parameters

Loss of main task model  
trained in simulation and  
evaluated on real data

Dataset generated by simulator

# Learning to Simulate (ICLR'19)

- Train the policy of selecting simulator parameters, using policy gradient, since the simulator is often non-differentiable



# Are better simulators enough?

Models overfit to any difference

High quality is expensive



## Virtual KITTI Dataset

**Multi-object tracking accuracy:**

**Sim: 63.7%**

**Real: 78.1%**

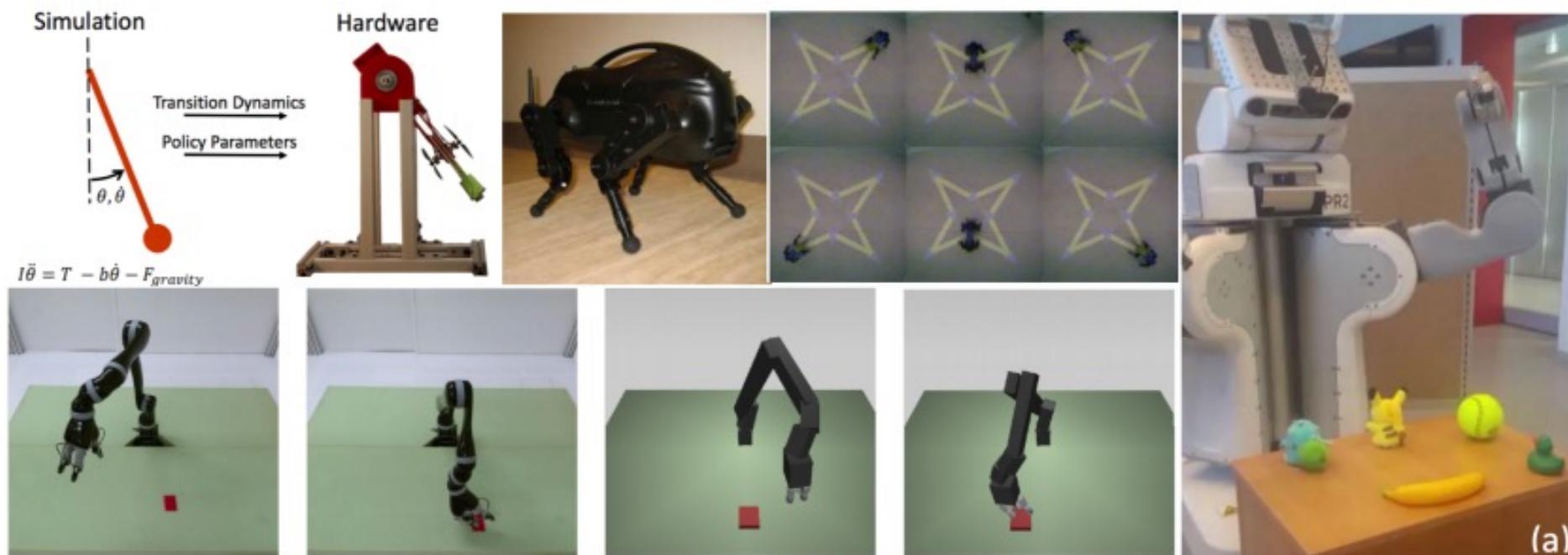
**Virtual Worlds as Proxy for Multi-Object Tracking Analysis**

[Gaidon\*, Wang\*, Cabon, Vig, 2016]

**Jungle Book:**  
**30M render hours**  
**19 hours per frame**  
**800 artist-years of effort**  
**Jungle Book, 2016**

# Supervised domain adaptation

## Fine-tuning



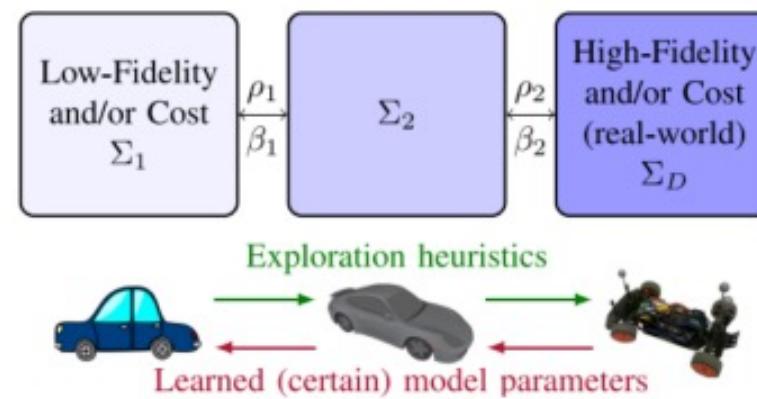
**Learning Omnidirectional Path Following Using Dimensionality Reduction** [Kolter, Ng, 2003]

**Efficient Reinforcement Learning for Robotics using Informative Simulated Priors** [Cutler, How, 2015]

**Sim-to-Real Robot Learning from Pixels with Progressive Nets** [Rusu et al. 2016]

**Deep Predictive Policy Training using Reinforcement Learning** [Ghadirzadeh, Maki, Kragic, Bjorkman, 2017]

## Iterative learning control



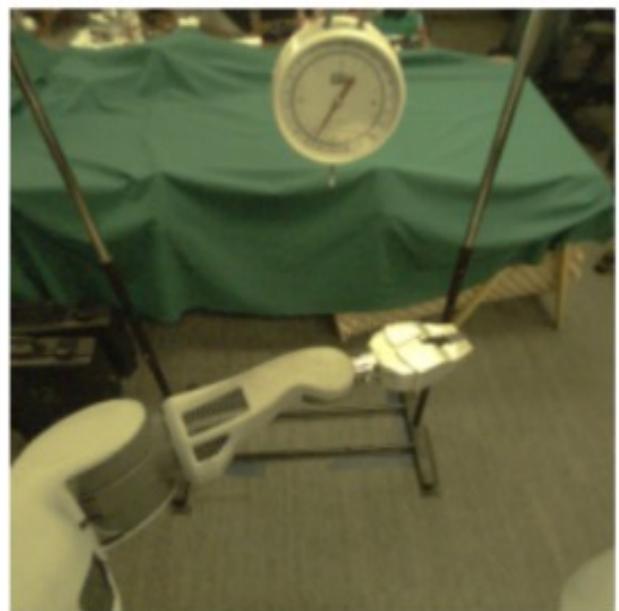
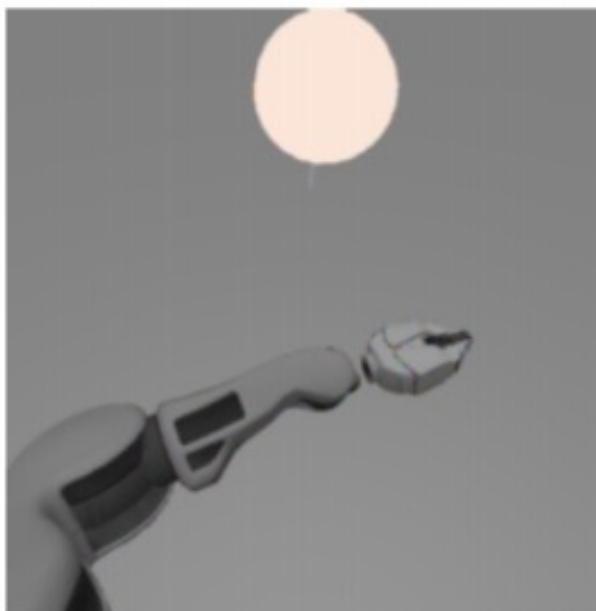
**Using inaccurate models in reinforcement learning** [Abbeel, Quigley, Ng, 2006]

**Reinforcement learning with multi-fidelity simulators** [Cutler, Walsh, How 2014]

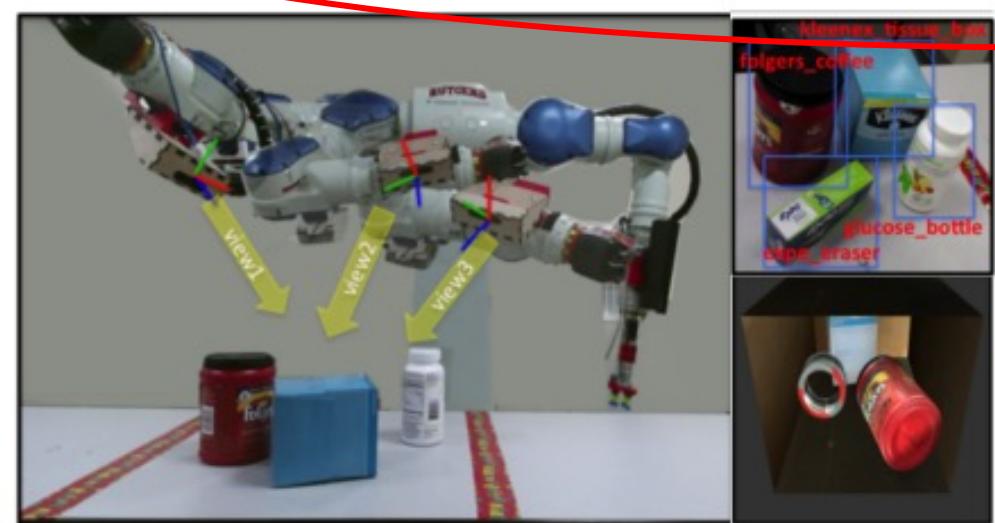
**Superhuman performance of surgical tasks by robots using iterative learning from human-guided demonstrations** [Van Den Berg, Miller, Duckworth, Hu, Wan, Fu, Goldberg, Abbeel, 2010]

# (Less) supervised domain adaptation

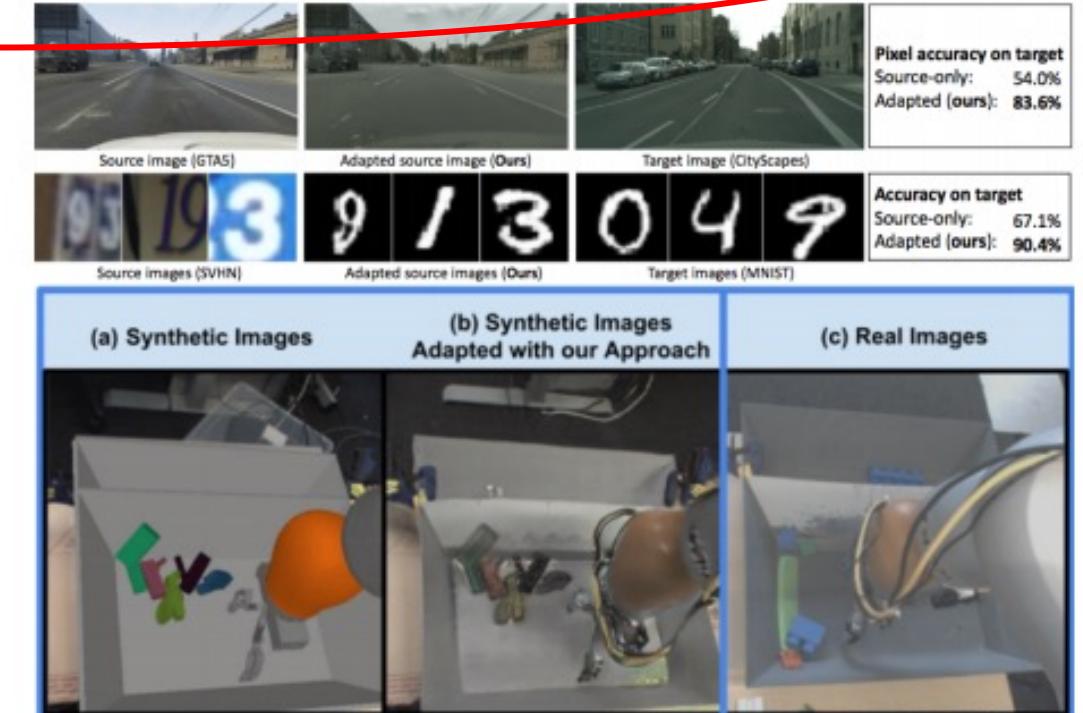
## Weakly Supervised



## Self-Supervised



## Unsupervised



**Adapting Deep Visuomotor Representations with Weak Pairwise Constraints** [Tzeng, Devin, Hoffman, Finn, Abbeel, Levine, Saenko, Darrell, 2016]

**A Self-supervised Learning System for Object Detection using Physics Simulation and Multi-view Pose Estimation** [Mitash, Bekris, Boualiyas, 2017]

**CyCADA** [Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Efros, Darrel, 2017]  
**Using Simulation and Domain Adaptation to Improve Efficiency of Deep Robotic Grasping** [Bousmalis et al., 2017]



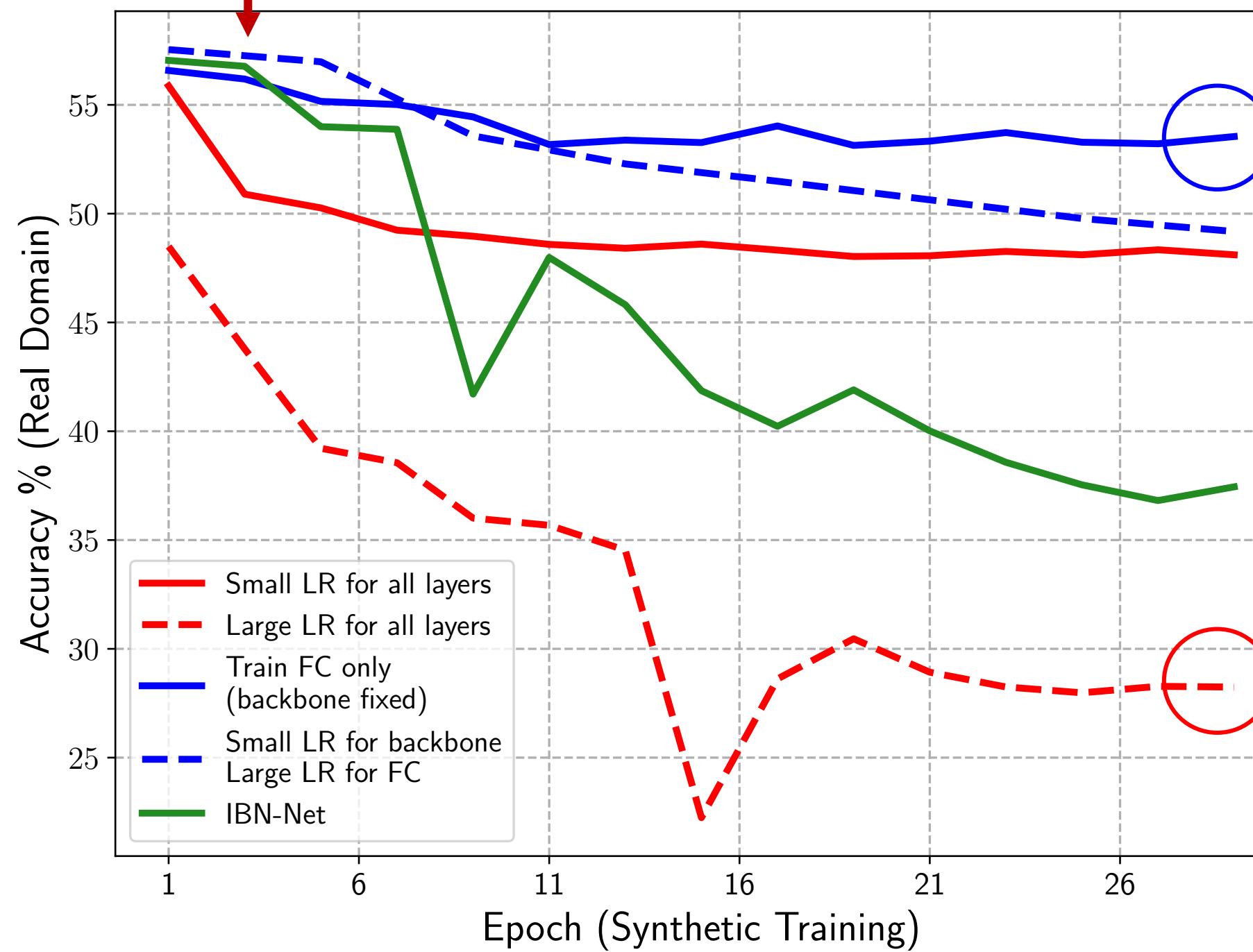
# Automated Synthetic-to-Real Generalization

ICML 2020

Wuyang Chen, Zhidong Yu, Zhangyang (Atlas) Wang, Anima Anandkumar

# Previous solutions: Heuristic Hand-tuning

1) Early stopping

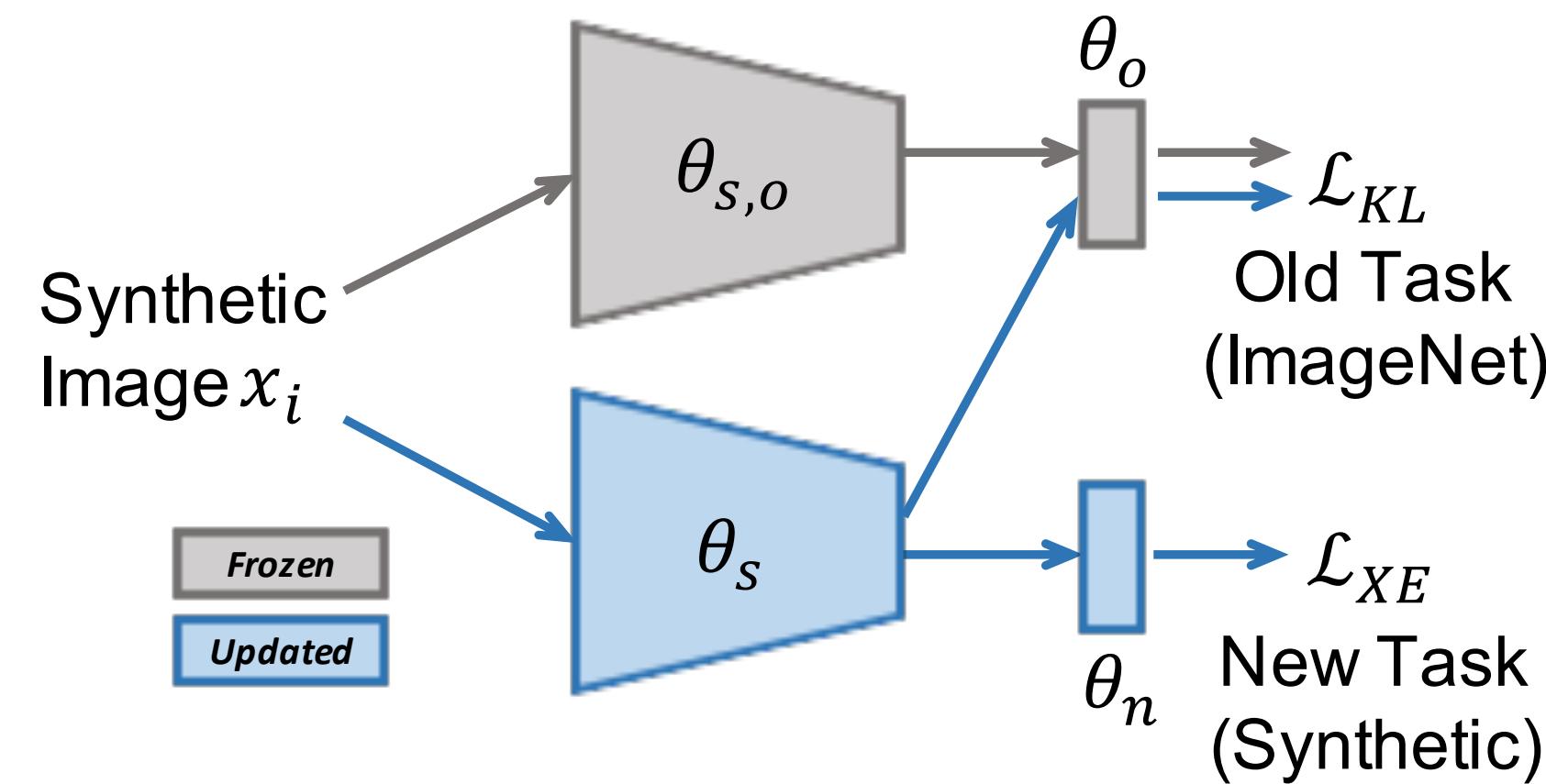


2) Small LR

Large LR

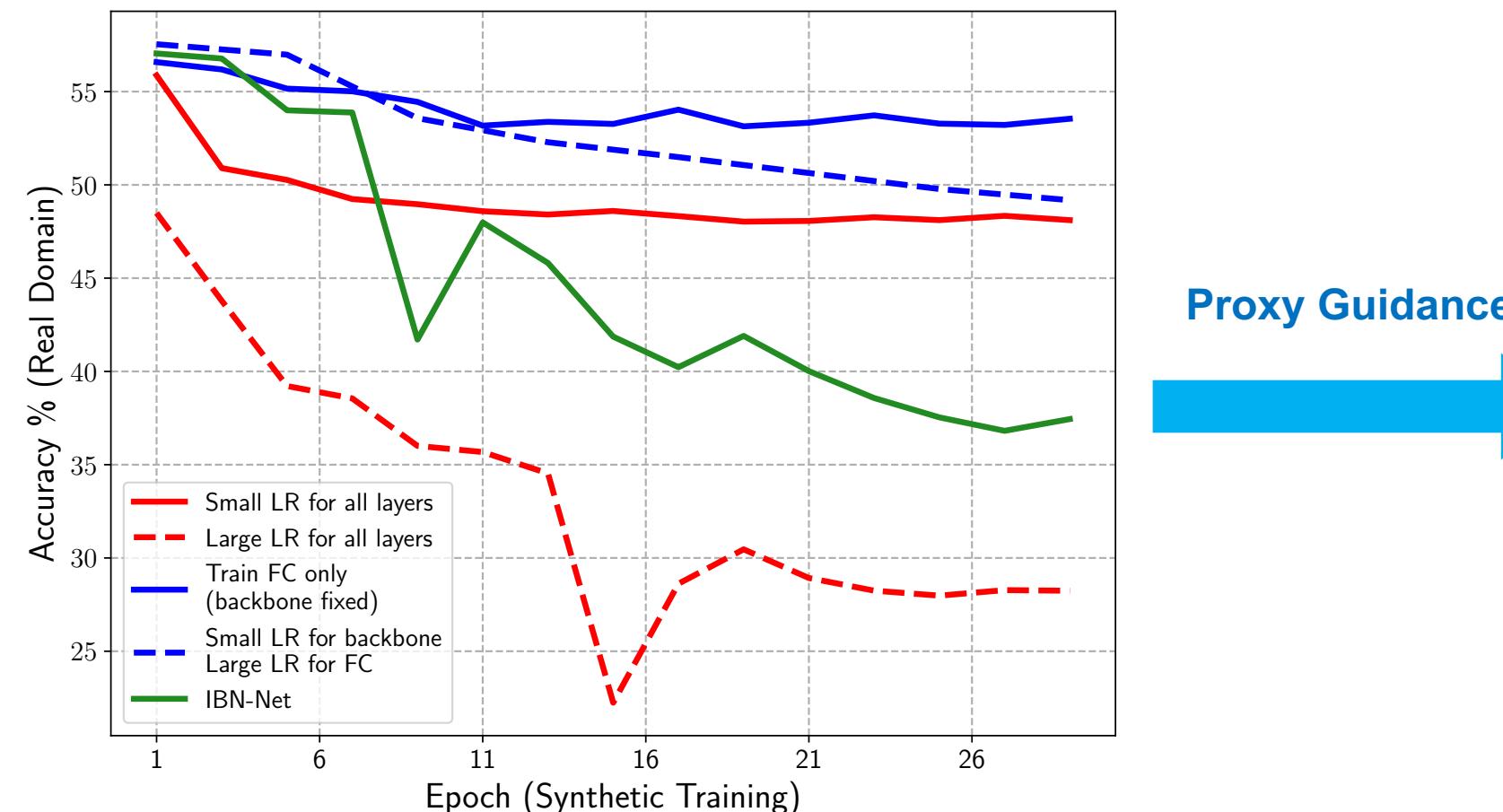
# ImageNet as Proxy Guidance

- Why early stopping?
  - Keep weights close to ImageNet initialization.
- We minimize  $\mathcal{L}_{KL}$ : new model vs ImageNet initialization
  - ImageNet as proxy guidance in syn2real training.

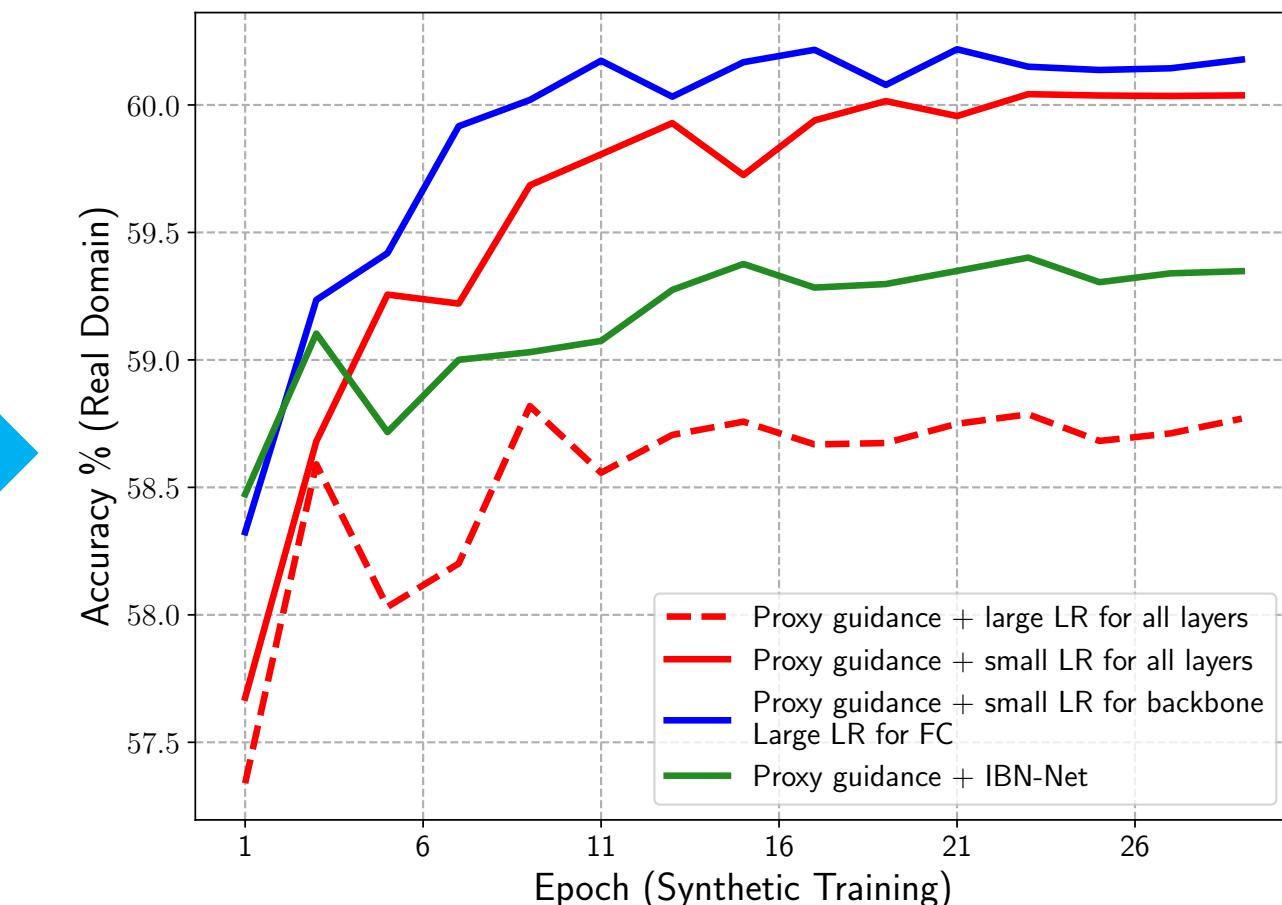


# ImageNet as Proxy Guidance

- Why early stopping?
  - Keep weights close to ImageNet initialization.
- We minimize  $\mathcal{L}_{KL}$ : new model vs ImageNet initialization
  - ImageNet as proxy guidance in syn2real training.

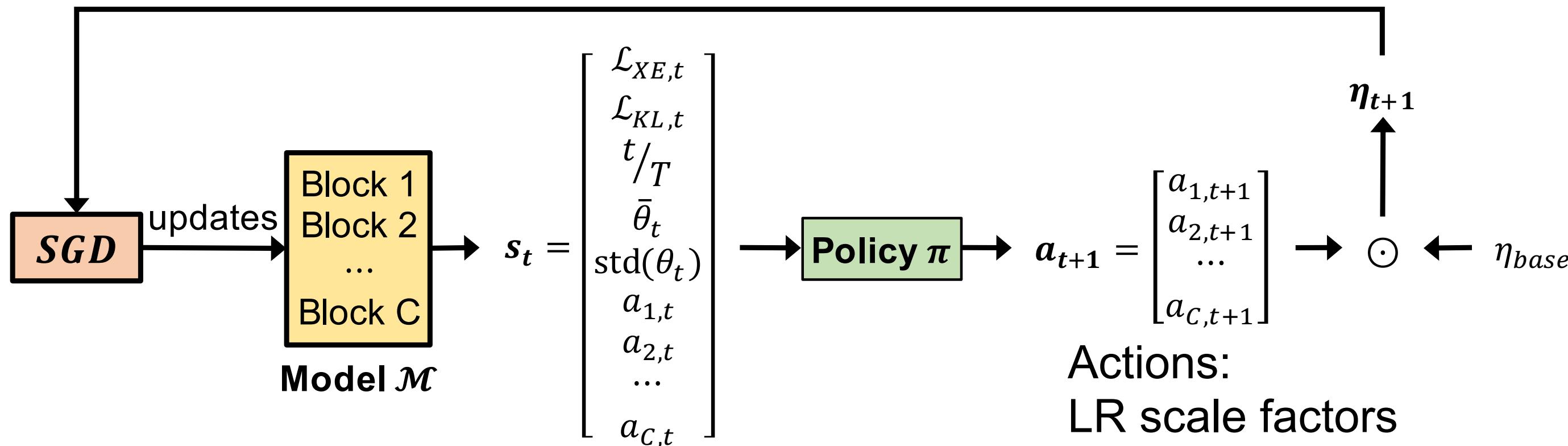


Proxy Guidance  
→



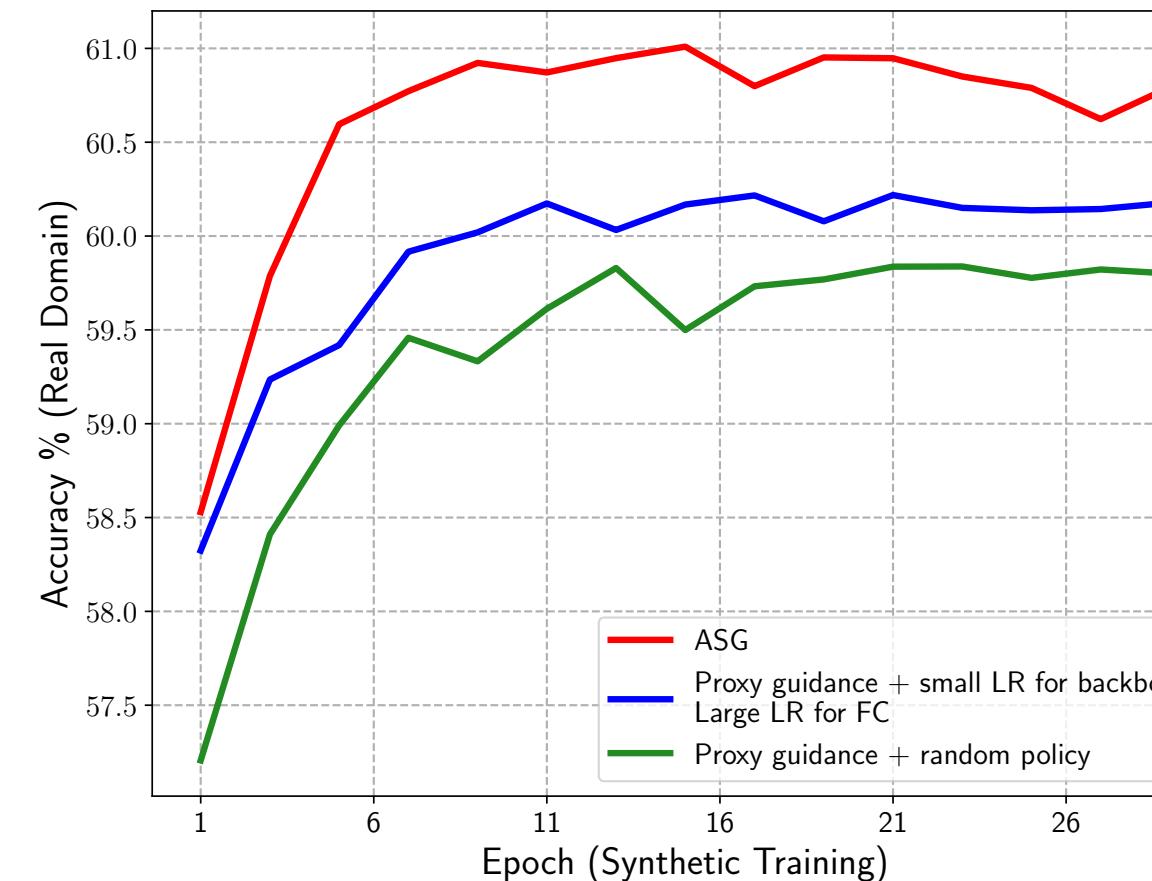
# L2O: automatic control of layer-wise learning rate

- Why small learning rate?
  - Keep weights close to ImageNet initialization.
- But how small for which layer?
  - L2O (learning-to-optimize): automatic control of layer-wise learning rate

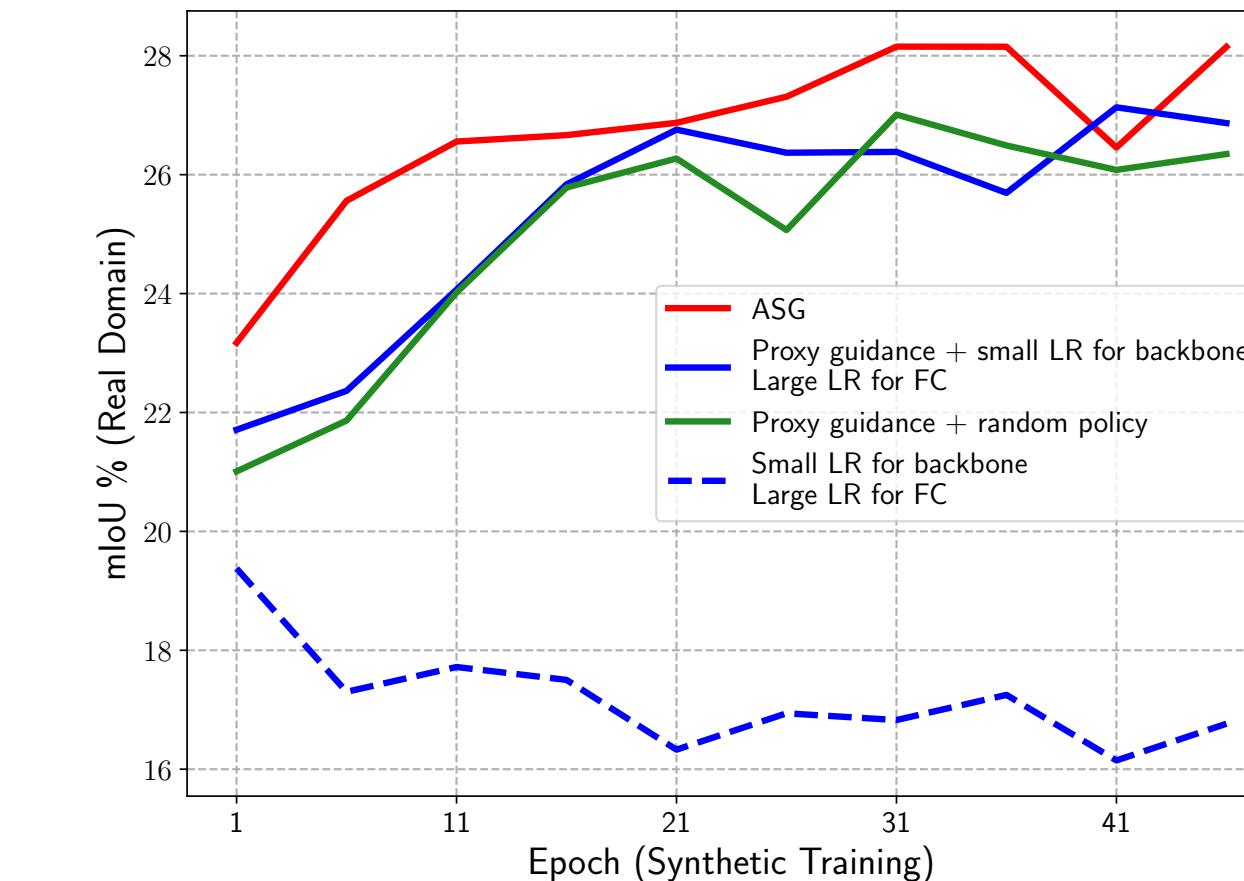


# Automated Synthetic-to-Real Generalization (ASG)

- Why small learning rate?
  - Keep weights close to ImageNet initialization
- But how small for which layer?
  - L2O (learning-to-optimize): automatic control of layer-wise learning rate



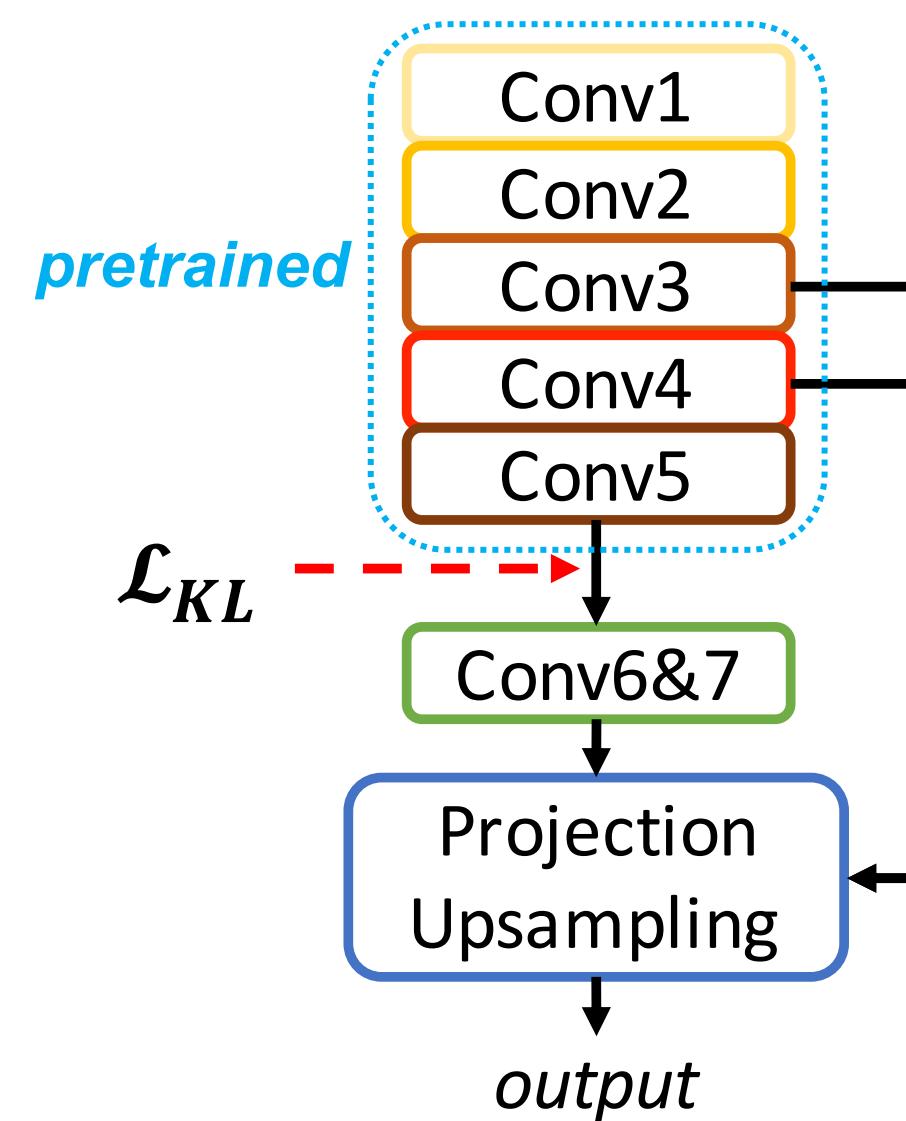
Visda17 → COCO



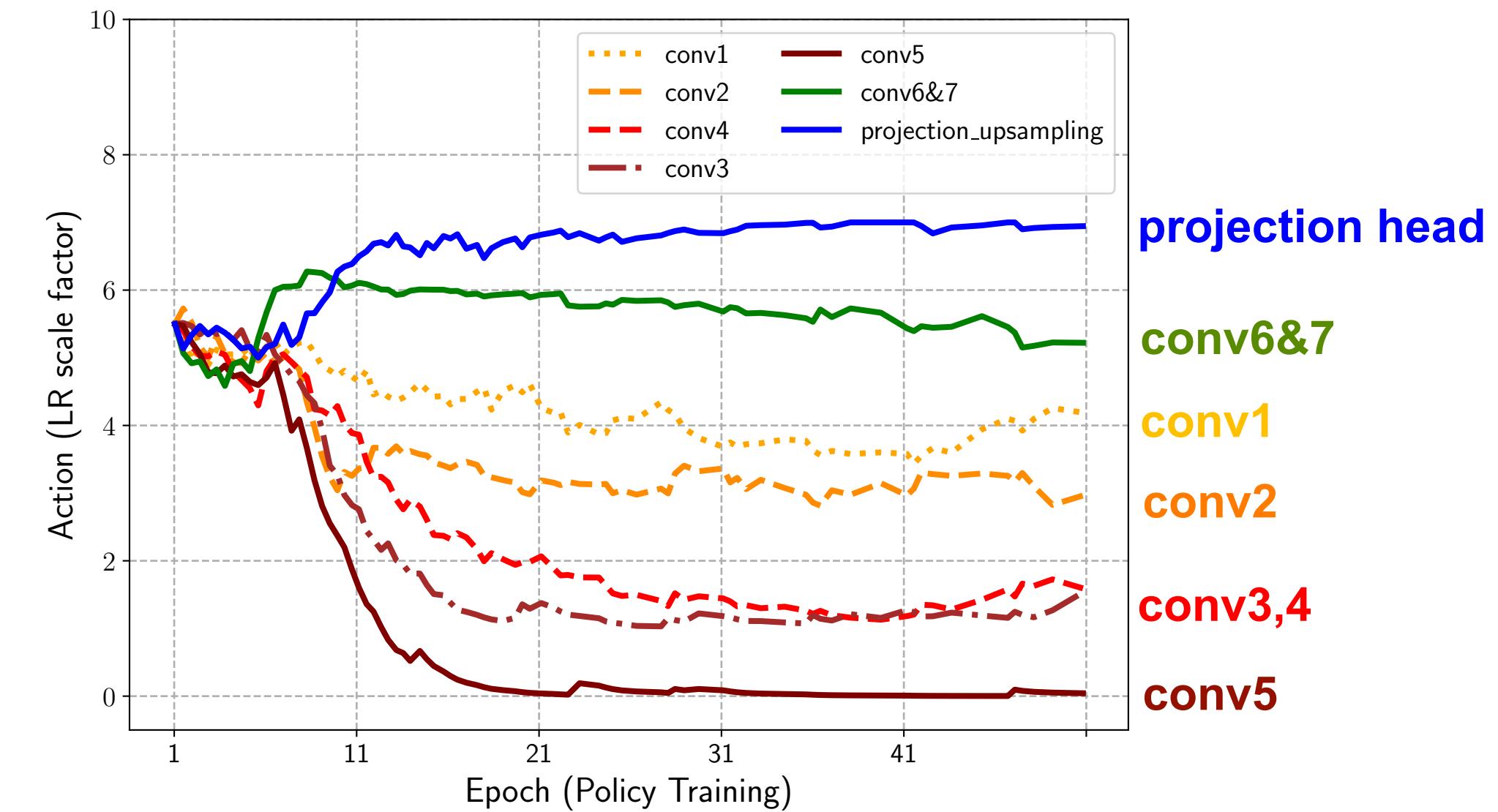
GTA5 → Cityscapes

# Action Behavior of RL-L2O Policy

- Backbone (ImageNet pretrained): closer to  $\mathcal{L}_{KL}$   $\rightarrow$  smaller LR
- Projection head: large LR



FCN-Vgg16



# Why ASG Works? Retaining ImageNet Information

#	Model	Visda-17	ImageNet
1.	Large LR for all layers	28.2	0.8
2.	+ our Proxy Guidance	58.7 (+30.5)	76.2 (+75.4)
3.	Small LR for backbone and large LR for FC	49.3	33.1
4.	+ our Proxy Guidance	60.2 (+10.9)	76.5 (+43.4)
5.	Oracle on ImageNet <sup>2</sup>	53.3 (+4.0)	<b>77.4</b>
6.	ROAD (Chen et al., 2018)	57.1 (+7.8)	<b>77.4</b>
7.	Vanilla L2 distance	56.4 (+7.1)	49.1
8.	SI (Zenke et al., 2017)	57.6 (+8.3)	53.9
9.	ASG (ours)	<b>61.5</b>	76.7

# ASG Improves Model Attention (GradCAM)

Input Image



Baseline



Skateboard ✗



Airplane ✓



Train ✗



Bus ✓



Motorcycle ✗



Horse ✓

# ASG Benefits Domain Adaptation & Self-Training

- ASG serves as better initialization

1. ImageNet → Self-training for DA

2. ImageNet → ASG → Self-training for DA

Method	Tgt Img	Accuracy
Source-Res101 (Zou et al., 2019)	✗	51.6
CBST (Zou et al., 2018)	✓	76.4 (0.9)
MRKLD (Zou et al., 2019)	✓	77.9 (0.5)
MRKLD + LRENT (Zou et al., 2019)	✓	78.1 (0.2)
ASG (ours)	✗	61.5
ASG + CBST	✓	82.5 (0.7)
ASG + MRKLD	✓	<b>84.6</b> (0.4)
ASG + MRKLD + LRENT	✓	84.5 (0.4)

## DISTRIBUTIONALLY ROBUST LEARNING FOR UNSUPERVISED DOMAIN ADAPTATION

**Haoxuan Wang** \*  
Shanghai Jiao Tong University  
hatch25@sjtu.edu.cn

**Anqi Liu** \*  
Caltech  
anqiliu@caltech.edu

**Zhidong Yu**  
NVIDIA  
zhidongy@nvidia.com

**Yisong Yue**  
Caltech  
yyue@caltech.edu

**Anima Anandkumar**  
Caltech & NVIDIA  
anima@caltech.edu  
aanandkumar@nvidia.com

Method	Mean
Source (Saito et al., 2018a)	52.4
MMD (Long et al., 2015)	61.1
MCD (Saito et al., 2018b)	71.9
ADR (Saito et al., 2018a)	74.8
CBST (Zou et al., 2020)	76.4
CRST (Zou et al., 2020)	78.1
AVH (Chen et al., 2020a)	81.5
<b>DRST (proposed)</b>	83.75
ASG (Chen et al., 2020b)	61.17
CBST-ASG (Chen et al., 2020b)	82.23
CRST-ASG (Chen et al., 2020b)	84.21
<b>DRST-ASG (proposed)</b>	<b>85.25</b>

# Summary

- NVIDIA: synthetic data powers vision applications.
- Domain gap: from synthetic to real data.
- ImageNet pre-training as proxy guidance.
- Automatic control of layer-wise learning rate.





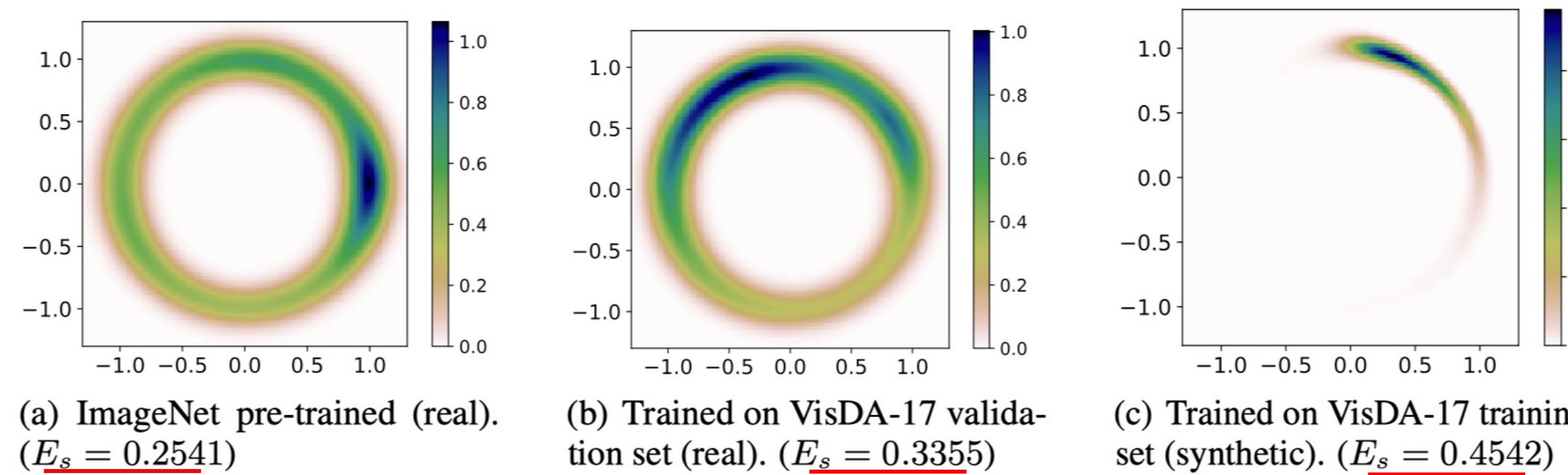
# Contrastive Syn-to-Real Generalization

ICLR 2021

Wuyang Chen, Zhiding Yu, Shalini De Mello, Sifei Liu, Jose M. Alvarez,  
Zhangyang (Atlas) Wang, Anima Anandkumar

# Deeper Look Into Domain Gap

- Synthetic images leads to collapsed feature space!

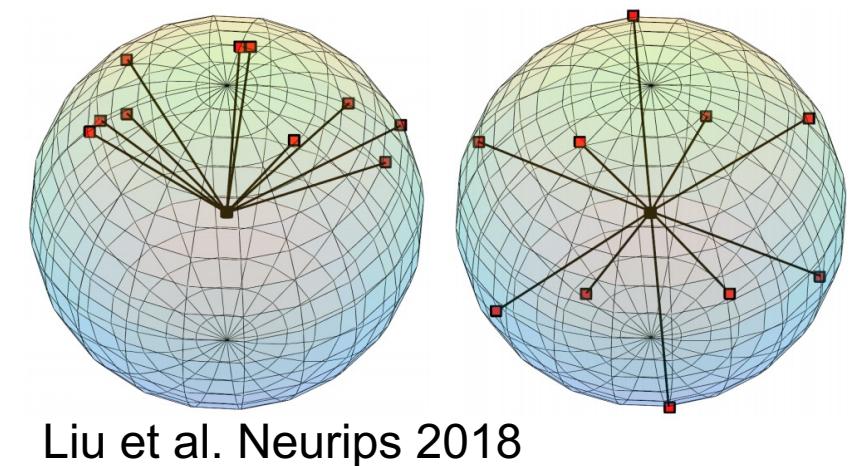


**Figure 2:** Feature diversity in  $\mathbb{R}^2$  with Gaussian kernel density estimation (KDE). Darker areas have more concentrated features.  $E_s$ : hyperspherical energy of features, lower the more diverse.

Hyperspherical Energy (HSE,  $E_s$ )

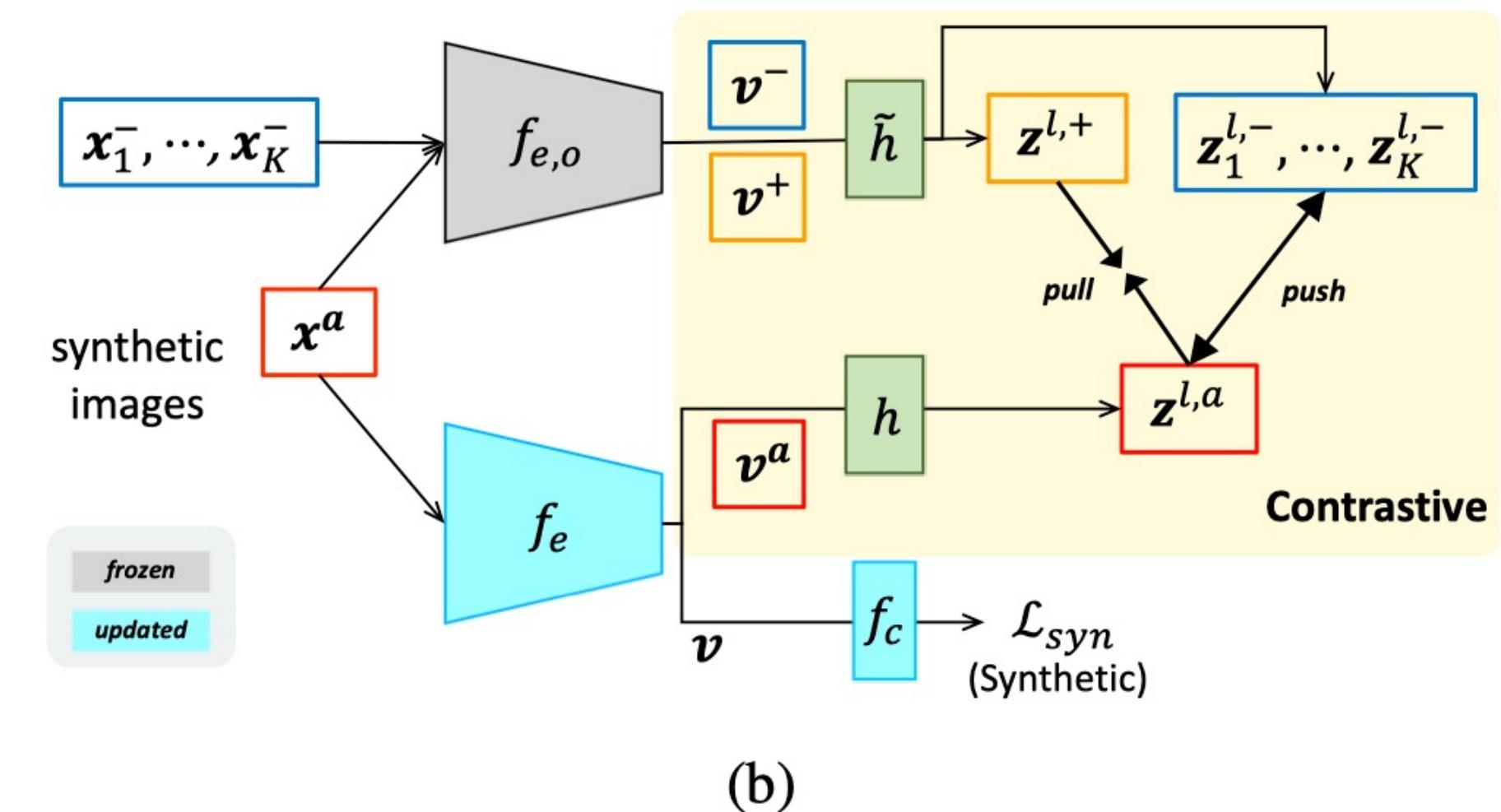
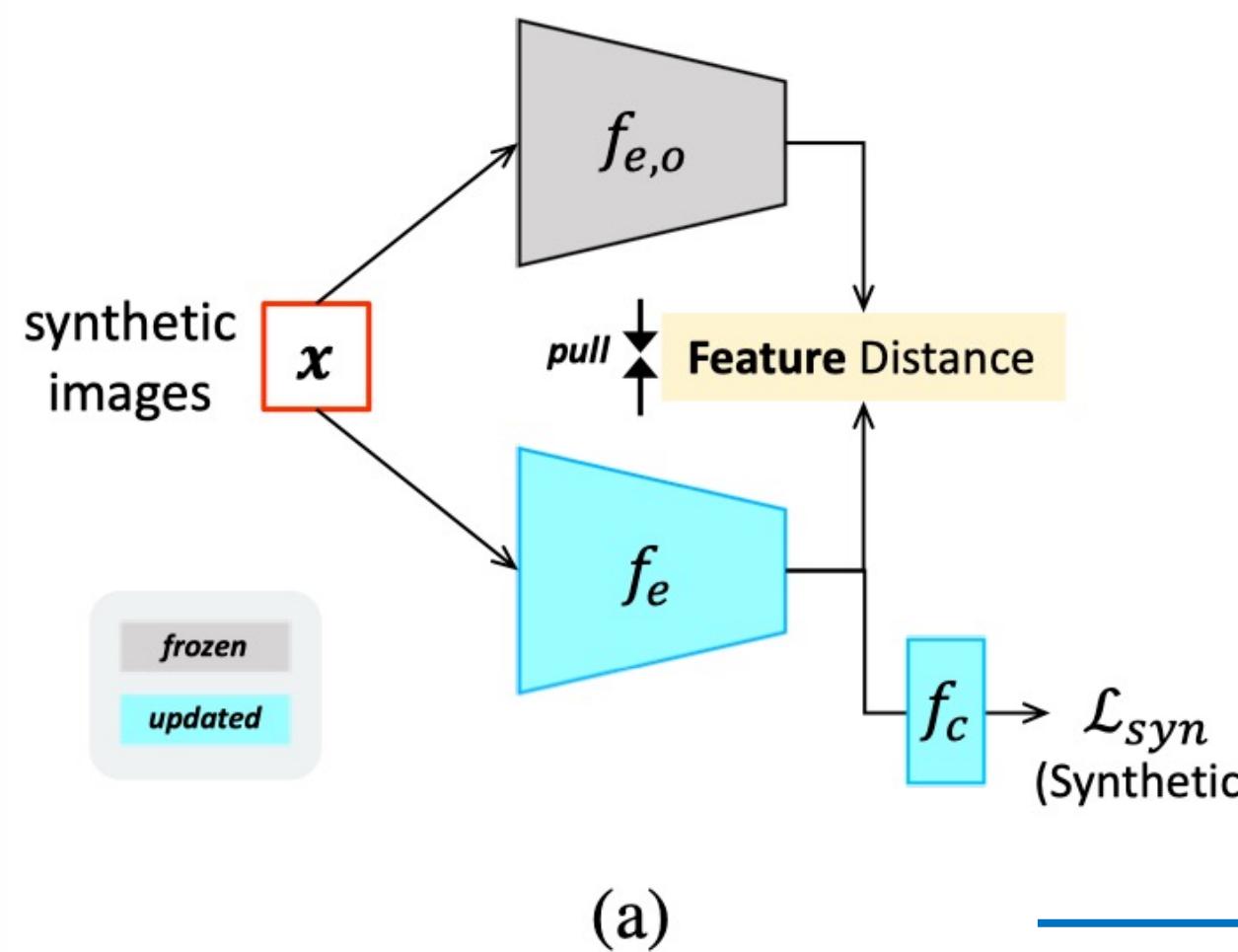
Low  $E_s \rightarrow$  diverse features

$$E_s \left( \bar{\mathbf{v}}_i |_{i=1}^N \right) = \sum_{i=1}^N \sum_{j=1, j \neq i}^N e_s (\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_j\|) = \begin{cases} \sum_{i \neq j} \|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_j\|^{-s}, & s > 0 \\ \sum_{i \neq j} \log (\|\bar{\mathbf{v}}_i - \bar{\mathbf{v}}_j\|^{-1}), & s = 0 \end{cases}$$



# ImageNet Distillation + Feature Diversity

- Synthetic-to-real with a “push and pull” strategy



# Contrastive Loss

- InfoNCE

$$\mathcal{L}_{\text{NCE}} = -\log \frac{\exp(\mathbf{z}^a \cdot \mathbf{z}^+ / \tau)}{\exp(\mathbf{z}^a \cdot \mathbf{z}^+ / \tau) + \sum_{\mathbf{z}^-} \exp(\mathbf{z}^a \cdot \mathbf{z}^- / \tau)}, \quad (3)$$

$$\mathcal{L} = \mathcal{L}_{\text{syn}} + \lambda \mathcal{L}_{\text{NCE}} \quad (4)$$

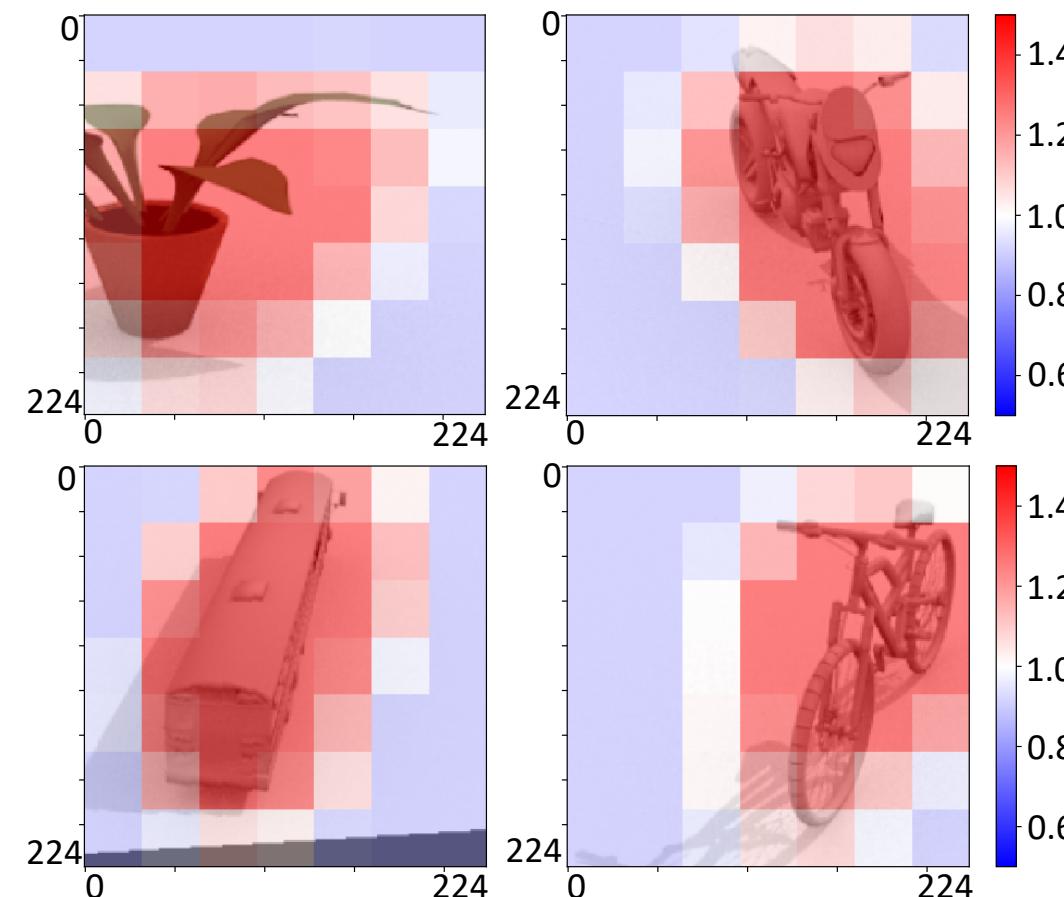
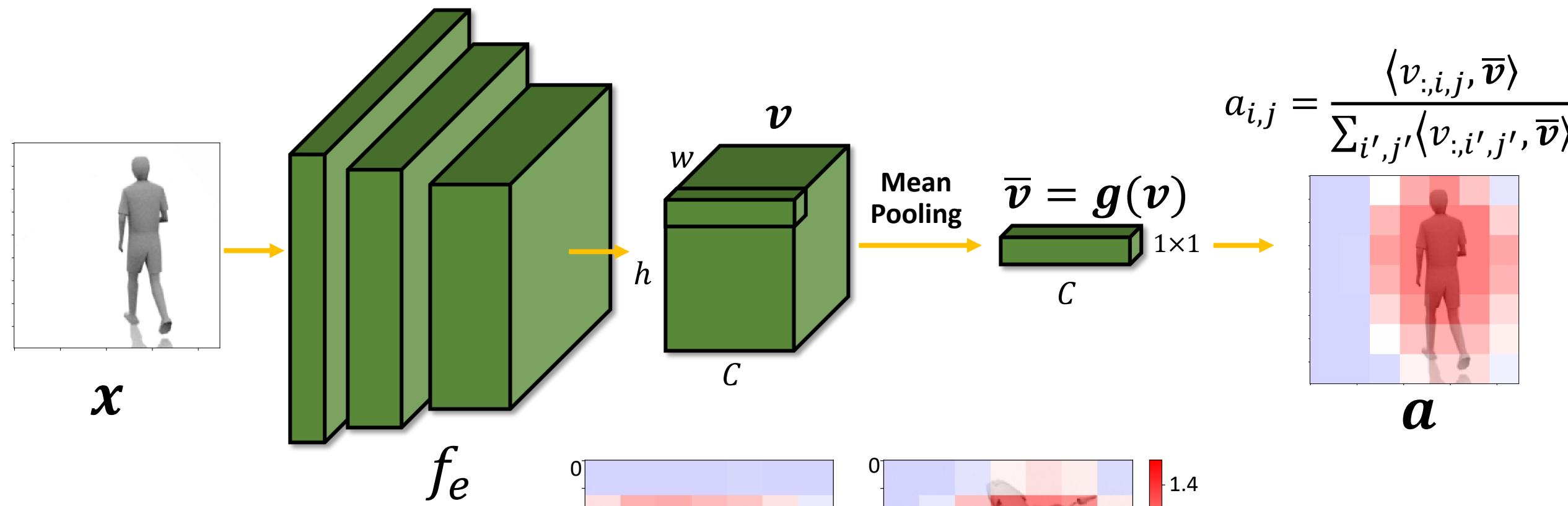
- Multi-layer InfoNCE

$$\mathcal{L}_{\text{NCE}} = \sum_{l \in \mathcal{G}} \mathcal{L}_{\text{NCE}}^l = \sum_{l \in \mathcal{G}} -\log \frac{\exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,+} / \tau)}{\exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,+} / \tau) + \sum_{\mathbf{z}^{l,-}} \exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,-} / \tau)} \quad (5)$$

- Dense InfoNCE (segmentation)

$$\mathcal{L}_{\text{NCE}} = \sum_{l \in \mathcal{G}} \sum_i \mathcal{L}_{\text{NCE}}^{l,i} = \sum_{l \in \mathcal{G}} \sum_i -\frac{1}{N_l} \log \frac{\exp(\mathbf{z}_i^{l,a} \cdot \mathbf{z}_i^{l,+} / \tau)}{\exp(\mathbf{z}_i^{l,a} \cdot \mathbf{z}_i^{l,+} / \tau) + \sum_{\mathbf{z}_i^{l,-}} \exp(\mathbf{z}_i^{l,a} \cdot \mathbf{z}_i^{l,-} / \tau)} \quad (6)$$

# Attention-guided Global Pooling



# Results: Feature Diversity vs Generalization

- Model preserves diverse features → generalize better on real domain

**Table 1:** Generalization performance and hyperspherical energy of the features extracted by different models (lower is better). Dataset: VisDA-17 (Peng et al., 2017) validation set. Model: ResNet-101.

Model	Power			Accuracy (%)
	0	1	2	
Oracle on ImageNet <sup>3</sup>	-	-	-	53.3
Baseline (vanilla synthetic training)	0.4245	1.2500	1.6028	49.3
Weight $l_2$ distance (Kirkpatrick et al., 2017)	0.4014	1.2296	1.5302	56.4
Synaptic Intelligence (Zenke et al., 2017)	0.3958	1.2261	1.5216	57.6
Feature $l_2$ distance (Chen et al., 2018)	0.3337	1.1910	1.4449	57.1
ASG (Chen et al., 2020b)	0.3251	1.1840	1.4229	61.1
CSG (Ours)	<b>0.3188</b>	<b>1.1806</b>	<b>1.4177</b>	<b>64.05</b>

# Results: Segmentation

## GTA5 → Cityscapes

**Table 5:** Comparison to prior domain generaliz

	road terrain	sidewalk sky	building person	wall rider	fence car	pole truck	traffic lgt bus	traffic sgn train	vegetation motorcycle	ignored bike
Methods										
No Adapt IBN-Net (Pan et al., 2018)										
No Adapt Yue et al. (Yue et al., 2019)										
No Adapt ASG (Chen et al., 2020b)										
No Adapt CSG (ours)										
No Adapt Yue et al. (Yue et al., 2019)										
No Adapt ASG (Chen et al., 2020b)			ResNet-101	27.94 32.79	4.85					
No Adapt CSG (ours)				28.94 38.88	<b>9.94</b>					

# Summary

- Synthetic data leads to collapsed learned features.
- Encouraging diverse features benefits synthetic-to-real generalization.
- Multi-layer & dense InforNCE + attention-guided pooling further boosts contrastive learning.

# Future Works

- More applications: Gaze, Detection, Robotics, etc.
- Joint training with domain adaptation.
- Better leveraging multiple sources
  - labeled real domain (ImageNet)
  - labeled synthetic domain
  - Unlabeled target real domain



The University of Texas at Austin  
**Electrical and Computer  
Engineering**  
*Cockrell School of Engineering*