

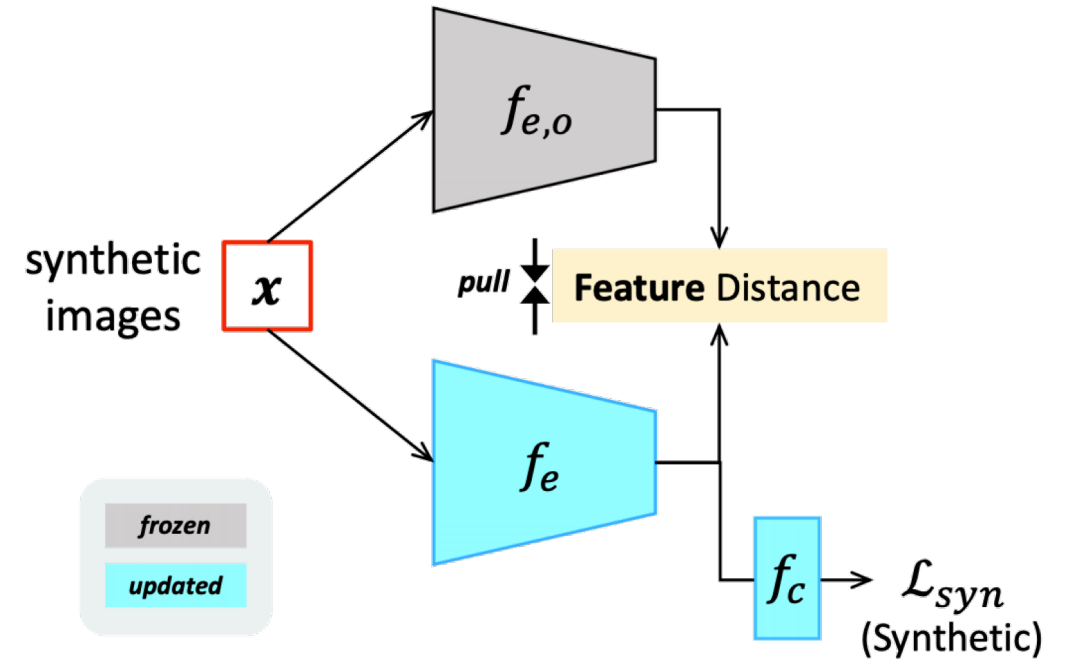
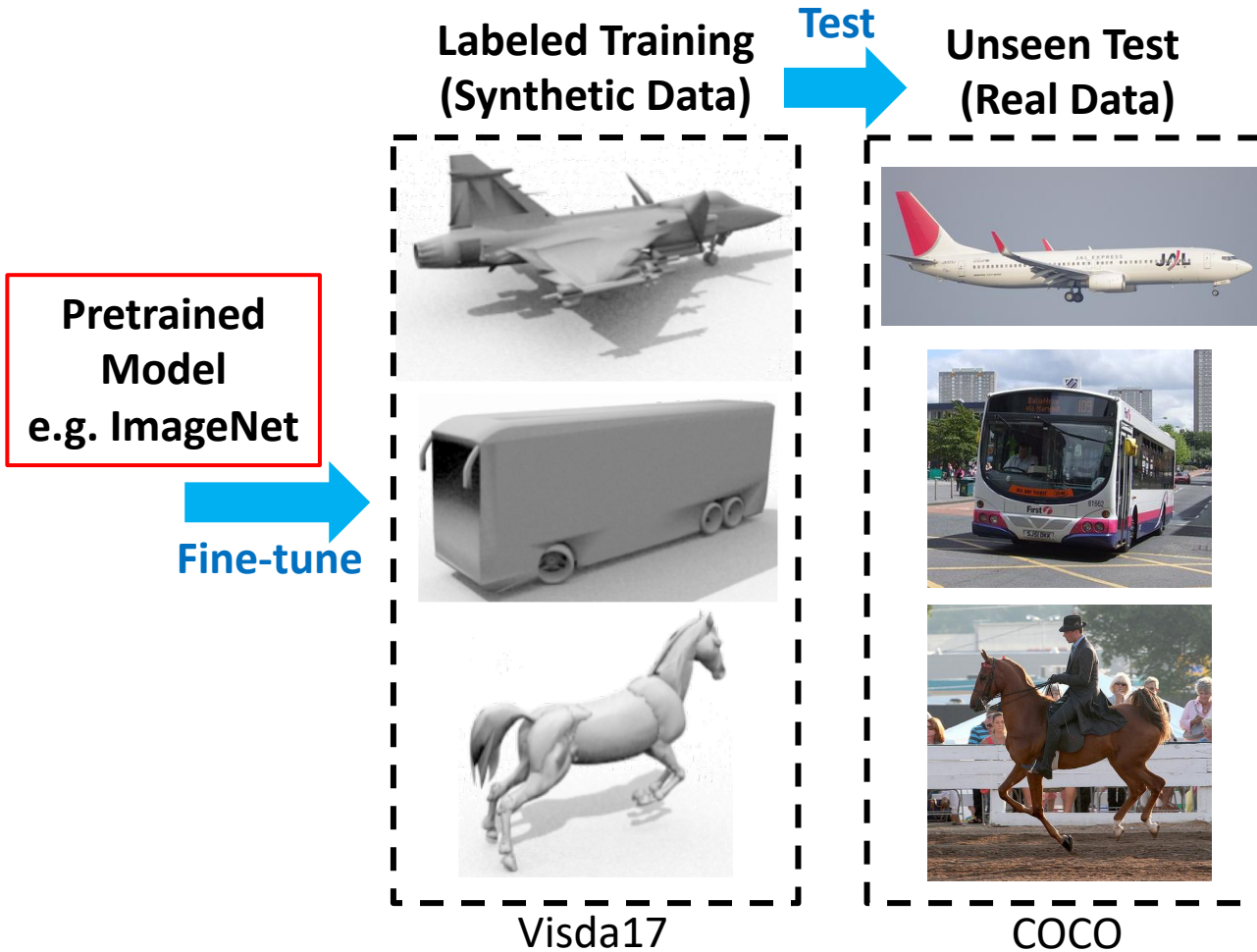
Contrastive Synthetic-to-Real Generalization

Wuyang Chen¹, Zhiding Yu², Shalini De Mello², Sifei Liu², Jose M. Alvarez²,
Zhangyang Wang¹, Animashree Anandkumar^{2,3}

¹UT Austin ²NVIDIA ³Caltech

* Work done during the author's research internship in NVIDIA Inc.

Syn-to-Real Generalization: Problem & Previous Solution

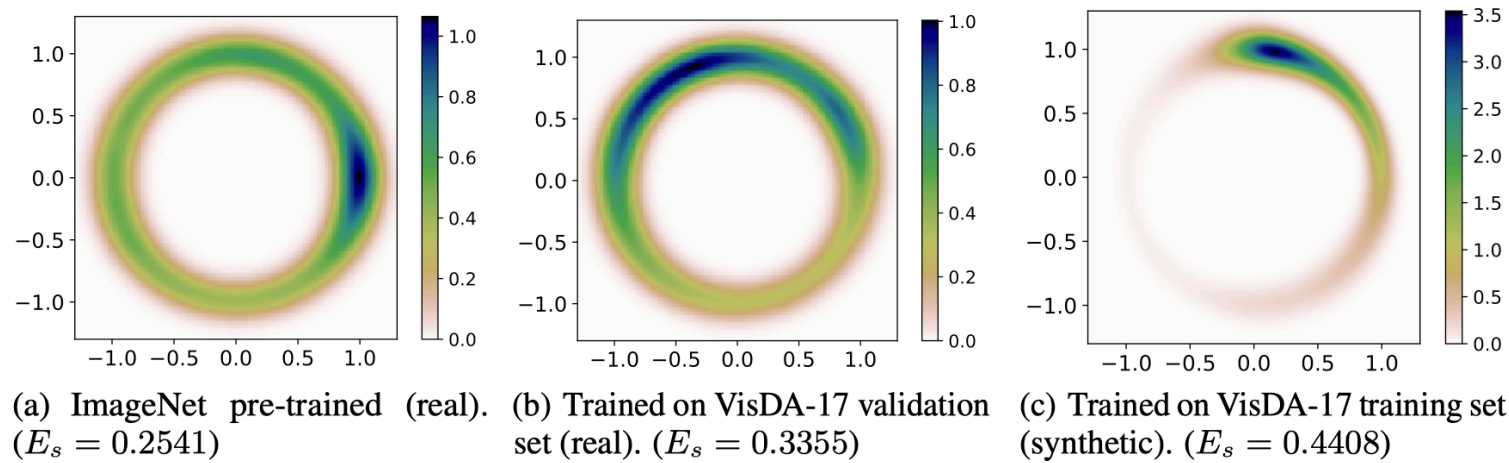


Chen, Wuyang, Zhiding Yu, Zhangyang Wang, and Animashree Anandkumar.
"Automated synthetic-to-real generalization." ICML 2020.

But why Synthetic Training Fails?

A Representation Learning Perspective

- Train model on natural images → diverse representations.
- Train model on **synthetic** images → **collapsed** representations!

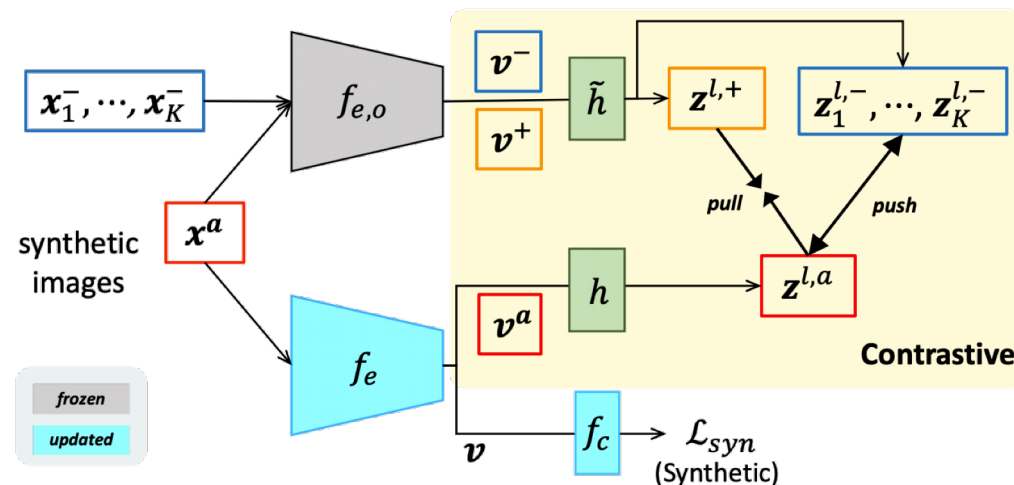


- E_s (Hyperspherical Energy): Lower the more diverse.

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

CSG: Contrastive Synthetic-to-Real Generalization

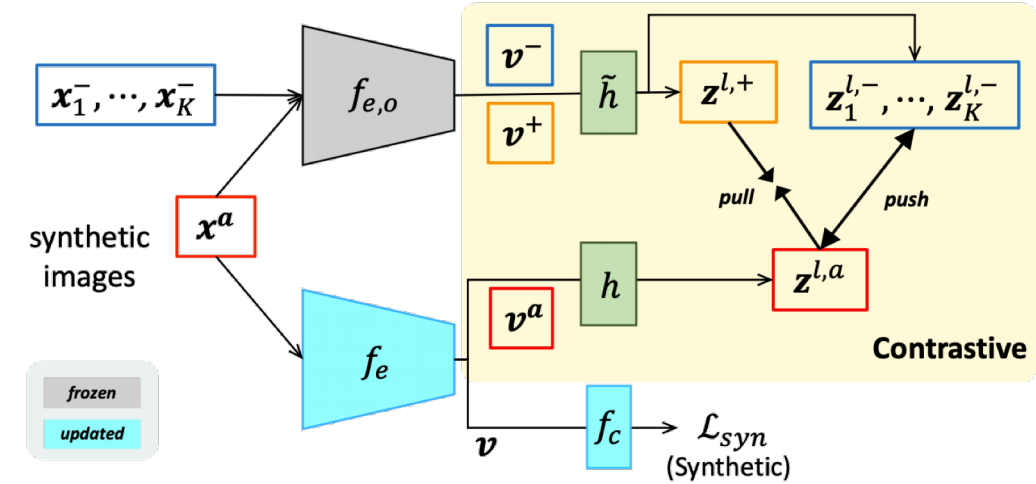
- How to transfer real domain knowledge + promote feature diversity?
- **Pull**: impose similarity b/w features from synthetic model v.s. ImageNet pre-trained model.
- **Push**: encourage feature diversity by pushing the feature embeddings away from each other across different images.



CSG: Contrastive Synthetic-to-Real Generalization

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

$$\mathcal{L} = \mathcal{L}_{\text{Task}} + \lambda \mathcal{L}_{\text{NCE}}$$



Multi-layer Contrastive Loss

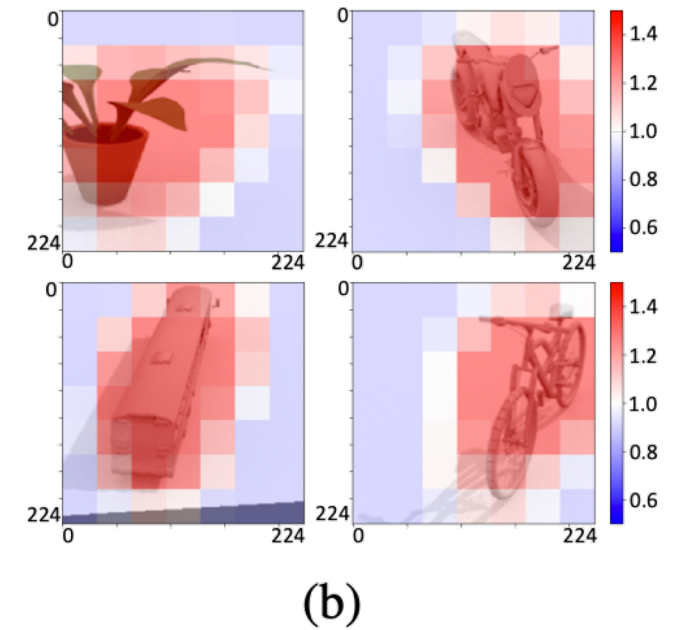
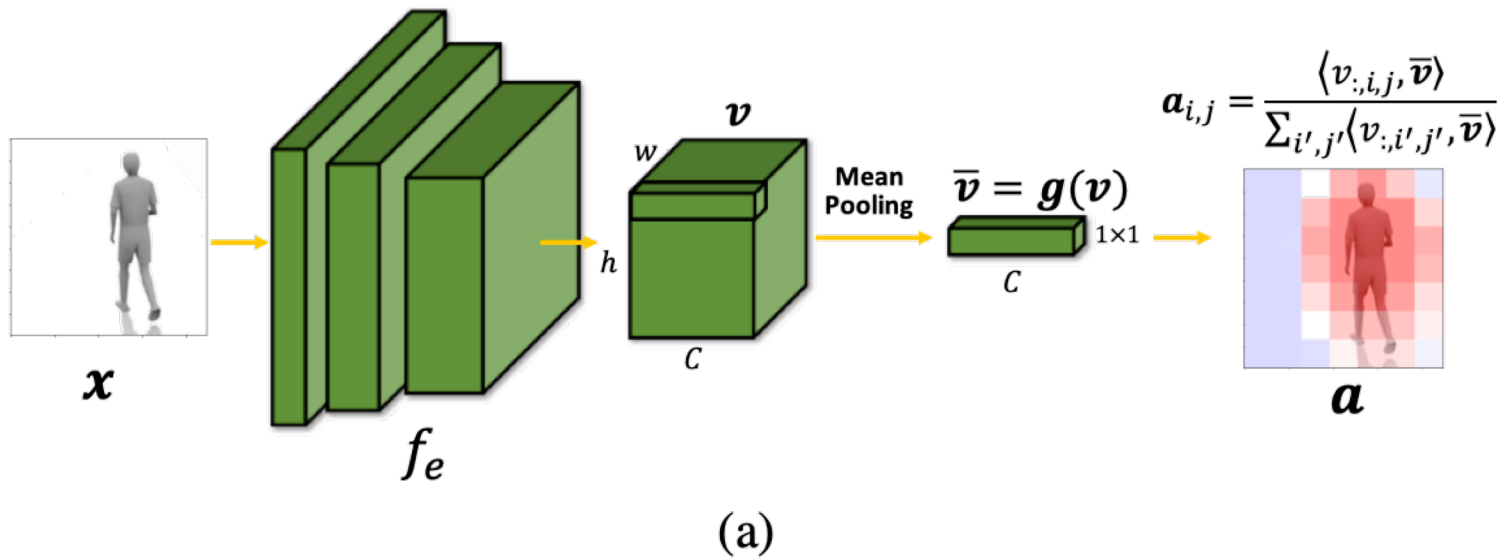
$$\mathcal{L}_{\text{NCE}} = \sum_{l \in \mathcal{G}} \mathcal{L}_{\text{NCE}}^l = \sum_{l \in \mathcal{G}} -\log \frac{\exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,+} / \tau)}{\exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,+} / \tau) + \sum_{\mathbf{z}^{l,-}} \exp(\mathbf{z}^{l,a} \cdot \mathbf{z}^{l,-} / \tau)}$$

Dense Contrastive Loss for Segmentation

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

A-Pool: Attentional Pooling for Improved Representation

- A-Pool for non-linear projection head $h(\cdot)$ focus contrastive learning on more semantically meaningful regions.



CSG: Best Performance Gain & Diverse Features

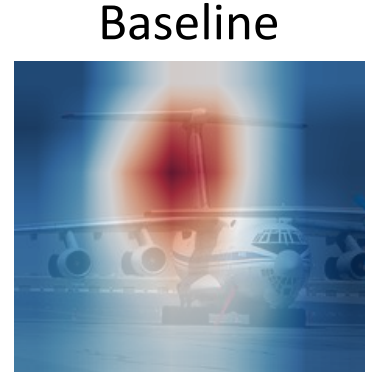
Classification: Visda17 → COCO

Model	Power			Accuracy (%)
	0	1	2	
Oracle on ImageNet ³	-	-	-	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., 2020c)	0.3251	1.1840	1.4229	61.1
CSG (Ours)	0.3188	1.1806	1.4177	64.05

Segmentation: GTA5 → Cityscapes

Methods	Backbone	mIoU %	mIoU ↑ %
No Adapt IBN-Net (Pan et al., 2018)	ResNet-50	22.17 29.64	7.47
No Adapt Yue et al. (Yue et al., 2019)		32.45 37.42	4.97
No Adapt ASG (Chen et al., 2020c)		25.88 29.65	3.77
No Adapt CSG (ours)		25.88 35.27	9.39
No Adapt Yue et al. (Yue et al., 2019)	ResNet-101	33.56 42.53	8.97
No Adapt ASG (Chen et al., 2020c)		29.63 32.79	3.16
No Adapt CSG (ours)		29.63 38.88	9.25

CSG Improves Model Attention (GradCAM)



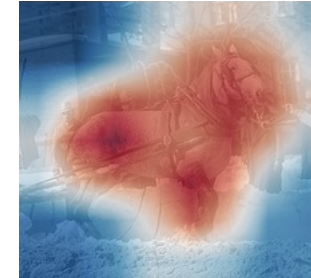
Train ✗

Airplane ✓



Motorcycle ✗

Bicycle ✓

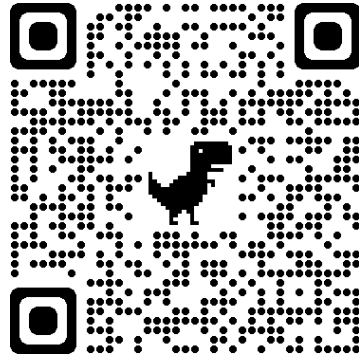


Plant ✗

Horse ✓

Thank you!

Code



Paper

