

# Progressive Transformation Learning for Leveraging Virtual Images in Training

Yi-Ting Shen<sup>\*,1</sup>

Hyungtae Lee<sup>\*,2</sup>

Heesung Kwon<sup>2</sup>

Shuvra S. Bhattacharyya<sup>1</sup>

<sup>\*</sup> equal contribution

<sup>1</sup> University of Maryland, College Park

<sup>2</sup> DEVCOM Army Research Laboratory

**CVPR 2023 HIGHLIGHT**

# Content

- Authors
- Background
- Method
- Experiments

# Content

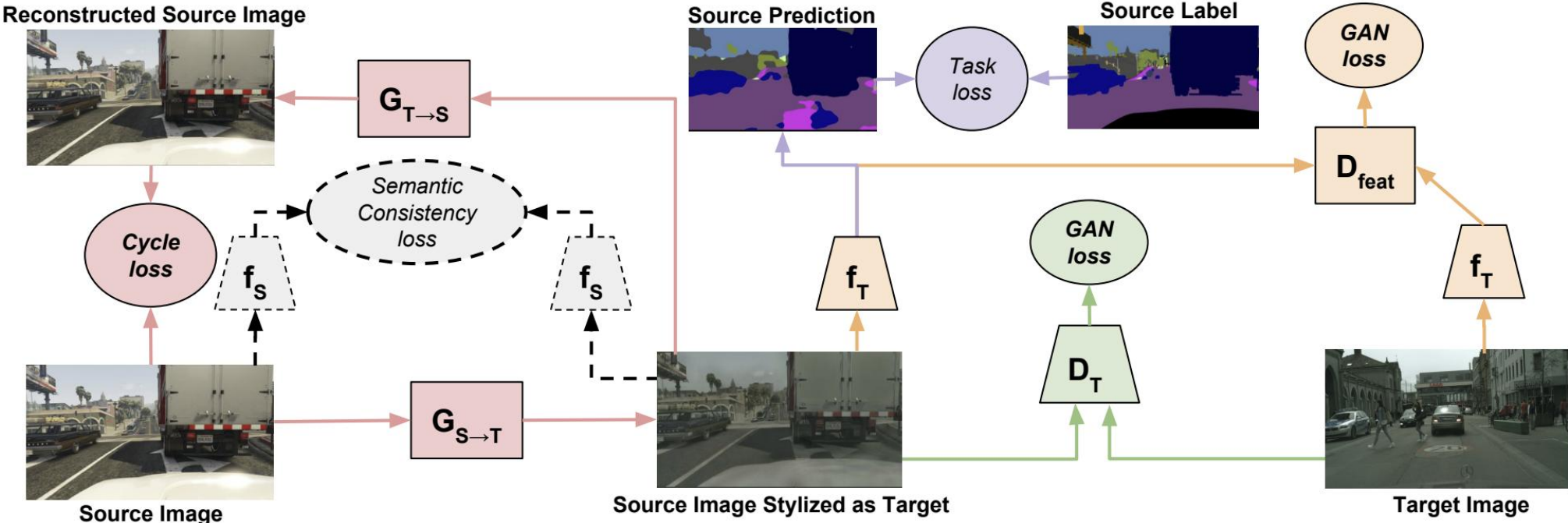
- Authors
- **Background**
- Method
- Experiments

# Background

CyCADA

# Background

## CyCADA



# Background

## CyCADA

$$\begin{aligned} \mathcal{L}_{\text{CyCADA}}(f_T, X_S, X_T, Y_S, G_{S \rightarrow T}, G_{T \rightarrow S}, D_S, D_T) \\ &= \mathcal{L}_{\text{task}}(f_T, G_{S \rightarrow T}(X_S), Y_S) \\ &+ \mathcal{L}_{\text{GAN}}(G_{S \rightarrow T}, D_T, X_T, X_S) + \mathcal{L}_{\text{GAN}}(G_{T \rightarrow S}, D_S, X_S, X_T) \\ &+ \mathcal{L}_{\text{GAN}}(f_T, D_{\text{feat}}, f_S(G_{S \rightarrow T}(X_S)), X_T) \\ &+ \mathcal{L}_{\text{cyc}}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T) + \mathcal{L}_{\text{sem}}(G_{S \rightarrow T}, G_{T \rightarrow S}, X_S, X_T, f_S). \end{aligned}$$

# Background

## Deep Mahalanobis detector

- Mahalanobis distance:

- Single sample:  $D_M(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$

- Two samples:  $D_M(x, y) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$

# Background

## Deep Mahalanobis detector

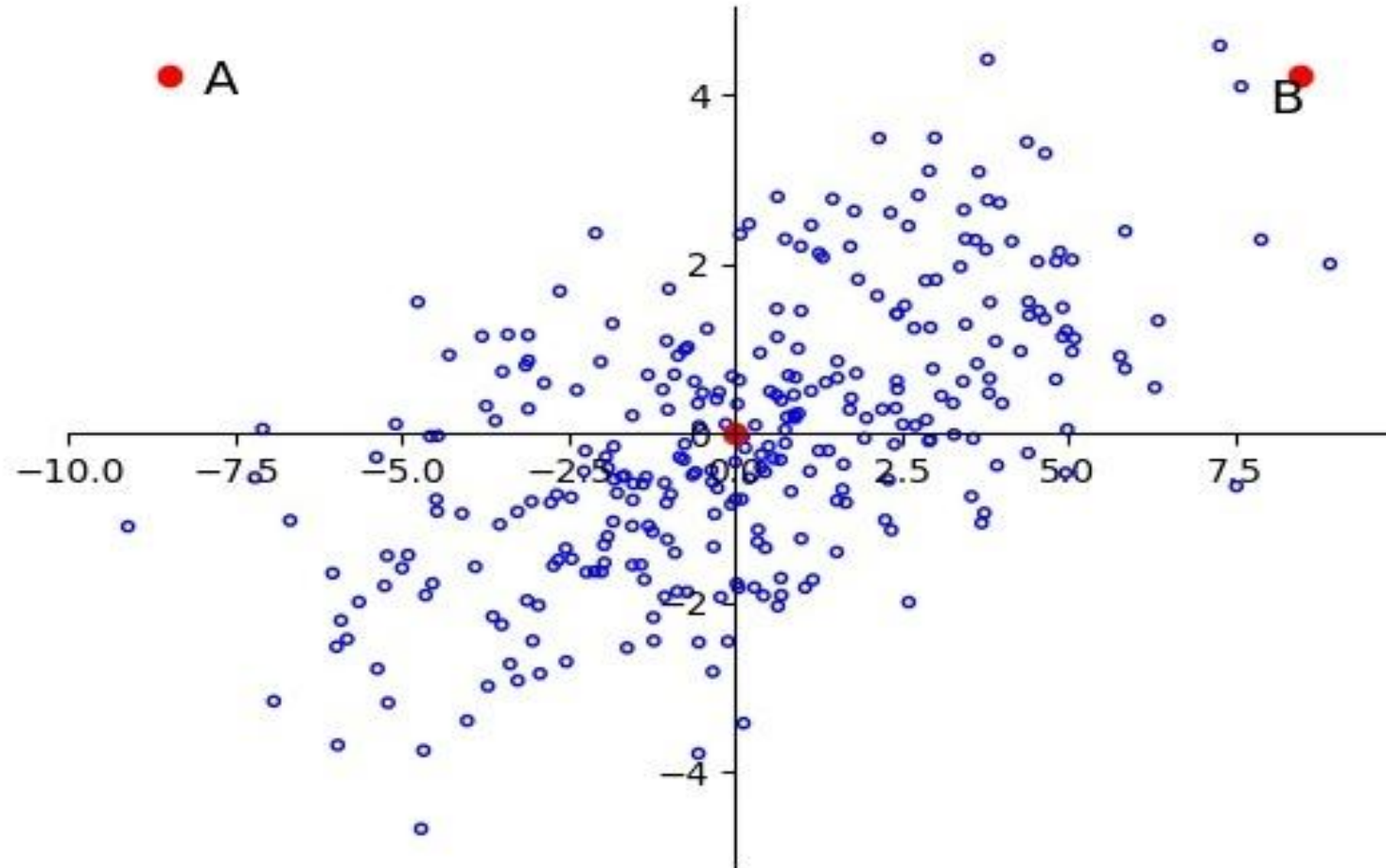
- Mahalanobis distance:

- Single sample:  $D_M(x) = \sqrt{(x - \mu)^T \overset{\text{covariance matrix}}{\Sigma^{-1}} (x - \mu)}$  mean of samples

- Two samples:  $D_M(x, y) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$



# Background



# Background

- Mahalanobis distance:  $Z = Y \begin{bmatrix} \frac{1}{\sqrt{\sigma_1}} & & \\ & \ddots & \\ & & \frac{1}{\sqrt{\sigma_F}} \end{bmatrix} \Leftrightarrow Y\Lambda \quad \Sigma = U^T Q U \quad Y = XU$

$$\begin{aligned}
 \text{Dis}_{\text{Euclidean}}(z_1, z_2) &= \sqrt{(z_1 - z_2)(z_1 - z_2)^T} \\
 &= \sqrt{(y_1 \Lambda - y_2 \Lambda)(y_1 \Lambda - y_2 \Lambda)^T} \\
 &= \sqrt{(x_1 U \Lambda - x_2 U \Lambda)(x_1 U \Lambda - x_2 U \Lambda)^T} \\
 &= \sqrt{(x_1 - x_2)U \Lambda \Lambda U^T (x_1 - x_2)^T} \\
 &= \sqrt{(x_1 - x_2)\Sigma^{-1}(x_1 - x_2)^T} \\
 &= \text{Dis}_{\text{mahalanobis}}(x_1, x_2)
 \end{aligned}$$

# Background

- Mahalanobis distance:  $Z = Y \begin{bmatrix} \frac{1}{\sqrt{\sigma_1}} & & \\ & \ddots & \\ & & \frac{1}{\sqrt{\sigma_F}} \end{bmatrix} \Leftrightarrow Y\Lambda \quad \Sigma = U^T Q U \quad Y = XU$

$$\begin{aligned} \text{Dis}_{\text{Euclidean}}(z_1, z_2) &= \sqrt{(z_1 - z_2)(z_1 - z_2)^T} \\ &= \sqrt{(y_1 \Lambda - y_2 \Lambda)(y_1 \Lambda - y_2 \Lambda)^T} \\ &= \sqrt{(x_1 U \Lambda - x_2 U \Lambda)(x_1 U \Lambda - x_2 U \Lambda)^T} \\ &= \sqrt{(x_1 - x_2)U \Lambda \Lambda U^T (x_1 - x_2)^T} \\ &= \sqrt{(x_1 - x_2)\Sigma^{-1}(x_1 - x_2)^T} \\ &= \text{Dis}_{\text{mahalanobis}}(x_1, x_2) \end{aligned}$$

**Hypothesis:** The features of each layer of convolutional neural networks follow a multivariate Gaussian distribution

# Background

We can calculate the mean value and the covariance matrix of all samples of each category in each convolution layer.

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} f(x_i)$$

$$\hat{\Sigma} = \frac{1}{N} \sum_c \sum_{i:y_i=c} (f(x_i) - \hat{\mu}_c)(f(x_i) - \hat{\mu}_c)^T$$

For a new sample  $x$ , its degree of category confidence:

$$M(x) = \max_c -(f(x) - \hat{\mu}_c)^T \hat{\Sigma}^{-1} (f(x) - \hat{\mu}_c)$$

# Background

## Deep Mahalanobis detector

In-dist (model)	OOD	Validation on OOD samples			Validation on adversarial samples		
		TNR at TPR 95%	AUROC	Detection acc.	TNR at TPR 95%	AUROC	Detection acc.
		Baseline [13] / ODIN [21] / Mahalanobis (ours)			Baseline [13] / ODIN [21] / Mahalanobis (ours)		
CIFAR-10 (DenseNet)	SVHN	40.2 / 86.2 / <b>90.8</b>	89.9 / 95.5 / <b>98.1</b>	83.2 / 91.4 / <b>93.9</b>	40.2 / 70.5 / <b>89.6</b>	89.9 / 92.8 / <b>97.6</b>	83.2 / 86.5 / <b>92.6</b>
	TinyImageNet	58.9 / 92.4 / <b>95.0</b>	94.1 / 98.5 / <b>98.8</b>	88.5 / 93.9 / <b>95.0</b>	58.9 / 87.1 / <b>94.9</b>	94.1 / 97.2 / <b>98.8</b>	88.5 / 92.1 / <b>95.0</b>
	LSUN	66.6 / 96.2 / <b>97.2</b>	95.4 / 99.2 / <b>99.3</b>	90.3 / 95.7 / <b>96.3</b>	66.6 / 92.9 / <b>97.2</b>	95.4 / 98.5 / <b>99.2</b>	90.3 / 94.3 / <b>96.2</b>
CIFAR-100 (DenseNet)	SVHN	26.7 / 70.6 / <b>82.5</b>	82.7 / 93.8 / <b>97.2</b>	75.6 / 86.6 / <b>91.5</b>	26.7 / 39.8 / <b>62.2</b>	82.7 / 88.2 / <b>91.8</b>	75.6 / 80.7 / <b>84.6</b>
	TinyImageNet	17.6 / 42.6 / <b>86.6</b>	71.7 / 85.2 / <b>97.4</b>	65.7 / 77.0 / <b>92.2</b>	17.6 / 43.2 / <b>87.2</b>	71.7 / 85.3 / <b>97.0</b>	65.7 / 77.2 / <b>91.8</b>
	LSUN	16.7 / 41.2 / <b>91.4</b>	70.8 / 85.5 / <b>98.0</b>	64.9 / 77.1 / <b>93.9</b>	16.7 / 42.1 / <b>91.4</b>	70.8 / 85.7 / <b>97.9</b>	64.9 / 77.3 / <b>93.8</b>
SVHN (DenseNet)	CIFAR-10	69.3 / 71.7 / <b>96.8</b>	91.9 / 91.4 / <b>98.9</b>	86.6 / 85.8 / <b>95.9</b>	69.3 / 69.3 / <b>97.5</b>	91.9 / 91.9 / <b>98.8</b>	86.6 / 86.6 / <b>96.3</b>
	TinyImageNet	79.8 / 84.1 / <b>99.9</b>	94.8 / 95.1 / <b>99.9</b>	90.2 / 90.4 / <b>98.9</b>	79.8 / 79.8 / <b>99.9</b>	94.8 / 94.8 / <b>99.8</b>	90.2 / 90.2 / <b>98.9</b>
	LSUN	77.1 / 81.1 / <b>100</b>	94.1 / 94.5 / <b>99.9</b>	89.1 / 89.2 / <b>99.3</b>	77.1 / 77.1 / <b>100</b>	94.1 / 94.1 / <b>99.9</b>	89.1 / 89.1 / <b>99.2</b>
CIFAR-10 (ResNet)	SVHN	32.5 / 86.6 / <b>96.4</b>	89.9 / 96.7 / <b>99.1</b>	85.1 / 91.1 / <b>95.8</b>	32.5 / 40.3 / <b>75.8</b>	89.9 / 86.5 / <b>95.5</b>	85.1 / 77.8 / <b>89.1</b>
	TinyImageNet	44.7 / 72.5 / <b>97.1</b>	91.0 / 94.0 / <b>99.5</b>	85.1 / 86.5 / <b>96.3</b>	44.7 / 69.6 / <b>95.5</b>	91.0 / 93.9 / <b>99.0</b>	85.1 / 86.0 / <b>95.4</b>
	LSUN	45.4 / 73.8 / <b>98.9</b>	91.0 / 94.1 / <b>99.7</b>	85.3 / 86.7 / <b>97.7</b>	45.4 / 70.0 / <b>98.1</b>	91.0 / 93.7 / <b>99.5</b>	85.3 / 85.8 / <b>97.2</b>
CIFAR-100 (ResNet)	SVHN	20.3 / 62.7 / <b>91.9</b>	79.5 / 93.9 / <b>98.4</b>	73.2 / 88.0 / <b>93.7</b>	20.3 / 12.2 / <b>41.9</b>	79.5 / 72.0 / <b>84.4</b>	73.2 / 67.7 / <b>76.5</b>
	TinyImageNet	20.4 / 49.2 / <b>90.9</b>	77.2 / 87.6 / <b>98.2</b>	70.8 / 80.1 / <b>93.3</b>	20.4 / 33.5 / <b>70.3</b>	77.2 / 83.6 / <b>87.9</b>	70.8 / 75.9 / <b>84.6</b>
	LSUN	18.8 / 45.6 / <b>90.9</b>	75.8 / 85.6 / <b>98.2</b>	69.9 / 78.3 / <b>93.5</b>	18.8 / 31.6 / <b>56.6</b>	75.8 / 81.9 / <b>82.3</b>	69.9 / 74.6 / <b>79.7</b>
SVHN (ResNet)	CIFAR-10	78.3 / 79.8 / <b>98.4</b>	92.9 / 92.1 / <b>99.3</b>	90.0 / 89.4 / <b>96.9</b>	78.3 / 79.8 / <b>94.1</b>	92.9 / 92.1 / <b>97.6</b>	90.0 / 89.4 / <b>94.6</b>
	TinyImageNet	79.0 / 82.1 / <b>99.9</b>	93.5 / 92.0 / <b>99.9</b>	90.4 / 89.4 / <b>99.1</b>	79.0 / 80.5 / <b>99.2</b>	93.5 / 92.9 / <b>99.3</b>	90.4 / 90.1 / <b>98.8</b>
	LSUN	74.3 / 77.3 / <b>99.9</b>	91.6 / 89.4 / <b>99.9</b>	89.0 / 87.2 / <b>99.5</b>	74.3 / 76.3 / <b>99.9</b>	91.6 / 90.7 / <b>99.9</b>	89.0 / 88.2 / <b>99.5</b>

Table 2: Distinguishing in- and out-of-distribution test set data for image classification under various validation setups. All values are percentages and the best results are indicated in bold.

# Content

- Authors
- Background
- **Method**
- Experiments

# Motivation

Unmanned aerial vehicle (UAV)  $\longleftrightarrow$  Ground-based cameras

Wide variety of camera viewing angles

Expand the training set  $\longrightarrow$  Domain gap

# Motivation

Unmanned aerial vehicle (UAV)  $\longleftrightarrow$  Ground-based cameras

Wide variety of camera viewing angles

Expand the training set  $\longrightarrow$  Domain gap

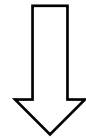
How to measure accurately the domain gap?



# Method

Softmax-based classifier: representation space

The distribution of each category——multivariate Gaussian distribution.



Sigmoid-based detector: representation space

The distribution of each category——multivariate Gaussian distribution.

# Method

- The representation space of the sigmoid-based detector:

$$P(f(\mathbf{x})|y_c = 1) \sim \mathcal{N}(f(\mathbf{x})|\mu_c, \Sigma_c), \quad (1)$$

where

$$\begin{aligned} \mu_c &= \frac{1}{|\mathbf{D}_c|} \sum_{\mathbf{x} \in \mathbf{D}_c} f(\mathbf{x}), \\ \Sigma_c &= \frac{1}{|\mathbf{D}_c|} \sum_{\mathbf{x} \in \mathbf{D}_c} (f(\mathbf{x}) - \mu_c) (f(\mathbf{x}) - \mu_c)^\top, \quad (2) \end{aligned}$$

# Method

- The domain gap between a new image  $\mathbf{x}_{new}$  and  $D_c$ :

$$d(\mathbf{x}_{new}) = (f(\mathbf{x}_{new}) - \mu_c)^\top \Sigma_c^{-1} (f(\mathbf{x}_{new}) - \mu_c). \quad (3)$$

To mitigate the effect of image size on measuring the domain gap—  
—multiple image scales:

$$d(\mathbf{x}_{new}) = \min_{s \in S} (\{d(\mathbf{x}_{new}^s)\}), \quad (4)$$

where  $S = \{128, 256, 384, 512\}$ .

# Method

Progressive Transformation Learning:

- Transformation candidate selection

$$w(\mathbf{x}) = \exp\left(-\frac{d(\mathbf{x})}{\tau}\right), \quad (5)$$

# Method

Progressive Transformation Learning:

- Virtual2Real transformation:  
Crop the person region——apply the transformation
- CycleGAN

# Method

$$\text{Set update: } \mathbf{R}^{t+1} = \mathbf{R}^t \cup \mathbf{C}_R^t \quad \mathbf{V}^{t+1} = \mathbf{V}^t / \mathbf{C}_V^t \quad (6)$$

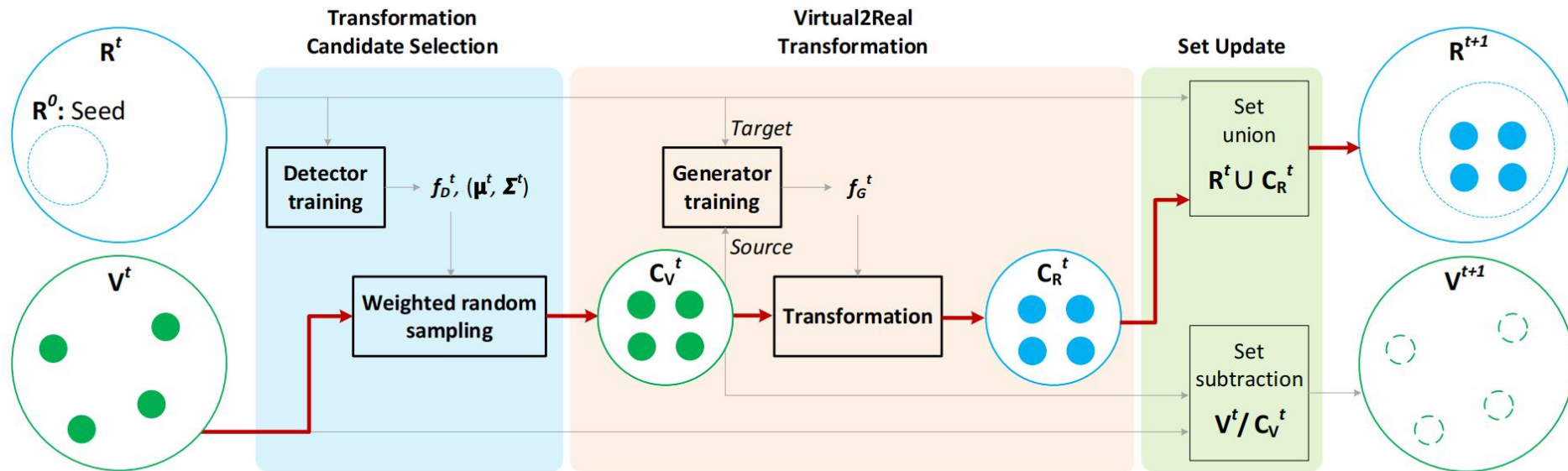


Figure 2. **Progressive Transformation Learning (PTL) pipeline.** The red arrow indicates the processing flow of the virtual images selected to be added to the training set.

# Experiments

Real UAV-based datasets:

- VisDrone
- Okutama Action
- ICG



Virtual dataset

- Archangel Synthetic

Metrics

- $AP@.5$
- $AP@[.5:.95]$



# Experiments

## Detector—— RetinaNet

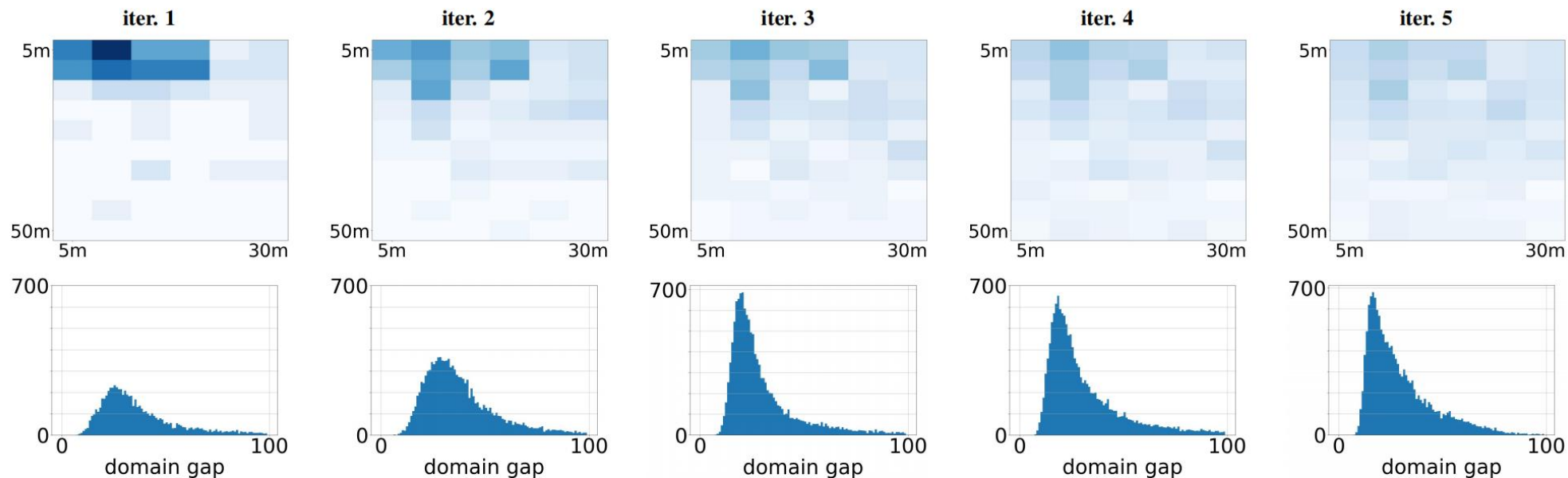


Figure 3. **Analysis of the use of virtual images when PTL progresses.** The figures in the top row show the accumulated distribution of transformation candidates with respect to camera locations (i.e., altitude and rotation circle radius from the target human in  $x$  and  $y$  axes, respectively) for each PTL iteration. Darker bins indicate that more virtual images have been added to the training set. The figures in the bottom row ( $x$  axis: domain gap,  $y$  axis: the corresponding number of virtual images) show the domain gap distribution of virtual images measured by eq. 4. These figures are collected from the experimental setup of using 50 real images from the VisDrone dataset for training.



# Experiments

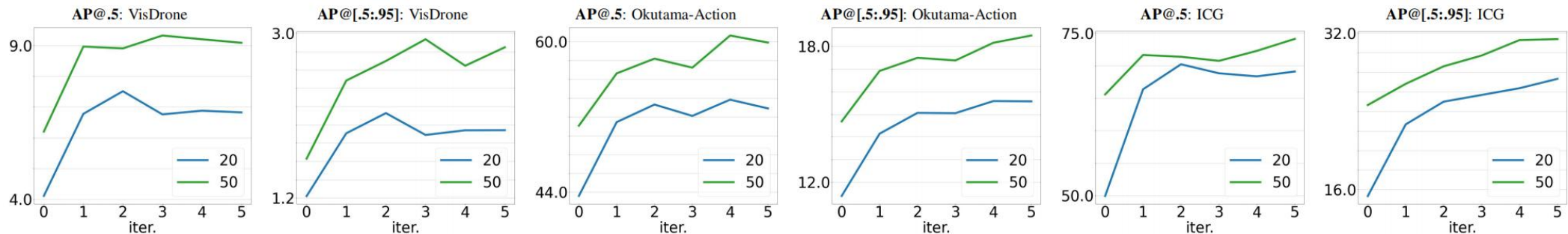


Figure 4. **Learning curves** of the two metrics (AP@.5 and AP@[.5:.95]) on the three datasets

method	train set	VisDrone		Okutama-Action		ICG	
		20	50	20	50	20	50
baseline	R	3.74/ 1.09	6.42/ 1.86	41.61/ 11.23	49.84/ 13.76	49.35/ 14.69	66.75/ 23.91
pretrain-finetune	R+V	4.99/ 1.46	6.25/ 1.99	44.57/ 12.78	49.06/ 15.08	66.92/ 26.67	68.41/ 29.73
naive merge	R+V	3.41/ 1.02	5.18/ 1.65	34.26/ 9.21	48.33/ 14.61	55.95/ 20.76	65.68/ 26.73
w/ transform	R+V	1.26/ 0.49	4.02/ 1.37	27.37/ 7.84	41.36/ 12.64	48.02/ 17.62	65.03/ 27.21
<b>PTL (5th itr.)</b>	R+V	6.83/ 1.94	9.09/ 2.85	52.89/ 15.57	59.90/ <b>18.48</b>	69.11/ <b>27.33</b>	<b>74.14/ 31.41</b>
		+3.09/+0.85	+2.67/+0.99	+11.28/ +4.34	+10.06/ +4.72	+19.76/+12.64	+7.39/ +7.50
<b>PTL (best)</b>	R+V	<b>7.52/ 2.13</b>	<b>9.33/ 2.94</b>	<b>53.82/ 15.59</b>	<b>60.65/ 18.48</b>	<b>70.23/ 27.33</b>	<b>74.14/ 31.41</b>
		+3.78/+1.04	+2.91/+1.08	+12.21/ +4.36	+10.81/ +4.72	+20.88/+12.64	+7.39/ +7.50

Table 1. **Low-shot learning accuracy** with 20 and 50 real images. AP@.5 and AP@[.5:.95] are reported in each bin. For PTL, the margin from the baseline accuracy is shown below the reported accuracy. The best accuracy for each setting is shown in bold. R and V denote the set of real images and the set of virtual images, respectively.

# Experiments

$\tau$	VisDrone	Okutama-Action	ICG
1	9.48/ 3.01	39.37/ 10.45	27.87/ 7.75
5	9.09/ 2.85	<b>42.39/ 11.41</b>	29.26/ 7.27
10	<b>9.68/ 2.87</b>	37.48/ 9.51	33.06/ 7.66
100	8.97/ 2.54	38.25/ 10.18	33.78/ 9.29
1000	8.90/ 2.63	39.15/ 10.00	<b>43.90/ 11.97</b>

Table 2. **Varying  $\tau$**  in PTL (after 5th iteration). Models are trained on the VisDrone dataset with 50 shot learning setup.

# img per itr	VisDrone	Okutama-Action	ICG
50	<b>9.59/ 2.90</b>	41.62/ 11.19	25.94/ 6.04
100	9.33/ <b>2.94</b>	<b>42.39/ 11.46</b>	30.01/ 7.36
200	9.04/ 2.91	41.29/ 11.18	<b>35.50/ 10.20</b>

Table 3. **Varying # of virtual images added to the training set per PTL iteration.** Models are trained on the VisDrone dataset with 50 shot learning setup. The reported accuracies are obtained by using the best PTL models.



Figure 5. **Sample Virtual2Real transformation output (VisDrone, 50-shot).** Each set consists of three images: original virtual image (left), transformed image (middle), and transformed image with background (right).



# Experiments

method	VisDrone				Okutama-Action				ICG			
	20		50		20		50		20		50	
	<i>Vis</i> → <i>Vis</i>				<i>Oku</i> → <i>Oku</i>				<i>ICG</i> → <i>ICG</i>			
baseline	3.74/	1.09	6.42/	1.86	41.61/	11.23	49.84/	13.76	49.35/	14.69	66.75/	23.91
PTL (5th itr.)	6.83/	1.94	9.09/	2.85	52.89/	15.57	59.90/	18.48	69.11/	27.33	74.14/	31.41
PTL (best)	7.52/	2.13	9.33/	2.94	53.82/	15.59	60.65/	18.48	70.23/	27.33	74.14/	31.41
	<i>Oku</i> → <i>Vis</i>				<i>Vis</i> → <i>Oku</i>				<i>Vis</i> → <i>ICG</i>			
baseline	1.62/	0.47	2.04/	0.57	17.13/	4.53	36.82/	9.87	2.92/	0.56	7.46/	1.83
PTL (5th itr.)	2.72/	0.94	3.05/	1.07	30.72/	7.45	42.39/	11.41	26.86/	7.22	29.26/	7.27
PTL (best)	3.00/	1.22	3.56/	1.17	33.25/	8.59	42.39/	11.46	29.60/	7.69	30.01/	7.36
	<i>ICG</i> → <i>Vis</i>				<i>ICG</i> → <i>Oku</i>				<i>Oku</i> → <i>ICG</i>			
baseline	0.54/	0.13	0.99/	0.26	3.56/	0.75	10.27/	2.49	5.37/	1.25	5.23/	1.20
PTL (5th itr.)	1.09/	0.33	1.61/	0.50	11.19/	2.58	14.20/	3.56	28.98/	8.14	25.39/	6.53
PTL (best)	1.58/	1.02	1.70/	0.63	12.82/	2.96	14.20/	3.71	28.98/	8.14	26.62/	6.53

Table 4. **Cross-domain detection accuracy.** The table shows experiments with  $3 \times 3$  cross-domain setups. For each setup, datasets shown to the left and right of the arrow are the training and test sets, respectively. The accuracies of PTL and the baseline without using virtual images for training are shown. Setups on the top use training and test images from the same domain, which provides a baseline accuracy in the cross-domain setups. All setups in each column are tested on the same dataset and the same low-shot regime.

Thanks For Listening