

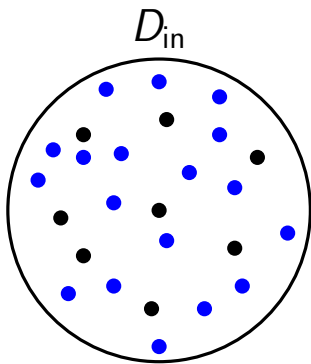


# Removing Undesirable Feature Contributions Using Out-of-Distribution Data

Saehyung Lee Changhwa Park Hyungyu Lee Jihun Yi Jonghyun Lee Sungroh Yoon

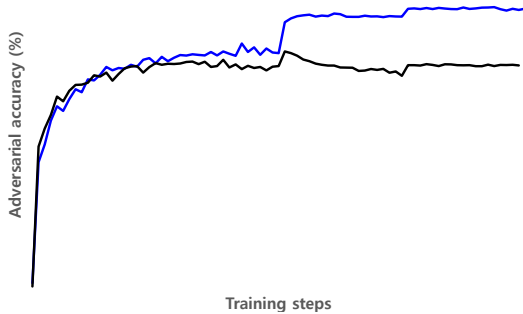


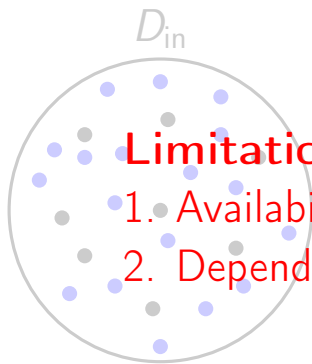
SEOUL  
NATIONAL  
UNIVERSITY



●: Training data

●: Unlabeled data



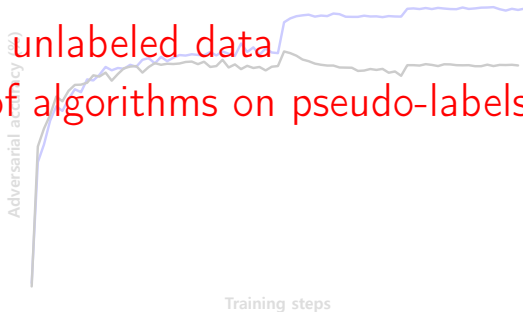


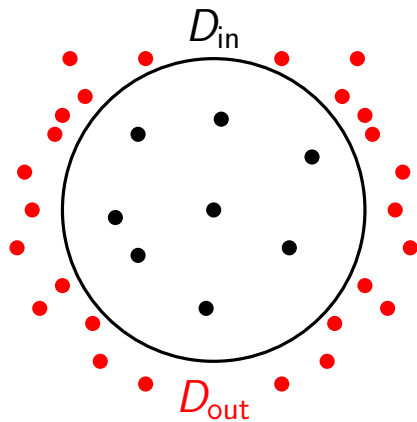
## Limitations:

1. Availability of unlabeled data
2. Dependence of algorithms on pseudo-labels

●: Training data

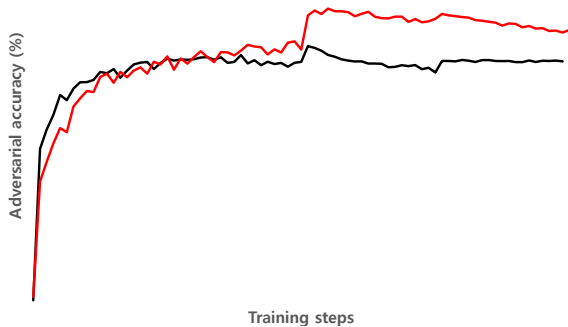
●: Unlabeled data





●: Training data

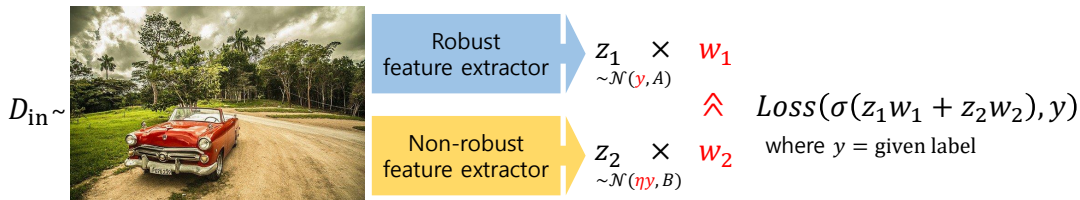
●: OOD data



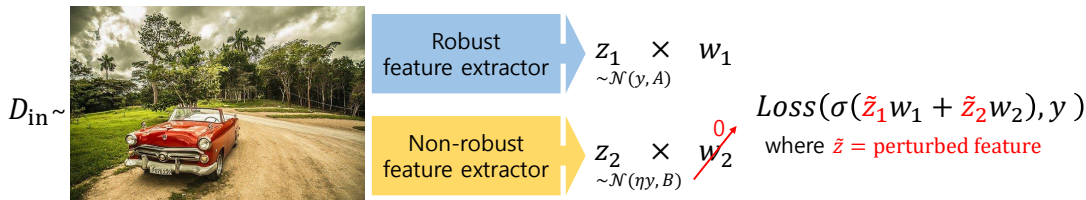
# Assumption

Undesirable features are shared among diverse image datasets.

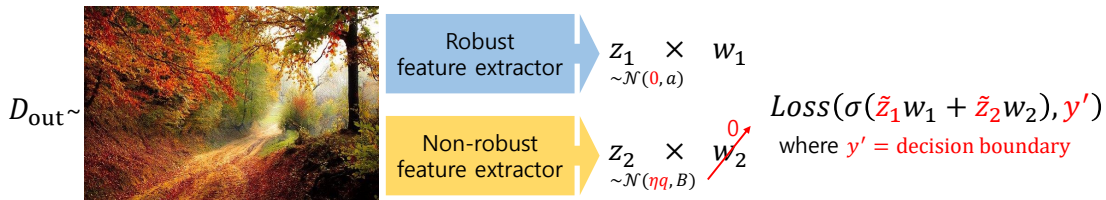
# Non-robust classifier



# Adversarial training

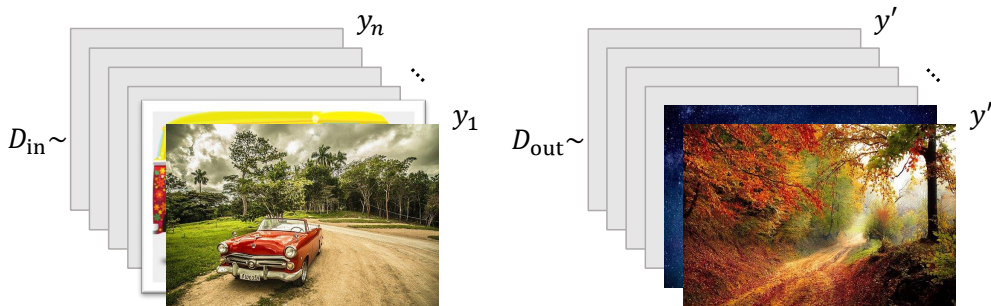


# Proposed method





# Proposed method



$$y' = \left[ \frac{1}{c}, \frac{1}{c}, \dots, \frac{1}{c} \right], \quad \text{where } c = \# \text{ of classes}$$

Training mini-batch

# Experimental details

- ▶ **80M-TI**: An OOD dataset created from the 80 Million Tiny Images dataset (Torralba et al., 2008) using confidence-based OOD detection algorithm (Carmon et al., 2019).
- ▶ Mainly compare the following setting:
  - ▶ **PGD**: The model trained using PGD-based adversarial training (Madry et al., 2017) on the target dataset.
  - ▶ **TRADES**: The model trained using TRADES (Zhang et al., 2019) on the target dataset.
  - ▶ **OAT<sub>PGD</sub>**: The model which is adversarially trained with OAT based on a PGD approach.
  - ▶ **OAT<sub>TRADES</sub>**: The model which is adversarially trained with OAT based on TRADES.

# Adversarial training results

Model	Target	OOD	Clean	PGD100	CW100	AA
Standard	CIFAR10	-	95.48	0.00	0.00	0.00
PGD		-	87.48	49.92	50.80	48.29
PGD+CutMix		-	89.35	53.39	52.35	49.05
TRADES		-	85.24	55.69	54.04	52.83
OAT <sub>PGD</sub>		80M-TI	86.63	<b>56.77</b>	<b>52.38</b>	<b>49.98</b>
OAT <sub>TRADES</sub>		80M-TI	<b>86.76</b>	<b>59.66</b>	<b>55.71</b>	<b>54.63</b>
Standard	CIFAR100	-	78.57	0.02	0.00	0.00
PGD		-	61.37	24.66	24.68	22.76
TRADES		-	58.84	30.24	27.97	26.91
OAT <sub>PGD</sub>		80M-TI	<b>61.54</b>	<b>30.02</b>	<b>27.85</b>	<b>25.36</b>
OAT <sub>TRADES</sub>		80M-TI	<b>63.07</b>	<b>34.23</b>	<b>29.02</b>	<b>27.83</b>
Standard	ImgNet10 (64 x 64)	-	86.03	0.11	0.06	0.00
PGD		-	82.80	48.77	48.86	48.34
OAT <sub>PGD</sub>		ImgNet990	81.91	<b>59.03</b>	<b>54.69</b>	<b>53.83</b>

# Adversarial training results

CIFAR10					ImgNet10 (64 x 64)		
OOD	None	SVHN	Simpson	Fashion	None	Places365	VisDA17
Clean	87.48	86.16	86.79	85.84	82.80	82.37	82.46
PGD20	50.41	<b>53.70</b>	<b>53.88</b>	<b>53.27</b>	49.00	<b>59.86</b>	<b>55.34</b>
CW20	51.11	<b>52.21</b>	<b>52.15</b>	<b>51.70</b>	48.91	<b>56.23</b>	<b>53.80</b>

# Standard training results

Dataset $N$	CIFAR10 2,500 / Full	CIFAR100 2,500 / Full
Standard	65.44 / 94.46	24.41 / 74.87
OAT <sub>SVHN</sub>	<b>68.56</b> / 94.45	<b>24.82</b> / <b>75.65</b>
OAT <sub>Simpson</sub>	<b>70.08</b> / 94.43	<b>27.04</b> / <b>76.03</b>
OAT <sub>80M-TI</sub>	<b>72.49</b> / <b>95.20</b>	<b>26.13</b> / <b>76.30</b>
Pseudo-label	- / 95.28	- / 77.24
Fusion	- / <b>95.53</b>	- / <b>77.36</b>

Dataset $N$	ImgNet10 (64 x 64) 100 / Full	ImgNet10 (160 x 160) 100 / Full
Standard	37.90 / 86.93	33.36 / 90.91
OAT <sub>VisDA17</sub>	36.21 / 86.71	<b>35.93</b> / <b>91.23</b>
OAT <sub>Places365</sub>	<b>41.84</b> / <b>88.37</b>	<b>40.11</b> / <b>91.42</b>
OAT <sub>ImgNet990</sub>	<b>42.18</b> / <b>87.88</b>	<b>40.41</b> / <b>91.87</b>

# Conclusion

- ▶ Propose **out-of-distribution data augmented training (OAT)** to leverage OOD data for adversarial and standard learning.
- ▶ Our **theoretical analyses** demonstrate how our proposed method can improve robust and standard generalization.
- ▶ The experimental results on CIFAR-10, CIFAR-100, and a subset of ImageNet suggest that OAT can **help reduce the generalization gap in adversarial and standard learning**.
- ▶ By applying OAT using various OOD datasets, it is shown that **undesirable features are shared among diverse image datasets**.

# Thank you!

<https://github.com/Saehyung-Lee/OAT>

**Acknowledgements:** This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (Ministry of Science and ICT) [2018R1A2B3001628], the BK21 FOUR program of the Education and Research Program for Future ICT Pioneers, Seoul National University in 2020, and AIR Lab (AI Research Lab) in Hyundai & Kia Motor Company through HKMC-SNU AI Consortium Fund.