

ViT에서 일반화의 shift 등 특성이 고려 연구.

### \* Introduction.

Transformer의 성능은 vision에서 큰 성능을 보였다.

TF는 generalization 능력의 이유는 조사되지 않았다.

보통 train, test는 같은 분포에 있다고 가정하지만, 그렇지 않아 많다.

따라서 Out-of-Distribution의 성능은 중요하다.

∴ TF의 Distribution Shift (DS)에 대해 조사.

DS는 전경과 배경이 따른 정도가 포함된다.

전경은 pixel, texture, shape, part가 있다.

DS는 semantically 다르므로 background, corruption, texture, style shift로 묶인다.

좋은 OOD를 위해선 인식이 비슷한 inductive bias를 가져야 함

↳ ∴ TF의 inductive bias가 있는지 조사.

↳ 1) ViT는 background, texture의 악한 bias를 가지고.  
shape, style의 큰 bias를 가짐  
↳ ↳ 인식과 비슷

∴ CNN보다 잘 일반화, OOD↑

↳ Model이 규칙이 따라, bias가 놓다.

↳ corruption, background shift가 증가하여, IID 및 OOD의 성능이 좋음.

3) 큰 patch로 훈련은 texture shift 미만 좋고, 나머지 성능↓

이후 논문에서 generalization의 강화된 GE-ViT 제시.

### \* Related Work.

## \* Distribution Shifts and Evaluation Protocols

### - Taxonomy of Distribution Shifts

background, corruption, texture, style shift로 나누어 분석 수행.

#### - Background Shift

auxillary cue로 구분됨

하지만 배경이 지배적일수 있으며 이를 빙자하기 위해서, background shift invariant 필요함.

#### - Corruption Shift

Image에 포함된 자연적인 불순물

↳ 이미지 훨씬, 처리 단계에서 발생

Pixel 수준에서 발생할수록 성능이 향상화 되게 됨.

#### - Texture Shift

색과 intensive spatial information을 제공하고 classification 향상.

#### - Style Shift.

Texture, shape 등을 포함하는 복잡한 Shift.

### - Model Zoo

#### - Vision Transformer.

#### - Big Transfer

### - Evaluation Protocols

classification encoder feature F & classifier C로 나뉨

$D_{train}, D_{iid}, D_{ood}$  3 dataset 이 구성됨

### - Accuracy on OOD Data

$$Acc(F, C; \mathcal{D}_{ood}) = \frac{1}{|\mathcal{D}_{ood}|} \sum_{(x,y) \in \mathcal{D}_{ood}} \mathbb{1}(C(F(x)) = y),$$

### - IID / OOD Generalize Gap

IID dataset 위에서의 OOD에서 얼마나 잘 작동하는가.

$\therefore$  generalize Gap 사용

$$Gap(F, C; \mathcal{D}_{iid}, \mathcal{D}_{ood}) = Acc(F, C; \mathcal{D}_{iid}) - Acc(F, C; \mathcal{D}_{ood}).$$

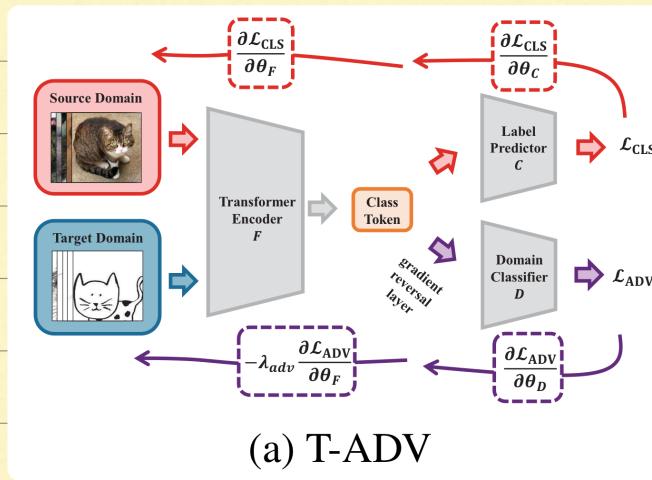
## \* Generalization - Enhanced ViT.

T-ADV, T-MME, T-SSL을 비교하여, GE-ViT 설계

### - Adversarial Learning

Domain-invariant 을 위한 domain discriminator 설계.

↳ Adv은 domain-confuse 목적



$$\mathcal{L}_{CLS} = \sum_{(x,y) \in \mathcal{D}_s} \mathcal{H}(\sigma(C(F(x))), y),$$

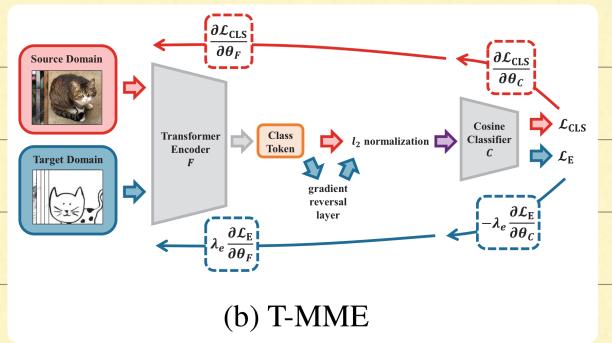
$$\mathcal{L}_{ADV} = \sum_{(x,y_d) \in \mathcal{D}_s, \mathcal{D}_t} \mathcal{H}(\sigma(D(F(x))), y_d),$$

$$\hat{\theta}_D = \arg \max_{\theta_D} \mathcal{L}_{ADV},$$

### → Gradient Reversal Layer (GRL)

### - Minimax Entropy

Target dataset conditional Entropy의 경우 minimax process를 이용하여, D를 학습하여 distribution gap을 줄임



$$\mathcal{L}_{CLS} = \sum_{(x,y) \in \mathcal{D}_s} \mathcal{H}(\sigma(C(F(x))), y),$$

$$\mathcal{L}_E = \sum_{x \in \mathcal{D}_t} \mathcal{H}(\sigma(C(F(x)))),$$

$$\hat{\theta}_C = \arg \min_{\theta_C} \mathcal{L}_{CLS} - \lambda_e \mathcal{L}_E,$$

→ cosine-similarity based  $C$ 를 사용

→ Cos weight of class,  $C \in l_2$  norm.

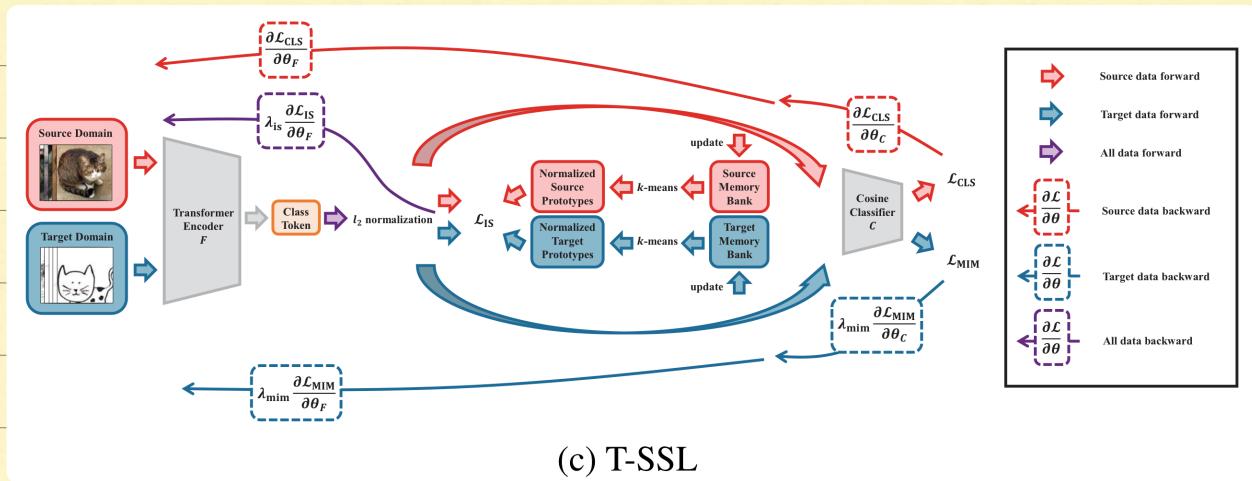
$$\hookrightarrow \frac{F(x)}{\|F(x)\|} \cdot \frac{1}{T} \frac{W^T F(x)}{\|F(x)\|}$$

unlabeled의 경우 사용한 labeled data와의 거리를 최소화하는 feature 찾기.

↳ entropy  $\rightarrow$  source domain 분리와 잘 clustering 되게 함

$\therefore$  FPLC 활용 minmax

### - Self-supervised learning



$$(\hat{\theta}_F, \hat{\theta}_C) = \arg \min_{\theta_F, \theta_C} \mathcal{L}_{CLS} + \lambda_{IS} \mathcal{L}_{IS} + \lambda_{MIM} \mathcal{L}_{MIM},$$

$$\mathcal{L}_{CLS} = \sum_{(x,y) \in \mathcal{D}_s} \mathcal{H}(\sigma(C(F(x))), y).$$

$$\mathcal{L}_{IS} = \sum_{i=1}^{|\mathcal{D}_s|} \mathcal{H}(P_i^s, c_s(i)) + \sum_{i=1}^{|\mathcal{D}_t|} \mathcal{H}(P_i^t, c_t(i)),$$

$$\mathcal{L}_{MIM} = \mathbb{E}_x [\mathcal{H}(p(y|x; \theta))] - \mathcal{H}(\mathbb{E}_{x \in \mathcal{D}_s \cup \mathcal{D}_t} [p(y|x; \theta)]).$$

## \* Systematic Study on ViTs Generalization.

### - In-Distribution Generalization.

1) Data scale이 증가할수록 모델의 성능 향상

2) DeiT-S/16이 BiT를 이길 수 있다.

### - Background Shifts Generalization Analysis

Background Shifts on the ImageNet-9 dataset

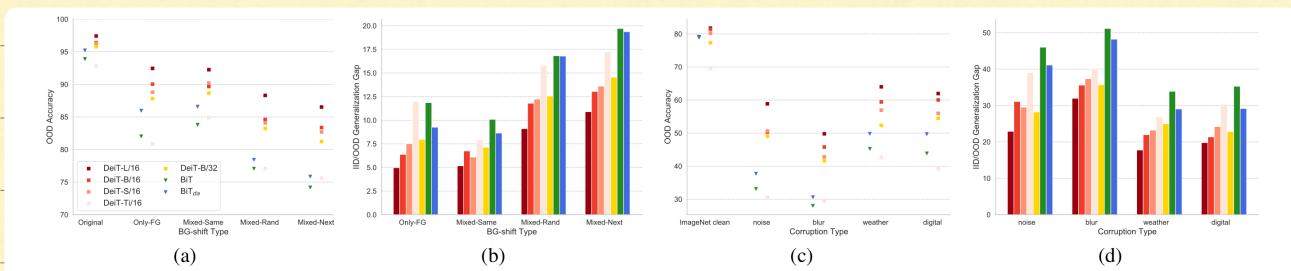


Figure 2. Results on ImageNet-9 and ImageNet-C. (a)-(b) and (c)-(d) respectively illustrate the OOD Accuracy and IID/OOD Generalization Gap for different models on ImageNet-9 and ImageNet-C datasets. From (a) and (b), we conclude that **1**) ViTs perform with a weaker background-bias than CNNs, **2**) a larger ViT extracts a more background-irrelevant representation. From (c) and (d), we draw the conclusions that **1**) ViTs deal with corruption shifts better than CNNs and generalize better along with model size scaling up, **2**) ViTs do benefit from diverse augmentation in enhancing generalization towards vicinal impurities, but their architectural advantage cannot be overlooked as well, **3**) patch size for training has little influence on ViTs' generalization ability.

### - ViTs perform with a weaker background-bias than CNNs

class와 관련이 있는 배경 mixed-same 및 neutral background의 mixed-rand의 성능 차이로 배경 의존도 개선.

Fig 1. a)

↳ ViT는 CNN보다 낮은 배경 편향.

### - A larger ViT extracts a more background-irrelevant representation.

ViT가 주제에 foreground를 더 집중

- Corruption Shifts Generalization Analysis → Fig 2. c, d

- ViT deal with Corruption Shift better than CNNs

and Generalize better along with model size scaling up

대부분 ViT가 BiT보다 성능이 좋음

↳ 더 넓은 OOD gap 성능

- ViTs benefit from diverse augmentations in enhancing generalization towards vicinal impurities.

but their architectural advantage cannot be overlooked.

BiT의 augmentation은 sensitivity가 높아서지만, ViT가 성능이 좋은데, architecture 덕분으로 볼 수 있다.

- Patch size for training has little influence on ViT's generalization ability.

Patch가 크면 일반화가 더 좋지만, model 내에서만 적용됨

- Texture Shifts Generalization Analysis → Fig 3. a, b

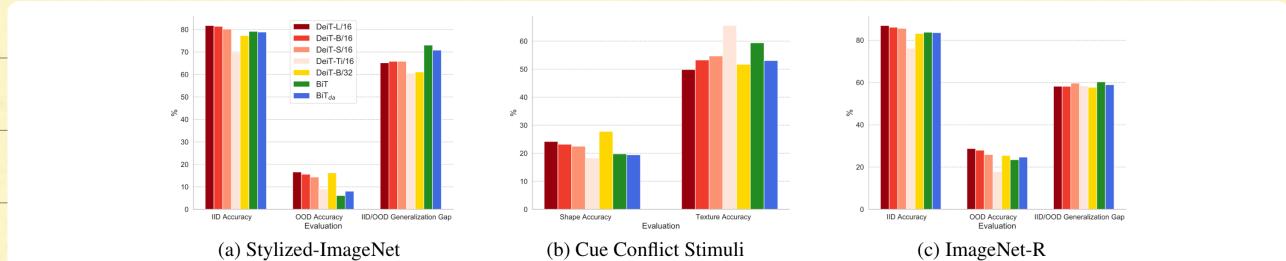


Figure 3. Results on Stylized-ImageNet, Cue Conflict Stimuli and ImageNet-R. (a), (b) and (c) respectively illustrate the OOD Accuracy and IID/OOD Generalization Gap for different models on Stylized-ImageNet, Cue Conflict Stimuli, and ImageNet-R data sets. From (a) and (b) we could draw the following conclusions that 1) ViTs' stronger bias towards shape enables them to generalize better under texture shifts and their shape biases have a positive correlation with their sizes, 3) ViTs with larger patch size exhibit a stronger bias towards the shape. From (c) we observe that most ViTs beat BiTs in OOD accuracy while having little difference in the IID/OOD generalization gap.

- ViTs' Stronger bias towards shape enables them to generalize better under texture shifts

and their shape biases have a positive correlation with their sizes

ViT의 shape이 강한 편인 때에는 texture는 일반화 잘됨

나곡복적인 영향을 덜 받음

- ViTs with larger patch size exhibit a stronger bias towards the shape.

근 patch 인수는 global에 짐작

- Style Shifts Generalization Analysis → Fig 3.C.

- ViTs have a diverse performance on IID/OOD generalization gap under Style Shift.



Figure 4. Results on DomainNet. From the results, we can conclude that 1) DeiT-S/16 performs better on the small-scale datasets in IID conditions. Thus, the model easily outperforms BiTs in OOD accuracy, 2) when inspecting the IID/OOD generalization gap, the results differ a lot. When models are trained on clipart and painting, there is no obvious difference of gap between DeiT-S/16 and BiTs.

대체로 ViTs의 OOD가 BiTs의 OOD를 높다.

IID/OOD의 gap은 비슷

∴ 예술 작품화 등의 task에서 ViTs가 더 효과는 있다

그리고, ViTs 경우 pre-trained를 더 잘 활용하고,

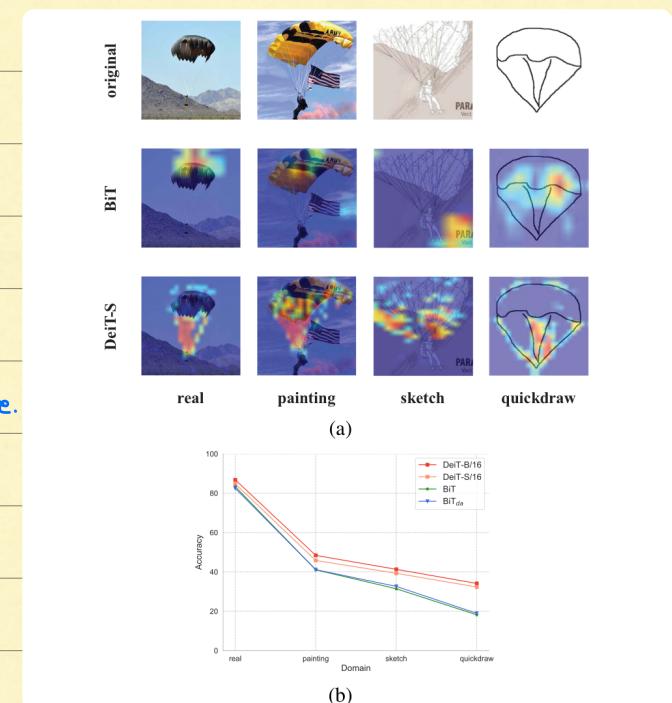
downstream이 좋다.

- ViTs shows stronger bias towards object structure.

Structure만 남은 degradate image의 경우는 좋다.

ViTs는 주로 구조 정보만 남아도 잘 catch하는데, BiT는 실패.

ViTs는 구조에 더 강하다.



- ViTs will eliminate different levels of DS in different layers.

ViTs every layer with DS를 한다

## \* Studies on Generalization-Enhanced ViTs

### - Setting

DominNet 사용

BiT의 1) 및 ViT의 2) 사용

### - Performance Analysis

1) GE-ViT는 4% 향상

2) 3D의 GE-ViT 모두 상위.

CNN은 self-supervised로 훈련.

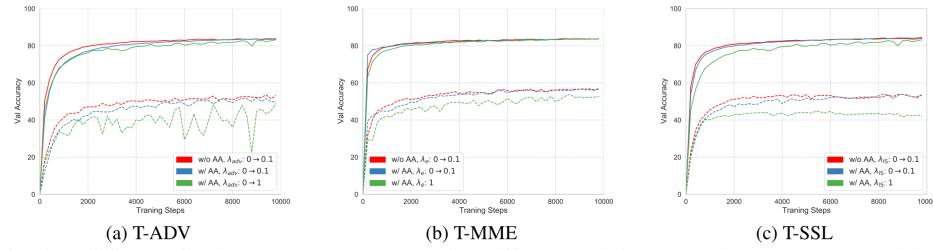
CNN과 ViT 모두 경쟁이 있는 듯

ViT가 더 어렵

**Table 2. Results of Generalization-enhanced methods.** Specifically, we compare three types of GE-ViTs with their corresponding CNNs. From the results we could conclude that 1) equipped with GE-ViTs, we achieve significant performance boosts towards out-of-distribution data by 4% from vanilla ViTs. 2) three GE-ViTs have almost the same improvement from vanilla models on OOD accuracy. 3) for the enhanced transformer models, larger ViTs still benefit more for the out-of-distribution generalization.

| Model     | Method    | R to C       | R to P       | P to C       | C to S       | S to P       | R to S       | P to R       | Avg.         |
|-----------|-----------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| DeiT-B/16 | -         | 54.64        | 48.40        | 40.37        | 45.69        | 36.75        | 41.31        | 55.33        | 46.07        |
|           | T-ADV     | 58.19        | 50.85        | 41.91        | 51.18        | 46.12        | 47.47        | 55.65        | 50.20        |
|           | T-MME     | <b>60.59</b> | <b>51.98</b> | 42.30        | 50.32        | 45.79        | <b>47.92</b> | 54.87        | 50.54        |
|           | T-SSL     | 56.80        | 49.06        | <b>45.96</b> | <b>51.79</b> | <b>46.95</b> | 45.95        | <b>60.98</b> | <b>51.07</b> |
| DeiT-S/16 | -         | 50.60        | 45.82        | 36.09        | 43.39        | 35.24        | 39.29        | 52.08        | 43.22        |
|           | T-ADV     | 53.60        | 47.84        | 37.99        | 47.10        | 41.61        | 41.94        | 52.82        | 46.13        |
|           | T-MME     | <b>56.86</b> | <b>49.15</b> | 38.97        | 46.48        | 42.95        | <b>42.07</b> | 52.49        | 47.00        |
|           | T-SSL     | 53.86        | 46.71        | <b>42.79</b> | <b>47.25</b> | <b>43.01</b> | 40.94        | <b>57.07</b> | <b>47.37</b> |
| BiT       | -         | 42.18        | 41.14        | 30.72        | 37.01        | 28.23        | 32.64        | 48.54        | 36.78        |
|           | DANN [10] | 45.20        | 42.86        | 32.96        | 40.44        | 36.63        | 35.26        | 49.25        | 40.37        |
|           | MME [24]  | 50.21        | <b>44.61</b> | 34.75        | 40.27        | 38.41        | 37.83        | 47.58        | 41.95        |
|           | SSL [34]  | <b>52.55</b> | 42.80        | <b>39.03</b> | <b>45.72</b> | <b>39.08</b> | <b>39.65</b> | <b>56.07</b> | <b>44.98</b> |
| VGG-16    | -         | 39.39        | 37.32        | 26.36        | 32.96        | 25.55        | 27.79        | 45.70        | 33.58        |
|           | DANN [10] | 43.26        | 40.09        | 28.68        | 36.22        | 31.63        | <b>35.45</b> | 44.73        | 37.15        |
|           | MME [24]  | 42.65        | <b>42.46</b> | 27.41        | <b>36.93</b> | 33.94        | 32.58        | <b>45.87</b> | 37.41        |
|           | SSL [34]  | <b>43.79</b> | 41.88        | <b>32.19</b> | 35.73        | <b>36.99</b> | 31.05        | 55.18        | <b>39.54</b> |

### - Smooth Feature Alignment.



**Figure 7. Investigation of Generalization-enhanced methods with different training strategies.** (a)-(c) show training curves on both source domain and target domain. From the results, we can conclude that classical training strategies (the green lines) on CNNs are not suitable for ViTs, which need smoother strategies (the red lines) to align features in both domains.