

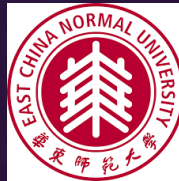
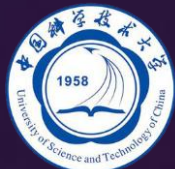


ICLR

Boosting Vanilla Lightweight Vision Transformers via Re-parameterization

Alibaba Cloud¹, Alibaba Group²,
University of Science and Technology of China³,
East China Normal University⁴,

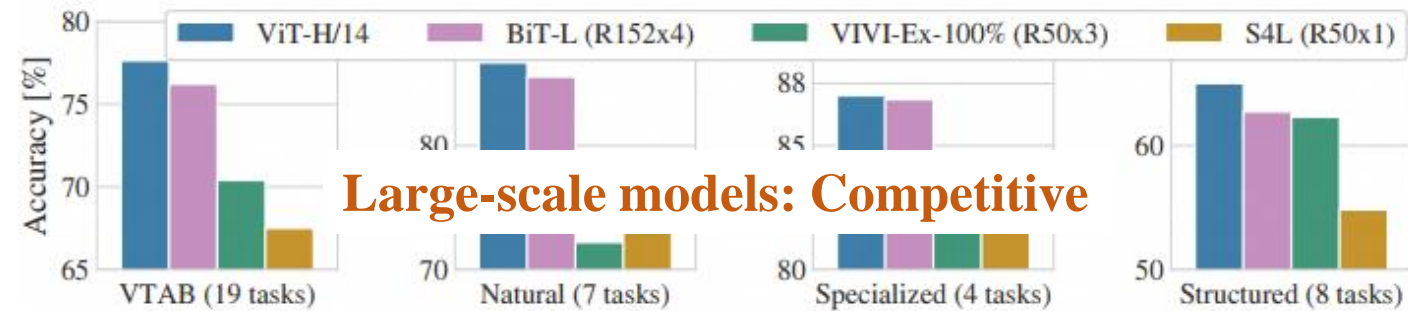
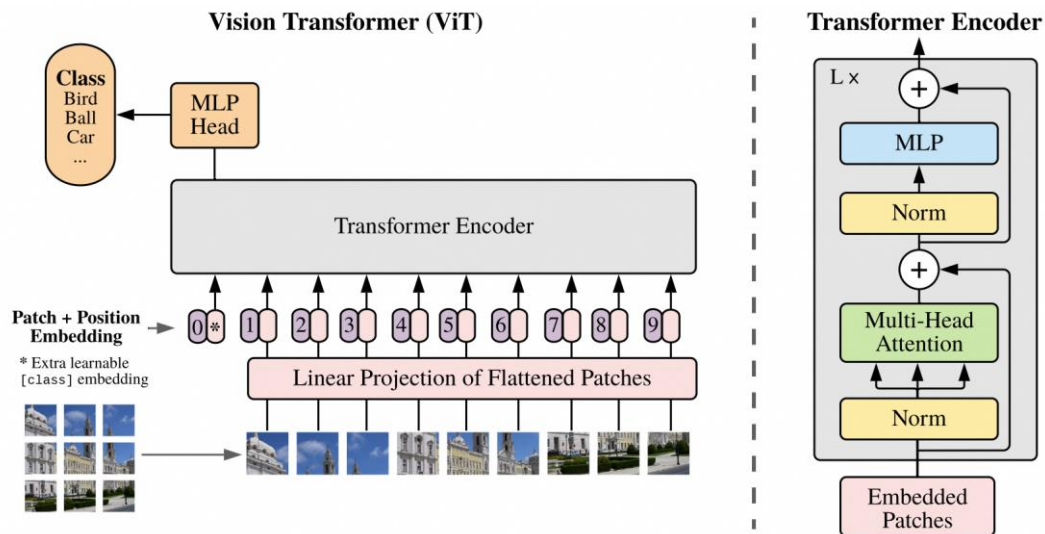
Zhentao Tan^{1,3}, Xiaodan Lin^{4,2}, Yue Wu¹, Qi Chu³, Le Lu², Nenghai Yu³, Jieping Ye¹



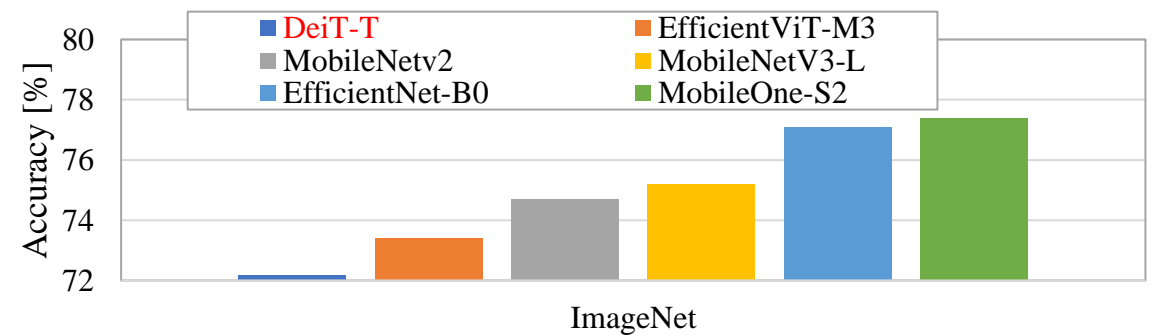
Background: ViTs



Performance difference between **large-scale** and **lightweight** vanilla vision Transformers.



Large-scale models: Competitive

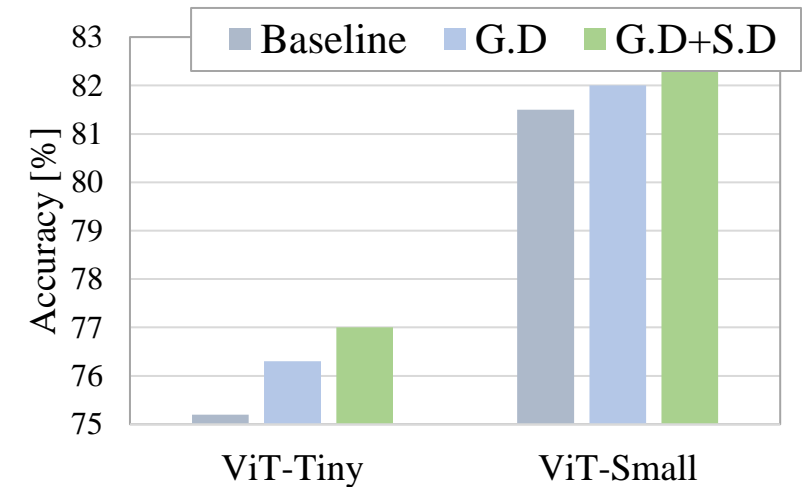
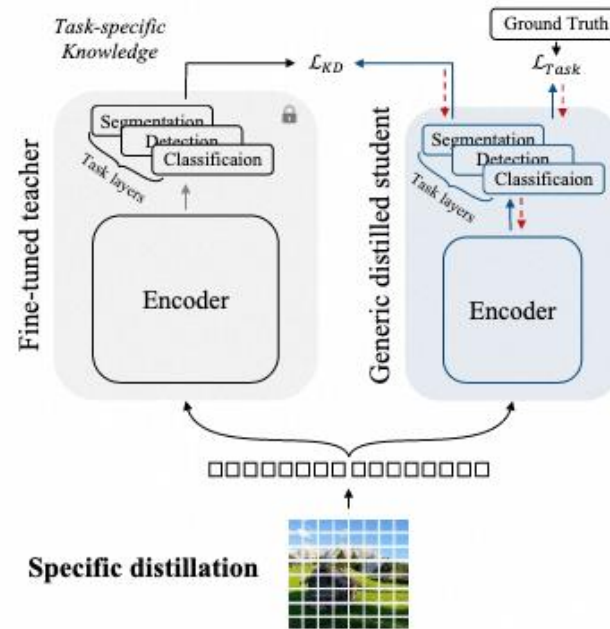
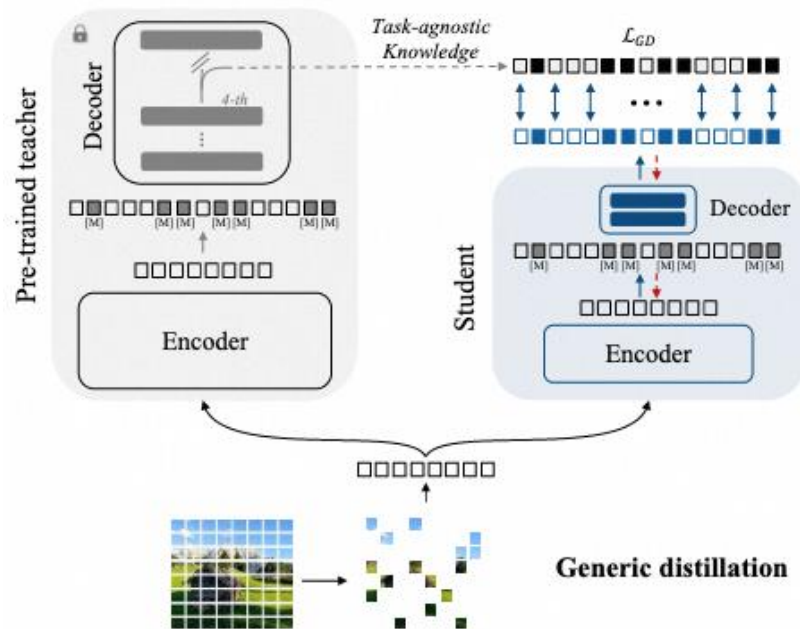


Lightweight models: Unsatisfactory

Background: Distillation

How can we improve the performance of **vanilla lightweight ViTs** **without any inference pipeline changes**?

Related works: Knowledge Distillation

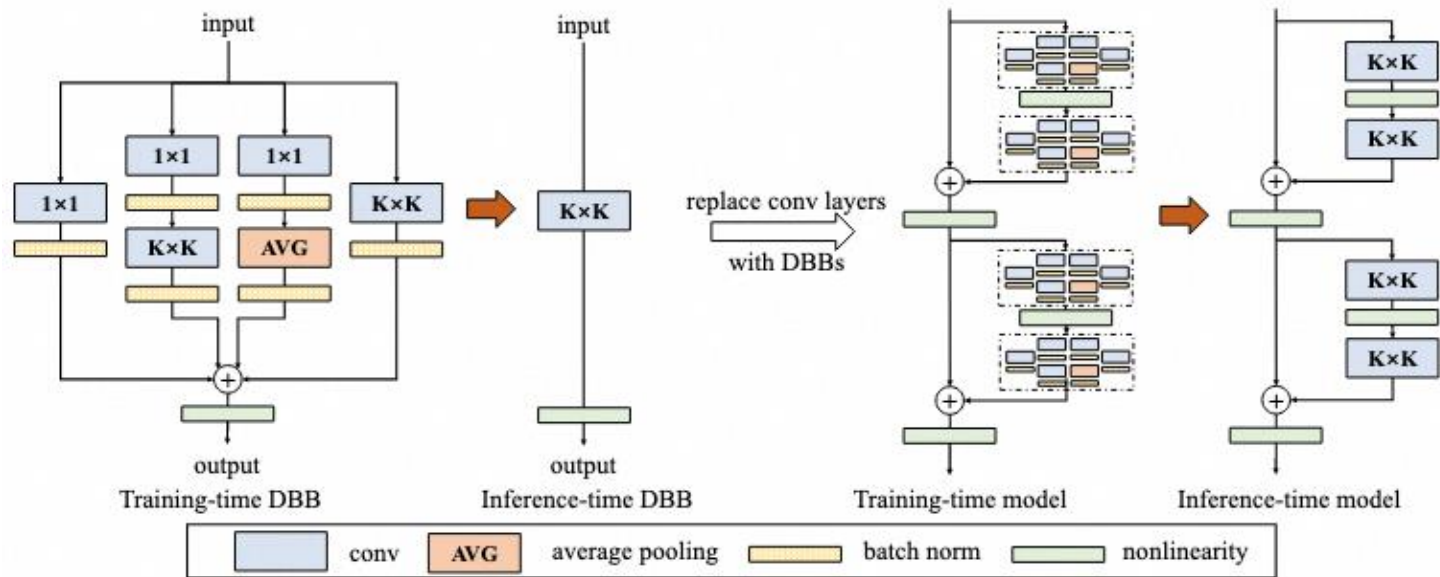


**Costs of pre-trained teachers
(especially for S.D)**

Background: Re-parameterization



Structure re-parameterization in CNNs

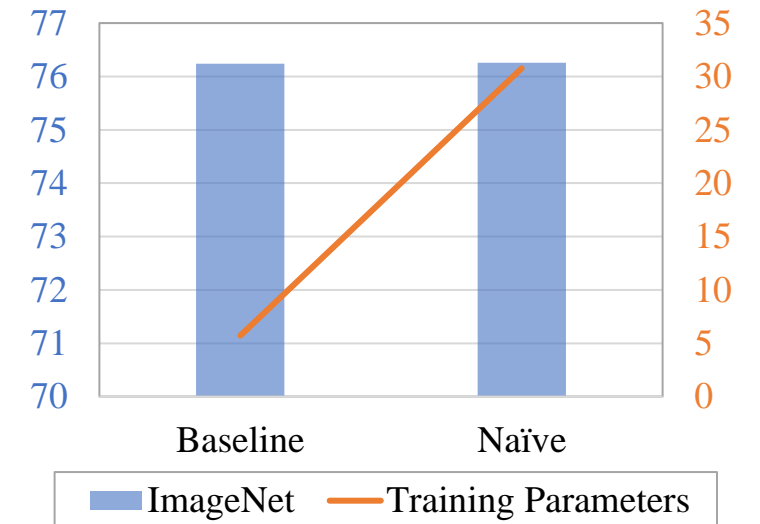
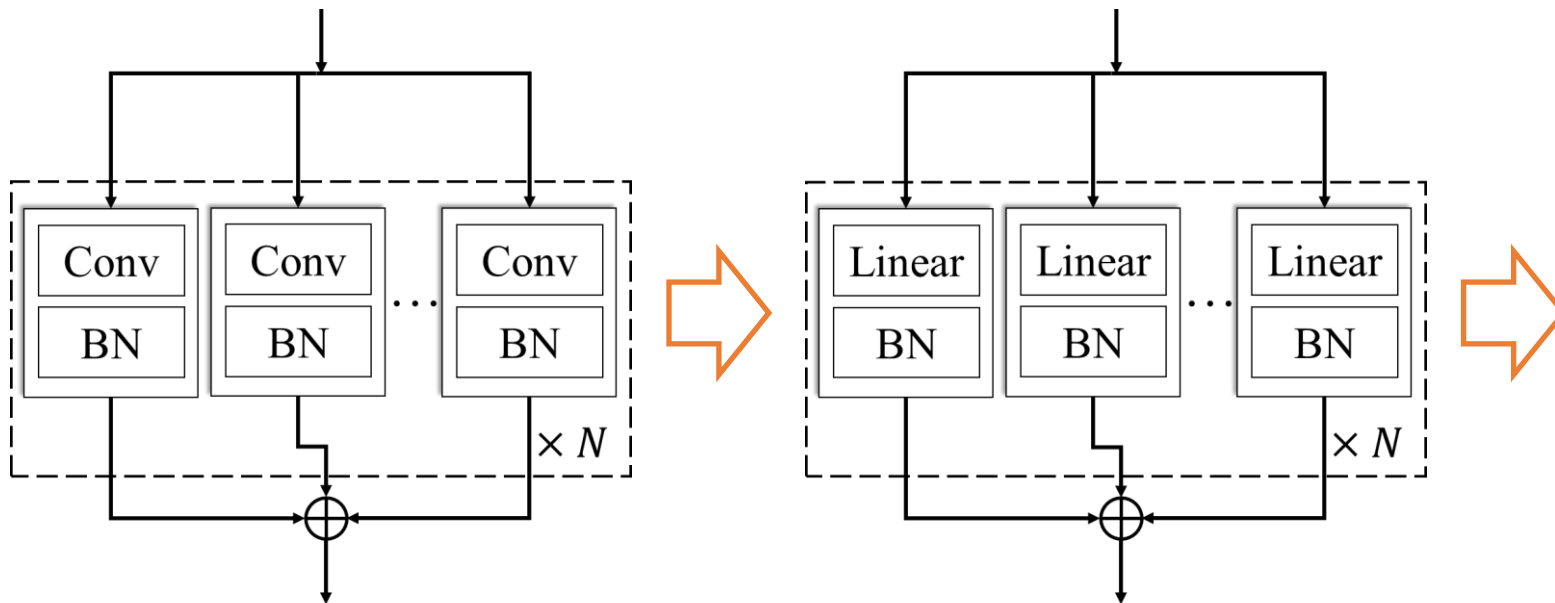


- **Training:**
 - multi-branch
 - over parameterization
- **Inference:**
 - single $K \times K$ convolution

Can we apply re-parameterization to vanilla ViTs?

Method: Single Experiments

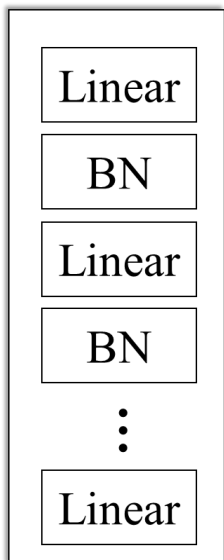
Single application: directly replacing convolution layers into linear layers.



Method: Linear Ensemble



Stacking linear layers with batch normalization in-between them.



- **Effectiveness**: similar to MLP which is appropriate to transformers and plays an important role to represent rich intra-token information.
- **Rationality**: batch normalization is still can be used in-between layers while keeping original layer normalization unchanged.
- **Operability**: it can be fused to a single linear layer after training.

1. Merging linear and batch normalization

$$\mathbf{W}'_{i,:} = \frac{\gamma_i}{\sigma_i} \mathbf{W}_{i,:}, \quad b'_i = \frac{(b_i - \mu_i)\gamma_i}{\sigma_i} + \beta_i,$$

2. Merging two adjacent linear layers

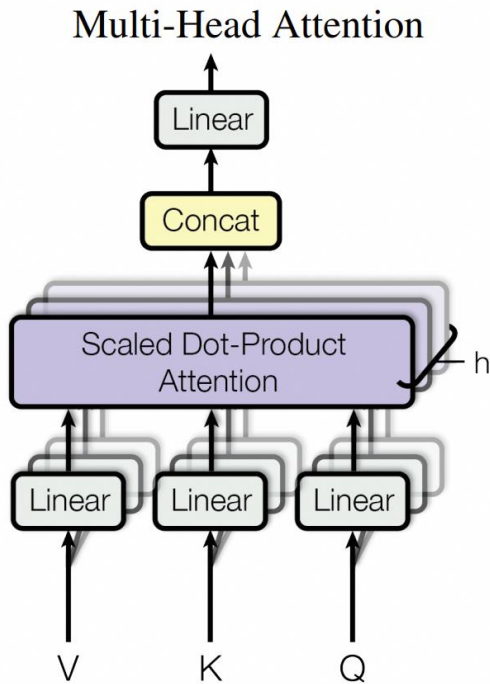
$$\begin{aligned} W'_{i,j}(l+1, l) &= \sum_{k=1}^{C_l} W_{i,k}^{l+1} W_{k,j}^l, \\ b'_i(l+1, l) &= \sum_{k=1}^{C_l} b_k^l W_{i,k}^{l+1} + b_i^{l+1}, \end{aligned}$$

Method: Distribution Rectification



Multi-branch re-parameterization will change the feature distribution.

Attention mechanism in vision transformers is **sensitive to this distribution changes.**



Normal Attention Operation:

$$Attention(Q, K, V) = softmax(\frac{A}{\sqrt{C_k}})V, \quad A = QK^T.$$

$$A_{i,j} = \sum_{k=1}^{k=C_k} Q_{i,k} K_{k,j} \quad \text{Variance scales to } C_k$$

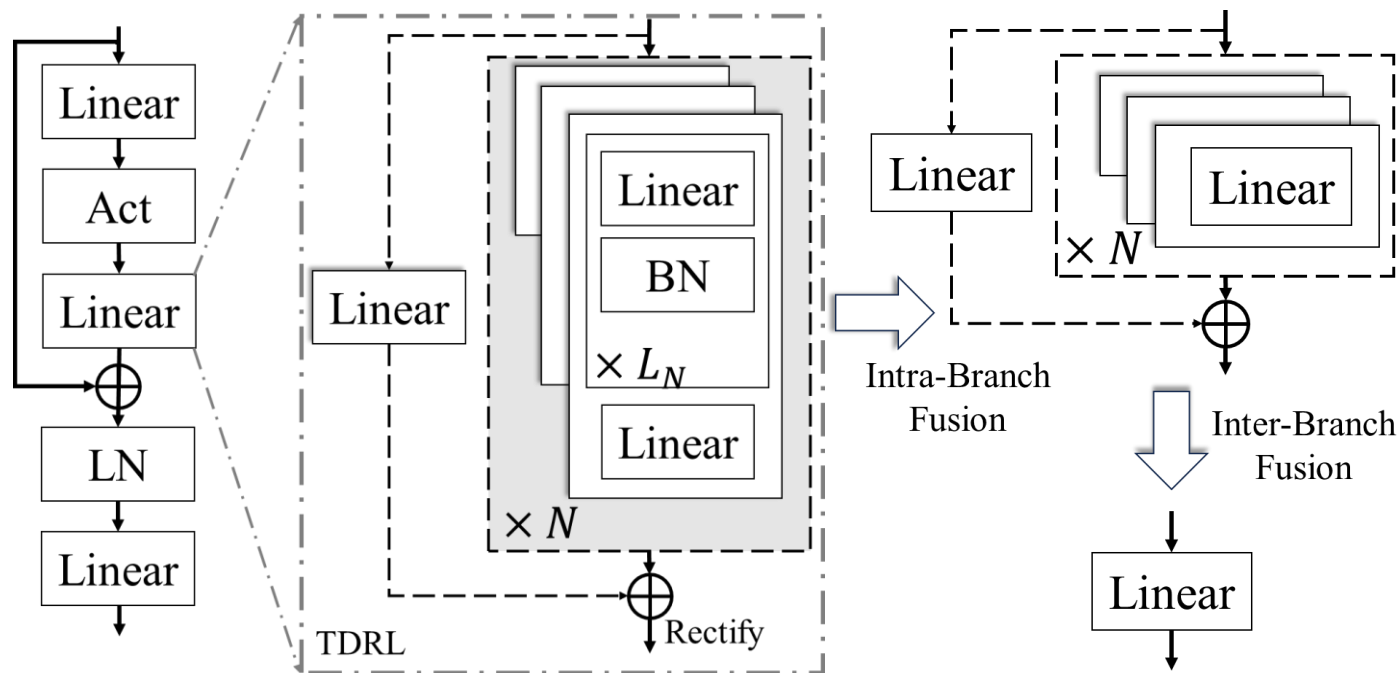
Re-parameterization Attention Operation:

$$A'_{i,j} = \sum_{k=1}^{k=C_k} (\sum_{n=1}^{N_Q} Q_{i,k}^n) (\sum_{m=1}^{N_K} K_{k,j}^m) = \sum_{k=1}^{k=C_k} \sum_{n=1}^{N_Q} \sum_{m=1}^{N_K} Q_{i,k}^n K_{k,j}^m.$$

$$R(x) = \begin{cases} BN(x), QKV \\ \frac{x}{\sqrt{C_k N_Q N_K}}, others \end{cases} \quad \text{Variance scales to } C_k N_Q N_K$$

Method: Structure

TDRL: Pyramid-wise Multi-branch



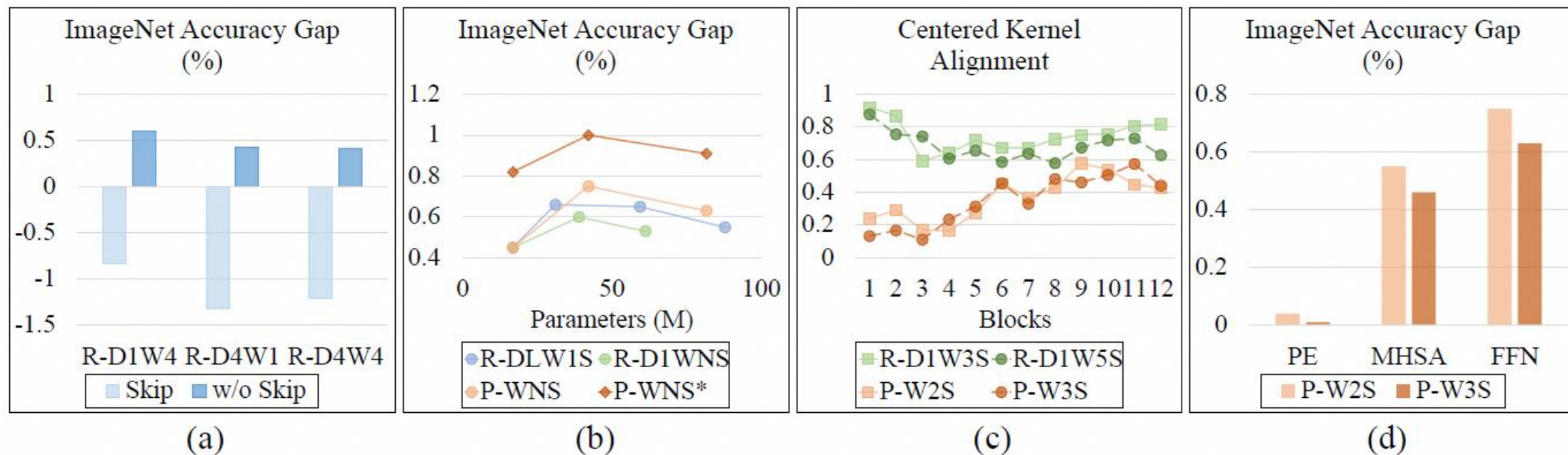
$$Y = \text{Rectify}(\text{Linear}(X) + \sum_{n=1}^N f_{n,L_n}(X)), \quad L_n = n,$$

- **Design:**
 - **Skip-branch**
 - **Rep-branch**
 - **Length: representation ability**
 - **Width: representation diversity**
- **Application:**
 - **Arbitrary linear layers**

Experiments: Ablations



➤ Structure Designs

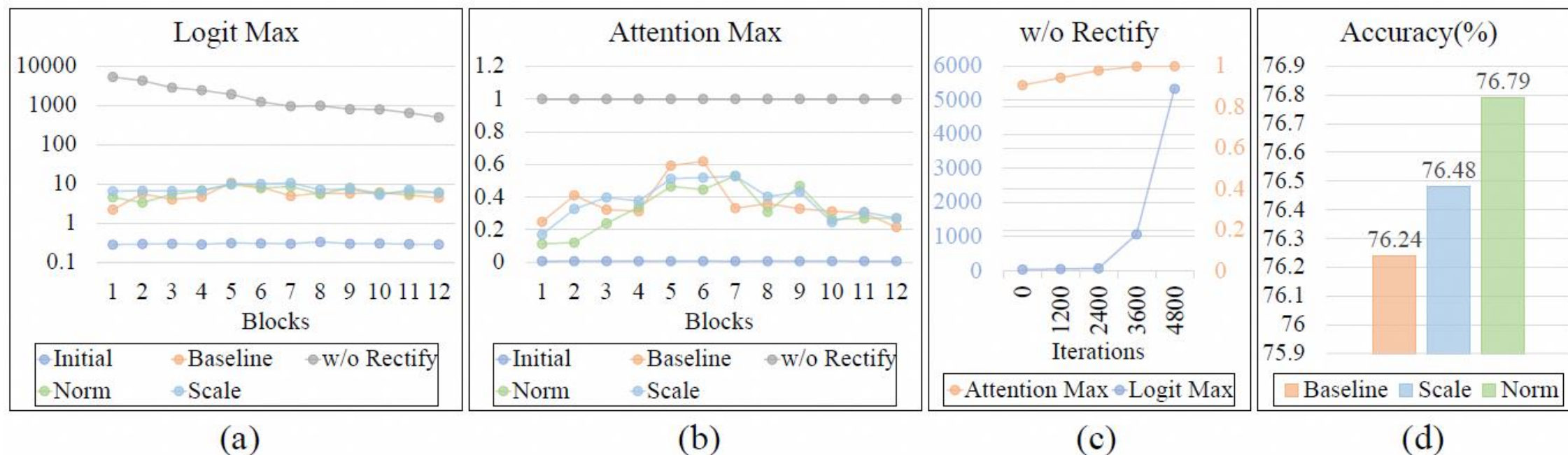


- **Skip connection plays an important role.**
- **More branches do not always lead to better performance.**
- **Pyramid-wise design leads to diversity representation between different branches.**
- **Applying TDRL to MHSA and FFN results in much more improvements.**

Experiments: Ablations



➤ Attention Distribution Rectification



- Without rectification, the maximum value in attention is prone to extreme values.
- Normalization is better than scaling as a rectification function.

Experiments: Classification



| Method | Network | Teacher | <i>FT</i> | <i>P</i> (M) | Acc(%) |
|---------------------------------------|---------|-----------|-----------|--------------|-------------------|
| Without Pre-training | | | | | |
| MobileNet-v3 (Howard et al., 2019) | CNNs | N/A | 600 | 6 | 75.2 |
| ConvNeXt-V1-F (Liu et al., 2022b) | CNNs | N/A | 600 | 5 | 77.5 |
| VanillaNet-5 (Chen et al., 2023) | CNNs | N/A | 300 | 15.5 | 72.5 |
| MobileViT-S (Mehta & Rastegari, 2021) | Hybrid | N/A | 300 | 6 | 78.3 |
| EfficientViT-M3 (Liu et al., 2023) | Hybrid | N/A | 300 | 7 | 73.4 |
| DeiT-Ti (Touvron et al., 2021) | ViTs | N/A | 300 | 5 | 72.2 |
| Manifold-Ti (Jia et al., 2021) | ViTs | CaiT-S24 | - | 6 | 75.1† |
| MKD-Ti (Liu et al., 2022a) | ViTs | CaiT-S24 | 300 | 6 | 76.4† |
| DeiT-Ti (Touvron et al., 2021) | ViTs | RegNetY | 300 | 6 | 74.5† |
| SSTA-Ti (Wu et al., 2022a) | ViTs | DeiT-S | 300 | 6 | 75.2† |
| ImageNet Pre-training | | | | | |
| DMAE-Ti (Bai et al., 2023) | ViTs | ViT-B | 100 | 6 | 74.9 |
| MAE-Lite (Wang et al., 2023) | ViTs | N/A | 100 | 6 | 76.2 |
| MAE-Ti (He et al., 2022) | ViTs | N/A | 200 | 6 | 75.2 |
| TinyMIM-Ti (Ren et al., 2023) | ViTs | TinyMIM-S | 200 | 6 | 75.8 |
| G2SD-Ti w/o S.D (Huang et al., 2023) | ViTs | ViT-B | 200 | 6 | 76.3 |
| G2SD-Ti (Huang et al., 2023) | ViTs | ViT-B | 200 | 6 | 77.0† |
| TDRL (ours) | ViTs | ViT-B | 200 | 6 | 78.3/78.6† |
| MAE-Lite (Wang et al., 2023) | ViTs | N/A | 300 | 6 | 78.0 |
| D-MAE-Lite (Wang et al., 2023) | ViTs | ViT-B | 300 | 6 | 78.4 |
| TDRL (ours) | ViTs | ViT-B | 300 | 6 | 78.7/79.1† |

- **TDRL achieves the best image classification accuracy under various epoch settings.**
- **When performing distillation during fine-tuning, the performance of TDRL can be further improved to 79.1%.**

†means performing distillation during fine-tuning.

Experiments: Dense Prediction



| Method | #Params (M) | Segmentation | Detection | |
|--------------------------------|-------------|--------------------------------------|--------------------------------------|--------------------------------------|
| | Seg/Det | mIoU | AP^{bbox} | AP^{mask} |
| Swin-T (Liu et al., 2021) | 59.9/47.8 | 44.5 | 46.0 \ddagger | 41.6 \ddagger |
| ConvNeXt-T (Liu et al., 2022b) | 60.0/48.1 | 46.0 | 46.2 \ddagger | 41.7 \ddagger |
| DINO-S (Caron et al., 2021) | 42.0/44.5 | 44.0 | 49.1 | 43.3 |
| iBOT-S (Zhou et al., 2021) | 42.0/44.5 | 45.4 | 49.7 | 44.0 |
| MAE-S (He et al., 2022) | 42.0/44.5 | 41.1/44.9 \dagger | 45.3 | 40.8 |
| MAE-Ti (He et al., 2022) | 11.0/27.7 | 36.9/42.0 \dagger | 37.9/43.5 \dagger | 34.9/39.0 \dagger |
| MAE-Lite (Wang et al., 2023) | 11.0/27.7 | 37.6 | 39.9* | 35.4* |
| D-MAE-Lite (Wang et al., 2023) | 11.0/27.7 | 42.0 | 42.3* | 37.4* |
| G2SD-Ti (Huang et al., 2023) | 11.0/27.7 | 41.4/44.5 \dagger | 44.0/46.3 \dagger | 39.6/41.3 \dagger |
| TDRL (ours) | 11.0/27.7 | 42.5/45.2\dagger | 46.5/47.4\dagger | 41.5/42.1\dagger |

- TDRL can also benefit dense prediction tasks such as semantic segmentation and object detection.

Experiments: Generality



ICLR

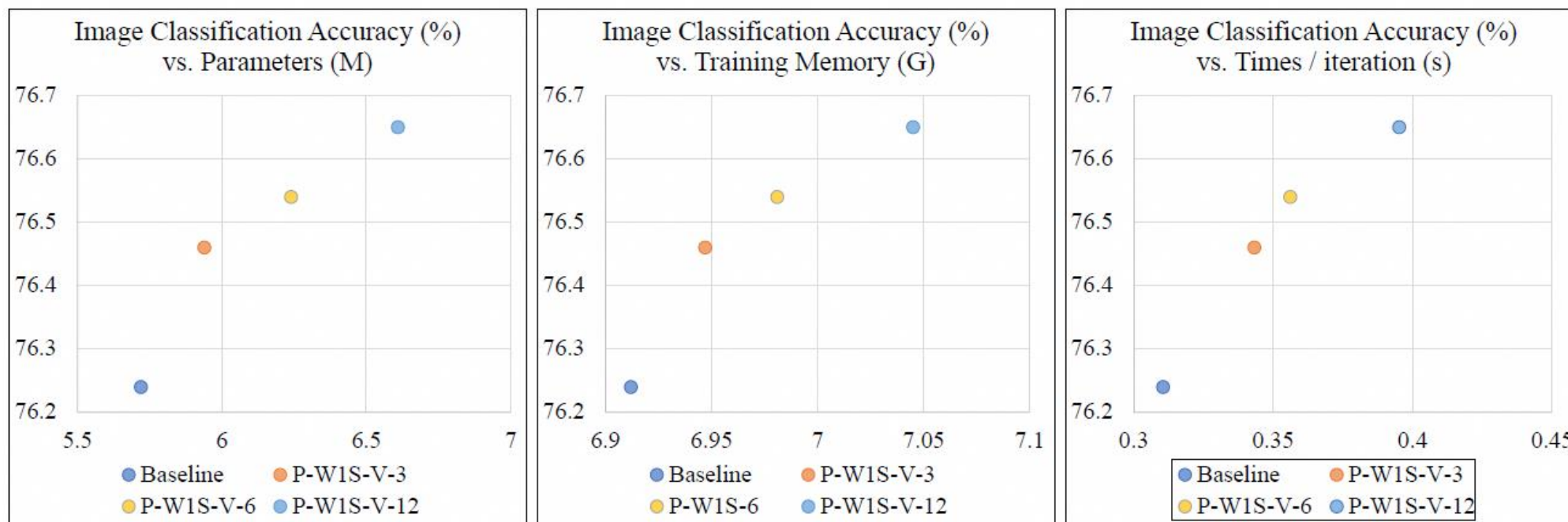
| TDRL | ViT-Small | Classification Accuracy (%) \uparrow | | | Image Generation FID \downarrow DDPM |
|--------------|-------------|--|--------------|--------------|---|
| | | Swin-Ti | Mobileone-S0 | VanillaNet-5 | |
| \times | 80.8 | 76.2 | 71.3 | 71.1 | 10.4 |
| \checkmark | 81.3 (+0.5) | 78.2 (+2.0) | 75.1 (+3.8) | 71.5 (+0.4) | 9.2 (+1.2) |

- **TDRL can be used in various networks and tasks:**
 - **Larger model (e.g., ViT-Small, Swin-Tiny)**
 - **CNN model (e.g., VanillaNet)**
 - **Hybrid model (e.g., MobileOne)**
 - **DDPM**

Experiments: Efficiency

| Method | Pre-train time (hours) | Pre-train epoch | Fine-tune epoch | Accuracy (%) |
|-----------------|------------------------|-----------------|-----------------|---------------|
| Baseline (G2SD) | 32.33 | 300 | 100 | 76.24 |
| TDRL (ours) | 32.02 | 220 | 100 | 76.73 (+0.59) |

- Under similar training costs, TDRL can still improve the performance of ViT-Tiny.





- **We propose a novel re-parameterization method, namely TDRL, to improve the performance of vanilla lightweight Vision Transformers.**
- **To improve the representation ability, we design a linear ensemble way and a pyramid-wise multi-branch structure.**
- **For stable training, we analyze the feature distribution change issues and propose a simple rectification method.**
- **Experiments on various tasks and backbone have demonstrated the effectiveness of our TDRL.**



ICLR

Thanks



<https://openreview.net/pdf?id=3rmpixOjPS>