# AdaShift: Learning Discriminative Self-Gated Neural Feature Activation With an Adaptive Shift Factor

**Sudong Cai**

Graduate School of Informatics, Kyoto University
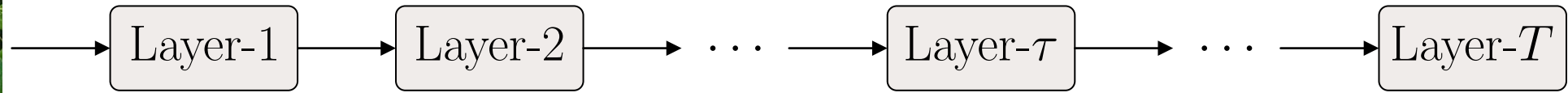
京都大学

KYOTO UNIVERSITY

1

# Introduction

- Preliminaries



$$\tilde{x} = \langle \boldsymbol{w}, \boldsymbol{x} \rangle$$

**Comprehensive scoring** (relaxed similarity)

$$\phi(\tilde{x}) = \rho(\tilde{x}) \cdot \tilde{x}$$

**Neural Activation: Feature selection by weights**
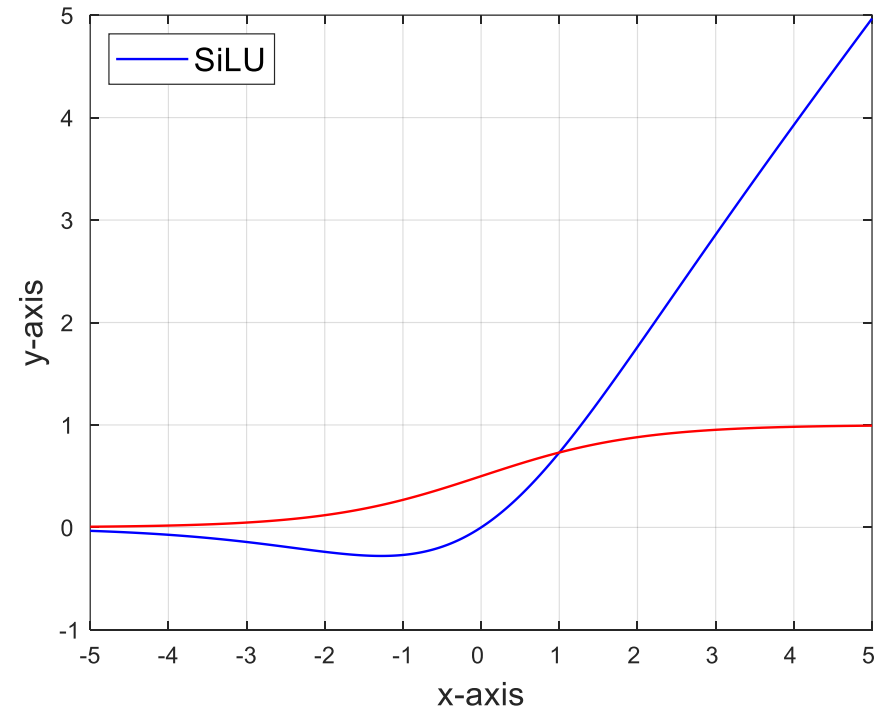
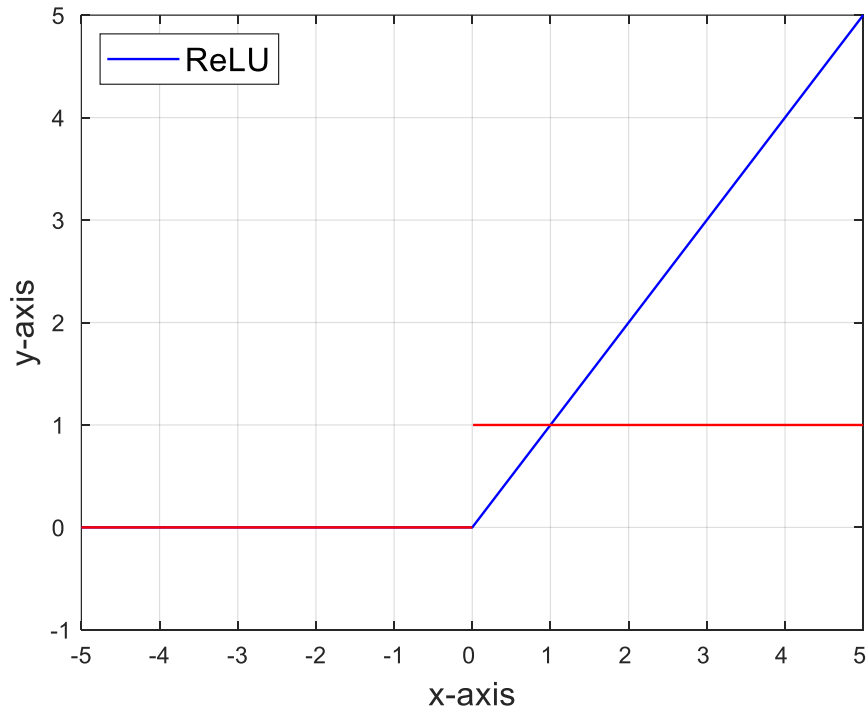**Importance score (weight)**

京都大学

# Introduction

- Self-Gated Activation Functions

$$\rho\left(\tilde{x}\right) = \begin{cases} 1, & \tilde{x} > 0\,; \\ 0, & \tilde{x} \leq 0\,. \end{cases}$$

Smoothing →

$$\rho\left(\tilde{x}\right) = \mathrm{sigmoid}\left(\tilde{x}\right)\,.$$

# Intuition

- **A critical problem:** *Mismatched Feature Scoring (MFS)*[1]



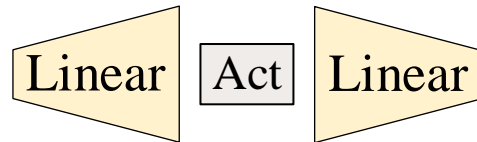$$\tilde{x} = \cos \theta_{\boldsymbol{w},\boldsymbol{x}} \|\boldsymbol{w}\| \|\boldsymbol{x}\|$$

1. Sudong Cai. IIEU: Rethinking Neural Feature Activation from Decision-Making. In *Proc. ICCV*, pages 5796-5806, 2023.

KYOTO UNIVERSITY

# Intuition

- Why AdaShift – Remaining problems in the modeling of IIEU

1. Efficiency of training

$$\phi(\tilde{x}) = \varsigma\left(\frac{\tilde{x}}{\|\boldsymbol{x}\|\|\boldsymbol{w}\| + \epsilon} + \nu\right) \cdot \tilde{x} \qquad \nabla_{\boldsymbol{w}} s(\boldsymbol{w}) = \frac{\|\boldsymbol{w}\|^2 \boldsymbol{x} - \boldsymbol{w}\boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}}{\|\boldsymbol{x}\|\|\boldsymbol{w}\|^3}$$

Linear | Act | Linear

Expensive FFN with high expansion ratio

2. Extensibility

$$\phi(\tilde{x}) = \varsigma\left(\frac{\tilde{x}}{\|\boldsymbol{x}\|\|\boldsymbol{w}\| + \epsilon} + \nu\right) \cdot \tilde{x}$$

$$\tilde{x} := \tilde{x} + \tilde{y}$$

How about if we also consider activating an un-projected feature (unit)?

**Can feature and filter norms play a useful role in feature re-weighting?**

京都大学

# Intuition of AdaShift

- Rethinking the meaning of feature and filter norms from a Softmax-based classification

  □ **(Intuition 5.1)** Feature and filter norms present local and dataset-level non-local cues, respectively

---

Pre-conditions: $\boldsymbol{w}\left(i\right):$ The representative filter of the class-$i$ $\qquad N$ classes in total

$\xrightarrow{\text{Let}}$ $\qquad \boldsymbol{w}\left(j\right):$ The representative filter of an arbitrarily class-$j$ $\quad \boldsymbol{x}$ is classified as the class-$i$

$\qquad\boldsymbol{x}:$ The representative feature of the input (image or pixel) $\qquad \forall \boldsymbol{w}\left(i\right) \neq 0, \forall \boldsymbol{x} \neq 0$

$\text{For } i \neq j$
$\xrightarrow{\phantom{For}}$
$$\frac{e^{\tilde{x}_i + b_i}}{\sum_{c=1}^{C} e^{\tilde{x}_c + b_c}} > \frac{e^{\tilde{x}_j + b_j}}{\sum_{c=1}^{C} e^{\tilde{x}_c + b_c}} \iff e^{\tilde{x}_i + b_i} > e^{\tilde{x}_j + b_j}$$

$$\longrightarrow \quad \|\boldsymbol{w}\left(i\right)\| \, \|\boldsymbol{x}\| \cos \theta_{\boldsymbol{w}(i),\boldsymbol{x}} + b_i \quad > \quad \|\boldsymbol{w}\left(j\right)\| \, \|\boldsymbol{x}\| \cos \theta_{\boldsymbol{w}(j),\boldsymbol{x}} + b_j$$

# Intuition of AdaShift

$$\|\boldsymbol{w}(i)\| \|\boldsymbol{x}\| \cos\theta_{\boldsymbol{w}(i),\boldsymbol{x}} + b_i \quad > \quad \|\boldsymbol{w}(j)\| \|\boldsymbol{x}\| \cos\theta_{\boldsymbol{w}(j),\boldsymbol{x}} + b_j$$

Let
$$\longrightarrow \quad \alpha = b_j - b_i$$

$$\longrightarrow \quad \|\boldsymbol{w}(i)\| \cos\theta_{\boldsymbol{w}(i),\boldsymbol{x}} - \|\boldsymbol{w}(j)\| \cos\theta_{\boldsymbol{w}(j),\boldsymbol{x}} > \frac{\alpha}{\|\boldsymbol{x}\|}$$

Case 1: $\|\boldsymbol{x}\| \gg |\alpha|$
$$\longrightarrow \qquad \text{Filter norms are influential \& Feature norms are not decisive}$$

Case 2: $\|\boldsymbol{x}\| \ll |\alpha|$
$$\longrightarrow \qquad \text{Filter norms are influential \& Feature norms are influential}$$

Case 3: Others
$$\longrightarrow \qquad \text{Filter norms are influential \& Feature norms matters}$$

**A**da**Shift Prototype**

Complementary information : tensor-level non-local cues

$$\longrightarrow \qquad \phi(\tilde{x}) = \varsigma(\tilde{x} + \Delta)\,\tilde{x}$$
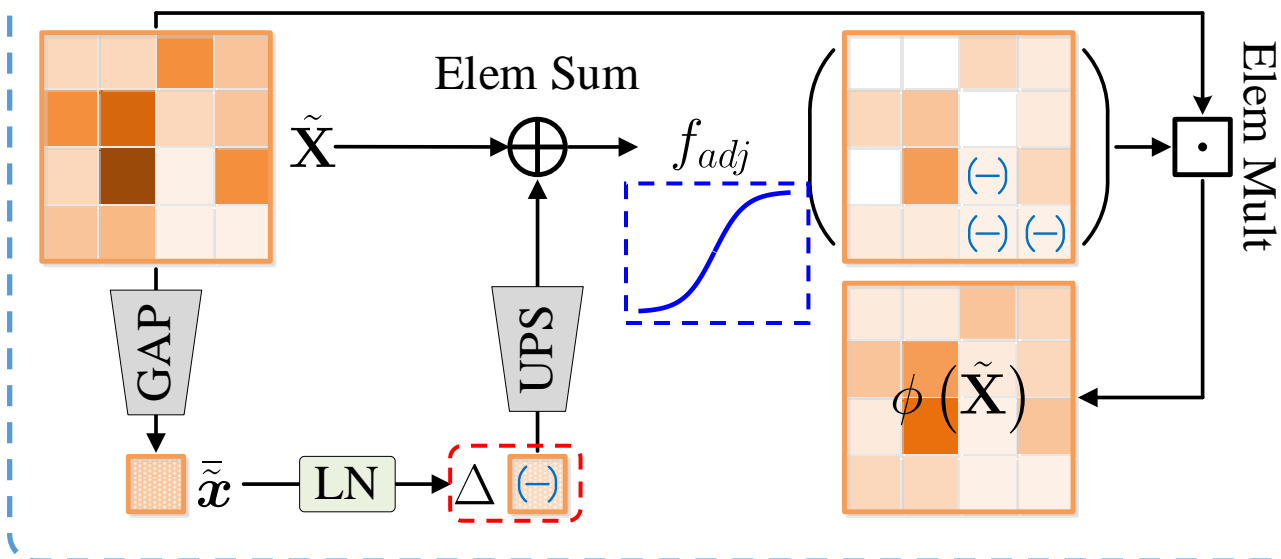
京都大学

# Method



**AdaShift Prototype**

$$\phi\left(\tilde{x}\right) = f_{\text{adj}}(\tilde{x} + \Delta) \cdot \tilde{x}$$

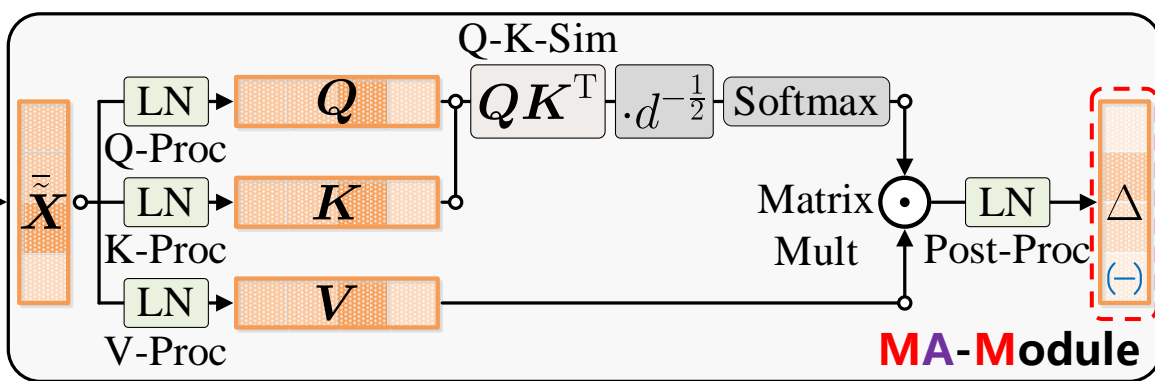Complementary information: tensor-level non-local cues

**AdaShift-B**

Elem Sum

$\tilde{X}$

$f_{adj}$

Elem Mult

$\phi\left(\tilde{\mathbf{X}}\right)$

GAP

UPS

$\bar{\tilde{x}}$ — LN — $\Delta$ (−)

RS Reshape   LAP Local Average Pooling   UPS Upsampling (Broadcast)

$\tilde{X}$

$$\Delta = \left[\text{MA}\left(\text{avgpool}_{K_H \times K_L}\left(\tilde{\mathbf{X}}\right)\right)\right]_c (h_k, l_k)$$

Elem Mult

$\phi\left(\tilde{\mathbf{X}}\right)$

LAP

$\tilde{\tilde{X}}$ — RS — $\tilde{\tilde{\mathbf{X}}}$

Q-Proc — LN — $\boldsymbol{Q}$

Q-K-Sim

$\boldsymbol{Q}\boldsymbol{K}^{\text{T}}$ · $d^{-\frac{1}{2}}$ Softmax

K-Proc — LN — $\boldsymbol{K}$

V-Proc — LN — $\boldsymbol{V}$

Matrix Mult

Post-Proc — LN — $\Delta$ (−)

RS UPS $\Delta$ (−)

**MA-Module**

Elem Sum

$f_{adj}$

**AdaShift-MA**

京都大学

# Experiment

- Dataset

  ◆ ImageNet [90]
  - The most popular large-scale visual (image) recognition benchmark dataset

  ◆ CIFAR(-100) [27]
  - A popular image recognition benchmark dataset of small-size images

  ◆ COCO [21]
  - A popular large-scale object detection benchmark dataset

  ◆ KITTI-Materials [23]
  - The benchmark dataset of RGB RMS

[90] J. Deng et al., Imagenet: A largescale hierarchical image database, *CVPR*, 2009
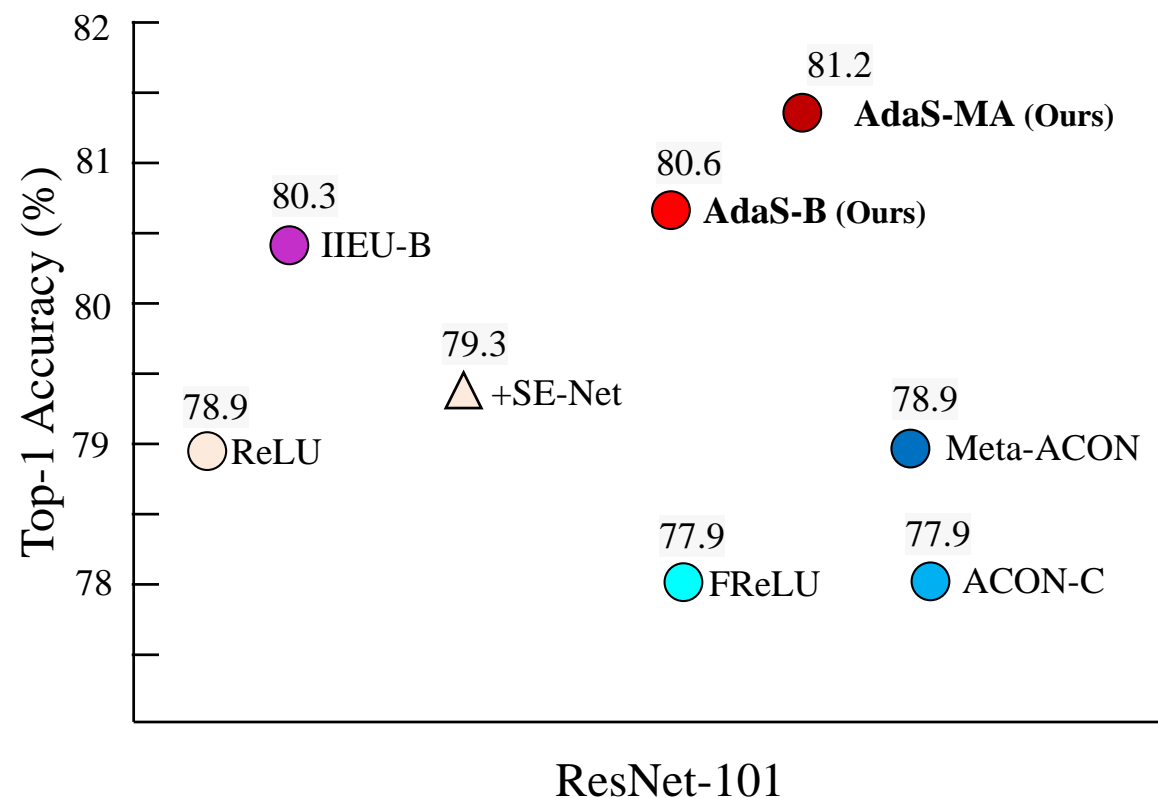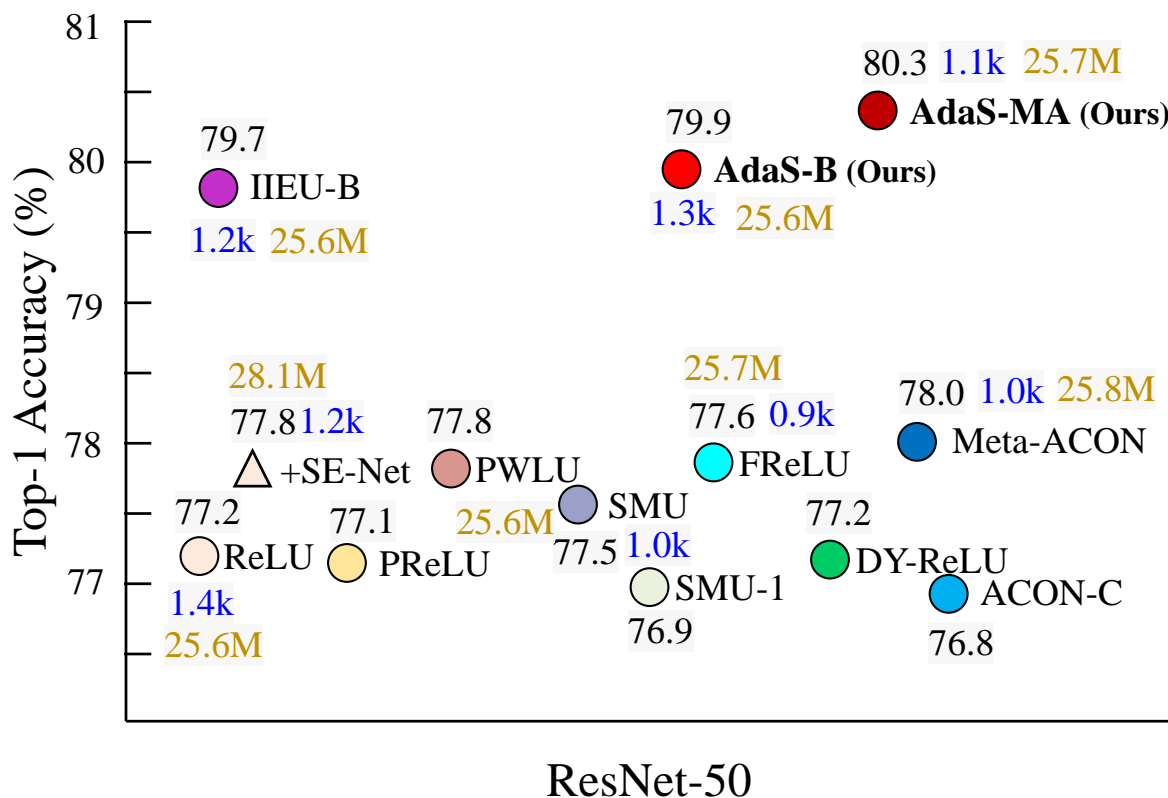
[27] A. Krizhevsky et al., Learning multiple layers of features from tiny images, *Master's thesis,* University of Toronto, 2009

[21] T.-Y. Lin et al., Microsoft coco: Common objects in context, *ECCV*, 2014

[23] Sudong Cai et al., RGB Road Scene Material Segmentation, *ACCV*, 2022

KYOTO UNIVERSITY

京都大学

# Results

Comparison of AdaShift-B and -MA with prevailing and (other) SOTA activation models on ImageNet using ResNet-50 and ResNet-101 backbones. In the comparison using the ResNet-50 backbone, the *number of parameters* and *throughput* of each model are indicated in Purple and Blue colors, respectively
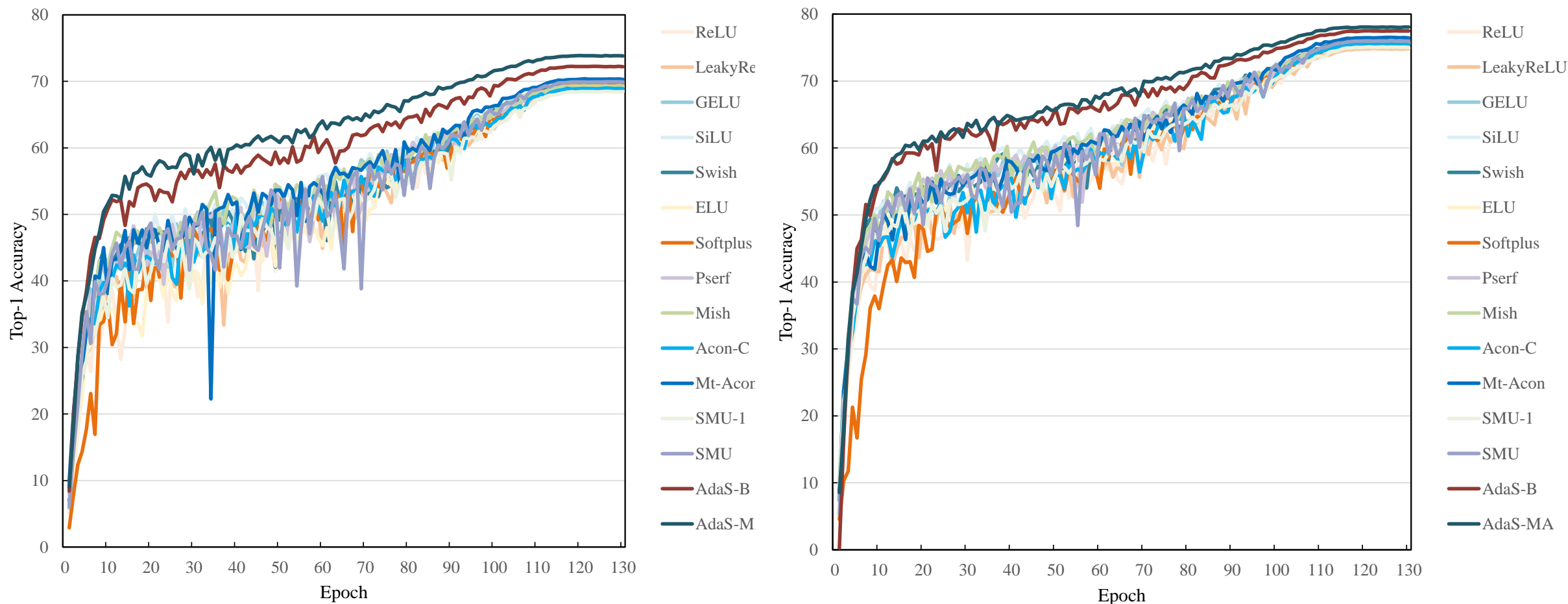


ResNet-50

ResNet-101

# Results

Comparison of AdaShift-enhanced ResNet-50(s) to representative vision Transformer counterparts. "✶" denotes the improved ViT trained with an extra regularization. Details can be found in Appendix G.2. (Training with Advanced Recipe)

| Network | Activation | Resolution | Train. Epoch | #Params. | FLOPs | Throughput | Top-1(%)↑ |
|---|---|---|---|---|---|---|---|
| ViT-B/16✶ | GELU | 224 × 224 | 300 | 86.6M | 16.9G | 775.6 | 79.7 |
| PoolFormer-S24 | GELU | 224 × 224 | 300 | 21.4M | 3.5G | 1144.6 | 80.3 |
| Swin-T | GELU | 224 × 224 | 300 | 28.3M | 4.5G | 1052.2 | **81.3** |
| ResNet-50 | **AdaS-B (ours)** | 224 × 224 | 300 | 25.6M | 4.1G | **1352.8** | 80.8 |
| ResNet-50 | **AdaS-Hyb (ours)** | 224 × 224 | 300 | 28.2M | 4.2G | **1201.6** | **81.7** |

京都大学

# Results

The Acc. curves (left) and loss curves (right) of ResNet-14 (Top) and ResNet-26 (Bottom) backbones with different activation models

# Gratitude for your time and patience!

## AdaShift: Learning Discriminative Self-Gated Neural Feature Activation With an Adaptive Shift Factor

**Sudong Cai**

Graduate School of Informatics, Kyoto University

KYOTO UNIVERSITY

京都大学