

# Cross-Architecture Adaptation: Cloud-Edge Continual Test-Time Adaptation with Dynamic Sampling and Heterogeneous Distillation

## Supplementary Material

### A. Detailed Implementation Settings

Due to space constraints in the main text, we provide further implementation details of our Cross-Architecture Adaptation (CAA) framework here. Both the cloud teacher (ViT-B) and the edge student (ResNet-18) are optimized using Stochastic Gradient Descent (SGD) with a momentum of 0.9.

For the edge-side distillation updates, we employ a dynamic learning rate schedule to prevent over-fitting to noisy pseudolabels. We initialize the learning rate at 0.03. To maintain stability, we monitor the cross-entropy loss  $\mathcal{L}_{CE}$  on the distillation batches. When  $\mathcal{L}_{CE}$  falls below the threshold  $\tau_{ce} = 0.23$  for  $T$  consecutive mini-batches, indicating that the local adaptation has converged for the current distribution, we decay the learning rate to 0.00025.

### B. Additional Experimental Results

#### B.1. Full Results on ImageNet-C (Severity Level 4)

In the main paper, we reported the comprehensive results for the most severe distribution shifts (Severity Level 5). In Table S1, we provide the complete breakdown of the adaptation accuracy across all 15 corruption types for Severity Level 4. Our CAA framework consistently outperforms existing state-of-the-art methods across the majority of corruptions.

#### B.2. Extensive Generalization Settings

To validate the generalization ability of CAA across diverse network architectures and datasets, we conducted additional experiments:

- **Model Diversity:** When deploying a heavier ResNet-50 on the edge while keeping the Cloud ViT-B unchanged,

CAA achieves 48.4% accuracy on ImageNet-C, outperforming the SOTA ETA (46.8%).

- **Homogeneous Architectures:** Our method is not strictly limited to heterogeneous setups. When utilizing a Cloud ViT-B and an Edge TinyViT (homogeneous Transformer structures), CAA achieves 46.6%, which is superior to the SOTA CEMA (43.5%). This demonstrates that CAA effectively coordinates inductive biases even in isomorphic scenarios.
- **Dataset Generalization:** We further evaluated CAA on the ImageNet-R dataset, which features diverse renditions of ImageNet classes. CAA achieves 28.3%, outperforming CEMA (27.6%).

### C. Hyperparameter Sensitivity

The hyperparameters introduced in our framework demonstrate strong robustness. We specifically analyze the crucial balancing parameter  $\alpha$  (which weights sample uncertainty against representativeness in the MDSCS module) and the cross-entropy threshold  $\tau_{ce}$  used for the learning rate schedule.

$\alpha$	0.1	0.3	<b>0.5</b>	0.7	0.9
Acc(%)	39.4	40.1	<b>41.2</b>	40.8	40.4
$\tau_{ce}$	0.10	0.20	<b>0.23</b>	0.30	0.40
Acc(%)	40.5	41.0	<b>41.2</b>	40.8	40.6

Table S2. Sensitivity analysis of hyperparameters  $\alpha$  and  $\tau_{ce}$  on ImageNet-C.

As shown in Table S2, we control for a comparable number of selected samples to ensure a fair comparison. CAA consistently outperforms state-of-the-art methods across different  $\alpha$  values, maintaining high stability with minimal

Severity Level=4	Noise			Blur				Weather				Digital				Avg.
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elastic	Pixel	JPEG	
ResNet18 (baseline)	7.8	6.2	6.5	18.6	12.1	16.2	22.1	17.5	21.9	28.4	57.9	14.2	39.3	26.4	43.6	22.6
• BN Adaptation <sup>†</sup>	29.8	24.9	28.2	26.8	25.7	33.0	39.2	31.4	35.0	50.5	62.2	44.6	54.4	49.0	49.0	38.9
• CoTTA	29.2	25.6	27.1	21.2	23.3	34.7	41.4	32.8	35.4	53.2	62.5	43.8	54.9	51.6	50.1	39.1
• ETA	41.4	41.9	42.0	32.9	37.0	40.3	44.9	39.7	39.9	<b>53.7</b>	59.5	<u>50.4</u>	56.3	54.1	53.1	45.6
• CEMA	<u>42.2</u>	<u>43.4</u>	<u>43.0</u>	32.0	<u>39.0</u>	<u>43.0</u>	<u>47.5</u>	<u>40.9</u>	<u>42.1</u>	<b>53.7</b>	<u>60.0</u>	49.2	<b>57.4</b>	<u>55.1</u>	<b>54.9</b>	<u>46.3</u>
• CoLA	41.2	42.0	42.2	<b>34.4</b>	38.4	40.8	45.0	38.7	40.2	53.0	58.9	<b>50.6</b>	56.6	54.4	53.5	45.3
• CAA(Ours)	<b>43.3</b>	<b>44.4</b>	<b>44.5</b>	<u>33.0</u>	<b>40.3</b>	<b>43.5</b>	<b>48.0</b>	<b>42.9</b>	<b>43.1</b>	<b>53.7</b>	<b>60.4</b>	49.0	<u>57.3</u>	<b>55.7</b>	<u>54.7</u>	<b>47.6</b>

Table S1. Comparisons with state-of-the-art methods on ImageNet-C (severity level 4) regarding **Accuracy (%)**. We adopt ViT-B as the foundation model and ResNet18 as the edge model. <sup>†</sup> denotes the TTA method that does not require any backward propagation and can be locally executed in edge devices.

deviation within the range of  $[0.3, 0.7]$ . Furthermore, the performance remains stable across various  $\tau_{ce}$  values, peaking at  $\tau_{ce} = 0.23$ .

#### **D. Computational Efficiency**

Although MDCS requires a dual forward pass to estimate uncertainty, the resulting precise sample selection drastically reduces the computational overhead associated with backpropagation. As a result, the per-epoch adaptation time on ImageNet-C for CAA is 49.6 ms, which is highly comparable to a full-sample update strategy without any selection (46.3 ms). Furthermore, our online clustering strategy involves a feature- and centroid-assignment process with a computational complexity of  $\mathcal{O}(N \times d)$  (where  $d = 512$ ), ensuring that the edge-side computation remains minimal.