

# SNAP-TTA: SPARSE TEST-TIME ADAPTATION FOR LATENCY-SENSITIVE APPLICATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Test-Time Adaptation (TTA) methods use unlabeled test data to dynamically adjust models in response to distribution changes. However, existing TTA methods are not tailored for practical use on edge devices with limited computational capacity, resulting in a latency-accuracy trade-off. To address this problem, we propose SNAP-TTA, a sparse TTA framework that significantly reduces adaptation frequency and data usage, delivering latency reductions proportional to adaptation rate. It achieves competitive accuracy even with an adaptation rate as low as 0.01, demonstrating its ability to adapt infrequently while utilizing only a small portion of the data compared to full adaptation. Our approach involves (i) Class and Domain Representative Memory (CnDRM), which identifies key samples that are both class-representative and domain-representative to facilitate adaptation with minimal data, and (ii) Inference-only Batch-aware Memory Normalization (IoBMN), which leverages representative samples to adjust normalization layers on-the-fly during inference, aligning the model effectively to changing domains. When combined with five state-of-the-art TTA algorithms, SNAP-TTA maintains the performances of these methods even with much-reduced adaptation rates from 0.01 to 0.5, making it suitable for edge devices serving latency-sensitive applications.

## 1 INTRODUCTION

Deep learning models often suffer from performance degradation under domain shifts caused by environmental changes or noise (Quiñonero-Candela et al., 2008). Test-Time Adaptation (TTA) offers a promising solution for domain shifts by utilizing only unlabeled test data without requiring source data. While TTA algorithms have advanced in complexity to improve accuracy in data streams (Wang et al., 2021; Niu et al., 2022; Wang et al., 2022; Yuan et al., 2023; Niu et al., 2023; Song et al., 2023), they are typically designed for resource-rich servers, overlooking the computational and memory limitations crucial for real-world deployment. Operations such as backpropagation, data augmentation, and model ensembling (Wang et al., 2022; Yuan et al., 2023; Zhang et al., 2022) result in substantial latency and memory consumption, making state-of-the-art (SOTA) TTA methods inefficient for practical use (Section 2).

For edge devices with limited computational power, such as mobile devices or IoT sensors, the adaptation latency from TTA methods becomes a critical bottleneck, particularly in latency-sensitive applications such as autonomous driving and real-time health monitoring. Moreover, the model must keep up with the data stream in those applications, but high computational overhead could cause it to miss critical samples, resulting in inference lags and reduced accuracy. This issue is exacerbated with fast data streams, such as high-frame-rate videos or high-performance sensors. For example, even a slight delay in processing sensor data can lead to dangerous situations in autonomous driving. A high adaptation latency that accumulates with each batch not only undermines real-time performance but also limits the potential of TTA algorithms in latency-sensitive applications.

In online TTA scenarios that require rapid response to incoming data streams on resource-constrained devices, *Sparse TTA (STTA)*, which adapts occasionally rather than at every batch, can offer a practical solution by reducing the adaption overhead. However, naïve STTA may result in performance degradation as it utilizes far less data (e.g., 0.1) for model adaptation (Figure 1). The

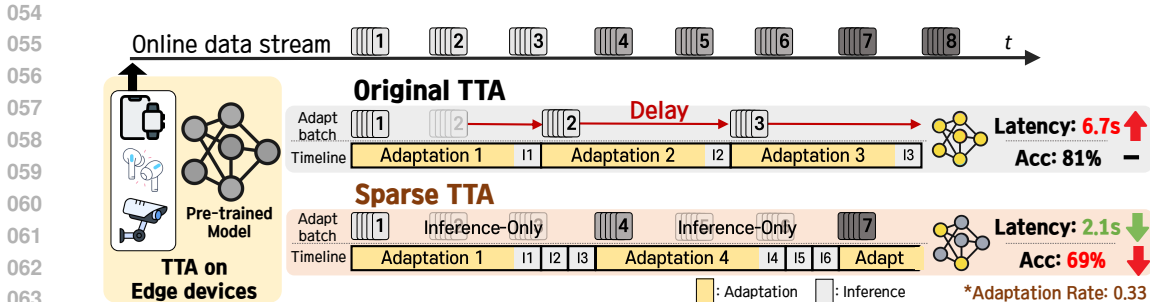


Figure 1: Comparison of average latency per batch and classification accuracy between the Original TTA and Sparse TTA approaches on edge devices processing an online data stream. With an adaptation rate of 0.33, adaptation occurs once every three batches, reducing latency relative to the adaptation rate but leading to a significant accuracy drop than fully adapting original TTA.

effectiveness of STTA hinges on selecting proper samples from a large pool, ensuring that the model maintains adequate performance with fewer updates (detailed analysis in Section 4).

Conventional TTA approaches that adopt sampling strategies are designed for non-i.i.d data (Gong et al., 2022; Niu et al., 2023; Yuan et al., 2023) or noisy data (Gong et al., 2023). They do not aim for data efficiency and thus yield high sample usage for updates. While EATA (Niu et al., 2022) excludes unreliable samples and utilizes fewer samples, it suffers from performance degradation when attempting more aggressive reductions. Data-efficient deep learning demonstrated that selecting easy, class-representative samples is effective when the sampling ratio is low (e.g., below 0.4) (Xia et al., 2022; Choi et al., 2024). However, these methods rely on ground-truth label information, which is typically unavailable in TTA scenarios.

We propose **SNAP-TTA: Sparse Network Adaptation for Practical Test-Time Adaptation**, a low-latency TTA framework designed for resource-constrained devices. SNAP-TTA addresses the challenge of balancing adaptation accuracy with computational efficiency in STTA, where only a small subset of data is used for updates. To that end, SNAP-TTA has two key technical enablers: First, it introduces a sampling strategy that combines *class-representative* and *domain-representative* samples. This approach enables the model to adapt effectively to domain shifts even with minimal data. Class and Domain Representative Memory (CnDRM) selects these critical samples by using pseudo-label confidence in a prediction-balanced manner for class-representative samples, and by identifying the domain-representative samples closest to the center of the target domain’s feature embedding (Section 3.1). Second, Inference-only Batch-aware Memory Normalization (IoBMN) refines the normalization process during inference by utilizing CnDRM’s class-domain representative statistics, leveraging the representativeness of these selected samples to correct skewed feature distributions at each inference step. This ensures that the model effectively adapts to domain shifts without back-propagation, maintaining alignment with the evolving data distribution (Section 3.2). These two components are integrated to perform adaptation, minimizing accuracy drop and latency in real-world domain-shifted scenarios.

SNAP-TTA is designed to work together with existing TTA methods orthogonally; thus, we evaluated SNAP-TTA integrated with existing SOTA TTA algorithms under diverse adaptation rates. Specifically, we evaluated SNAP-TTA with five SOTA TTA algorithms (Tent(Wang et al., 2021), EATA(Niu et al., 2022), SAR(Niu et al., 2023), CoTTA(Wang et al., 2022), and RoTTA(Yuan et al., 2023)) on three common TTA benchmarks (CIFAR10-C, CIFAR100-C (Hendrycks & Dietterich, 2019a), and ImageNet-C (Hendrycks & Dietterich, 2019b)). SNAP-TTA effectively reduces latency while minimizing performance drops in existing TTA methods. For instance, on our implementation in Raspberry Pi 4(Raspberry Pi Foundation, 2019) testbed, SNAP-TTA achieved up to 87.5% latency reduction at an adaptation rate of 0.1. In CIFAR10-C, SNAP-TTA-integrated methods consistently outperformed their original counterparts, showing up to 13.38% accuracy gain for CoTTA at an adaptation rate of 0.01. In addition, SNAP-TTA integration performed comparable accuracy to the original TTA methods under full adaptation settings. For instance, it achieved 77.12%~81.74% accuracy for Tent at various adaptation rates, whereas the full adaptation accuracy was 80.43% in CIFAR10-C.

## 2 PRELIMINARIES

We focus on the Test-Time Adaptation (TTA) latency challenges specific to edge devices, highlighting the constraints of adapting models in real-time environments with limited resources. Detailed related works are in Appendix A.

**Test-Time Adaptation and Its Latency Challenge on Edge Devices.** In unsupervised domain adaptation, the source domain data  $\mathcal{D}_S = \mathcal{X}^S, \mathcal{Y}$  is drawn from the distribution  $P_S(\mathbf{x}, y)$ , while the target domain data  $\mathcal{D}_T = \mathcal{X}^T, \mathcal{Y}$  follows  $P_T(\mathbf{x}, y)$ , typically without known labels  $y_j$ . Given a pre-trained model  $f(\cdot; \Theta)$  on the source domain  $\mathcal{D}_S$ , test-time adaptation (TTA) (Wang et al., 2021) adjusts the model to the target distribution  $P_T$  using only target instances  $\mathbf{x}_j$ , updating the parameters  $\Theta$  to reduce domain discrepancy.

When applied to resource-constrained devices, however, current TTA approaches face significant latency challenges. In real-time applications that require rapid inference, online TTA becomes impractical due to the need for adaptation at every batch (Figure 4, detailed latency tracking reported in Appendix E.3). Our experiment on Raspberry Pi 4 (Raspberry Pi Foundation, 2019) showed a minimum of 3.83 seconds latency per batch for existing TTA methods. This indicates existing methods could not handle real-time applications with fast data streams and strict latency requirements, such as autonomous driving (Tampuu et al., 2024; Liu et al., 2023). TTA methods such as CoTTA use computationally intensive operations such as data augmentations and ensemble models at the cost of increased latency. Relatively lightweight algorithms incur non-negligible latency from adaptation processes such as backpropagation, which becomes bottlenecks in resource-constrained devices without the parallel processing capabilities and memory bandwidth of GPUs.

A recent work (Alfarra et al., 2024), recognizing latency as a problem, proposed a TTA evaluation protocol that penalizes methods that are slower than the data stream rate. Instead of penalizing a model for being slow, we utilize Sparse TTA, where the model actively chooses to adapt at sparse intervals for the goal of maintaining a real-time inference rate. As real deployments involve devices with different computational capabilities and data streams of varying speeds, we believe a framework that effectively maintains various TTA methods’ performance across different latency requirements is crucial.

**Sparse Test-Time Adaptation and Adaptation rates.** Sparse Test-Time Adaptation (STTA) aims to efficiently adapt models by reducing both the frequency of updates and the number of samples used per update, which is essential for minimizing latency in edge devices. The concept of adaptation rate plays a central role in STTA, as it controls both the update frequency and the number of data points used. Unlike Original Test-Time Adaptation (TTA), which uses full batches of data and can create significant computational overhead, STTA employs an adaptation rate to limit updates and data usage proportionally, thus introducing sparsity (Figure 1).

By adjusting the *adaptation rate*, STTA can minimize latency and computational costs while maintaining adaptation performance. This rate defines how sparsely updates occur and the proportion of samples used for updates compared to the Original TTA, enabling efficient model adjustments to distribution shifts. The balance between adaptation accuracy and computational efficiency makes STTA particularly suitable for environments that demand both quick responses and minimal resource usage.

## 3 METHODOLOGY

SNAP-TTA framework resolves the high latency and inefficiency issue of existing Test-Time Adaptation (TTA) methods. By introducing a Sparse TTA (STTA) strategy combined with a novel sampling method, SNAP-TTA minimizes adaptation delays while maintaining accuracy. The overall system, illustrated in Figure 2, consists of two primary components: (i) Class and Domain Representative Memory (CnDRM) for efficient sampling and (ii) Inference-only Batch-aware Memory Normalization (IoBMN) to correct feature distribution shifts during inference. Together, these components enable effective STTA with minimal computational overhead.

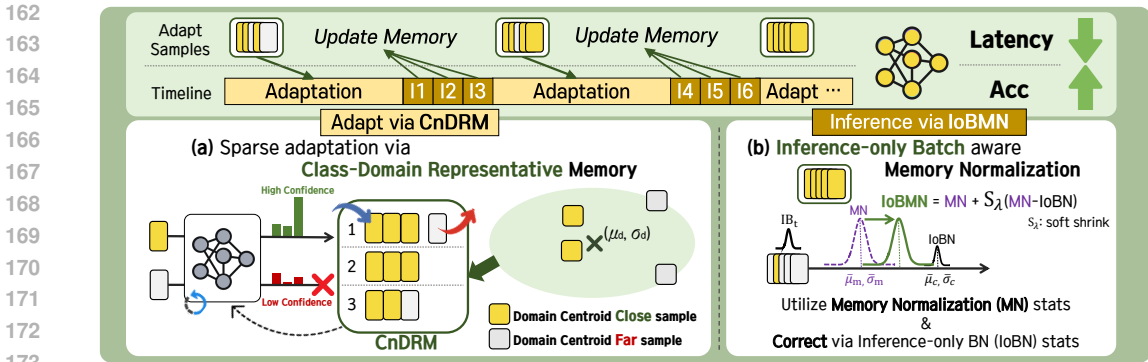


Figure 2: Design overview of SNAP-TTA. The framework consists of two primary components: (a) Class and Domain Representative Memory (CnDRM), which efficiently selects representative samples to minimize adaptation overhead, and (b) Inference-only Batch-aware Memory Normalization (IoBMN), which corrects feature distribution shifts during inference. Together, these components implement the Sparse TTA (STTA) strategy, reducing latency while maintaining model accuracy.

### 3.1 CLASS AND DOMAIN REPRESENTATIVE MEMORY (CnDRM)

CnDRM is a core component of SNAP-TTA that addresses the challenges of efficient data sampling for STTA. In STTA, the adaptation rate directly impacts the number of samples used, necessitating a careful sampling strategy to optimize performance with minimal data. Given this limited sampling ratio, CnDRM selects both class and domain-representative samples to maintain model performance while minimizing adaptation overhead.

**Motivation.** Data sampling is crucial in data-efficient deep learning, especially when working with a limited number of samples. In high data sampling ratio scenarios, score-based methods prioritize difficult or rare samples, often achieving performance comparable to full-dataset training. However, when the sampling ratio is low, selecting easy and class-representative samples becomes more effective (Choi et al., 2024). This method selects samples that minimize differences in loss gradients or curvature, ensuring that the generalizability is retained even with fewer samples. Similarly, the Moderate Coreset (Xia et al., 2022) paper demonstrates that at low sampling ratios of 0.2 to 0.4, the distance from the class center significantly impacts performance, with samples closer to the center being particularly effective in scenarios with high label noise. In the STTA setting, where ground truth labels are unavailable and the probability of incorrect predictions is high, selecting representative samples based on potentially incorrect predictions resembles a high label noise situation. Therefore, selecting class-representative easy samples could provide some benefit to STTA.

However, if the model must perform STTA at an even lower adaptation rate (e.g., 0.1) due to the latency limits, selecting class-representative samples alone would be insufficient (Table 4). Unlike traditional classification tasks, STTA is an unsupervised domain adaptation, which requires identifying target domain-representative samples that reflect the distributional shift between the source and target domains. In these cases, we argue that focusing on domain-representative instances is just as crucial, as selecting samples that best capture the domain shift can help the model retain generalizability with minimal data. Therefore, selecting both **class-representative and domain-representative samples** could enhance STTA performance in low-data environments, where each sample must contribute significantly to model adaptation.

**Criteria 1: Class Representation.** CnDRM selects samples with higher confidence scores to avoid the issues caused by low-confidence samples. Low-confidence samples are typically located near decision boundaries and are more likely to carry incorrect pseudo-labels. This strategy ensures that the adaptation process is guided by stable learning signals, which is important in the absence of ground-truth labels. By focusing on high-confidence samples, CnDRM mitigates the risk of propagating errors resulting from incorrect pseudo-labels, thereby supporting more effective and stable adaptation (Details in Appendix E.2). The confidence score  $C(\mathbf{x})$  for each sample  $\mathbf{x}$  is calculated as:  $C(\mathbf{x}) = \max_{y \in \mathcal{Y}} p(y|\mathbf{x}; \Theta)$  where  $p(y|\mathbf{x}; \Theta)$  is the softmax probability for class  $y$ . Only samples with confidence above a predefined threshold  $\tau_{conf}$  are retained. For a balanced representa-

tion across diverse classes, CnDRM selects these high-confidence samples in a prediction-balanced manner. This balance helps maintain the model’s overall classification capability and prevents bias towards certain classes when only a low sample ratio is available for adaptation. By leveraging both high confidence and prediction balance, CnDRM effectively selects class-representative samples that are diverse and reliable, even without access to ground-truth labels.

**Criteria 2: Domain Representation.**

In addition to class-representative sampling, CnDRM selects domain-representative samples to facilitate adaptation to new domain conditions. Building on the efficient class-representative sampling criteria, we argue that *selecting samples close to the domain centroid* would enhance performance in STTA. Our preliminary experiment results validate improved performance when selecting samples near the centroid (Figure 3).

For ImageNet-C Gaussian noise, TTA with the closest 20% of samples achieved 26.65% accuracy, whereas the farthest 20% showed a lower accuracy of 18.52%.

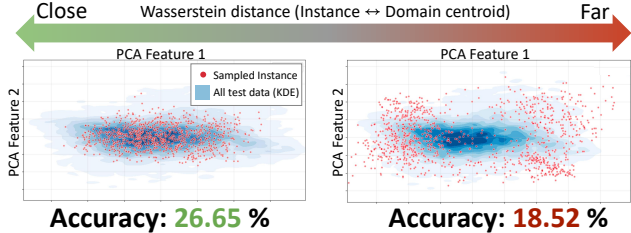


Figure 3: Sampling visualization and accuracy comparison between the closest 20% and farthest 20% samples from the domain centroid (based on Wasserstein distance) on ImageNet-C Gaussian noise.

As early layers in deep learning models tend to retain domain-specific features (Zeiler & Fergus, 2014; Lee et al., 2018; Segu et al., 2023), we utilize the hidden features of early layers to identify domain-representative samples (Appendix E.1). We use the feature statistics (mean and variance) of the first normalization layer to evaluate domain representation. This choice is made as domain discrepancies can be effectively reduced through normalization adjustments (Nado et al., 2020; Schneider et al., 2020). Domain discrepancies in hidden features are substantially reduced after passing through a single normalization layer, significantly minimizing domain shift differences (Li et al., 2016). While deeper layers provide detailed information, using the first layer balances capturing domain-specific information and maintaining computational efficiency.

The domain centroid  $\mathbf{c}_{domain}$  is computed using a momentum-based update of batch statistics from the normalization layer:  $\mu_{domain} \leftarrow (1 - \beta)\mu_{domain} + \beta\mu_t$  and  $\sigma_{domain}^2 \leftarrow (1 - \beta)\sigma_{domain}^2 + \beta\sigma_t^2$ , where  $\mu_t$  and  $\sigma_t^2$  are the mean and variance of the current batch  $t$ , and  $\beta$  is the momentum parameter. In our preliminary study, we found that using only the mean and standard deviation values before the first normalization was sufficient to calculate the domain centroid. The sampled instances effectively represented the domain and were correctly positioned in the embedding space for each criterion (Figure 3).

To determine domain-representative samples, CnDRM calculates the Wasserstein distance between each sample’s feature statistics and the domain centroid. The Wasserstein distance measures the similarity between two distributions by considering their mean and variance, evaluating how well a sample represents the domain. It is useful for capturing domain characteristics, leading to its wide use in domain generalization (Segu et al., 2023). For each sample  $\mathbf{x}_t$ , the feature statistics  $(\mu_{\mathbf{x}_t}, \sigma_{\mathbf{x}_t})$  are taken from the input to the normalization layer, and the Wasserstein distance  $W(\mathbf{x}_t, \mathbf{c}_{domain})$  is given by:

$$W(\mathbf{x}_t, \mathbf{c}_{domain}) = \sqrt{(\mu_{\mathbf{x}_t} - \mu_{domain})^2 + (\sigma_{\mathbf{x}_t} - \sigma_{domain})^2}. \tag{1}$$

**Memory Management Algorithm.** The memory management in CnDRM maintains efficiency without introducing additional overhead. To achieve this, the memory size is kept equal to the batch size for minimal resource usage. Within this fixed memory, samples are managed by balancing the number of samples per class based on predictions so that each class remains well-represented. For domain adaptation, samples in memory are periodically replaced with new samples that are closer to the domain centroid and meet the confidence threshold to retain only the most class-domain representative samples. Algorithm 1 has details.

**Algorithm 1** Class and Domain Representative Memory (CnDRM)

---

**Require:** test data stream  $x_t$ , memory  $M$  with capacity  $N$ , confidence threshold  $\tau_{conf}$ , sample unit for memory  $s$ , adaptation rate  $1/k$

- 1: **for** batch  $b \in \{1, \dots, B\}$  **do**
- 2:    $\hat{Y}_b \leftarrow f(b; \Theta)$
- 3:   **for** each sample  $x_t$  in batch  $b$  **do**
- 4:      $\hat{y}_t \leftarrow \hat{Y}_b[t]$
- 5:     confidence  $\leftarrow C(x_t; \Theta)$
- 6:      $\mathbf{c}_t(\mu_{\mathbf{x}_t}, \sigma_{\mathbf{x}_t}) \leftarrow$  mean and variance of early hidden feature
- 7:      $w_{x_t} \leftarrow W(x_t, \mathbf{c}_{domain})$
- 8:     **if** confidence  $> \tau_{conf}$  **then** ▷ Class-representative samples
- 9:       Add  $\mathbf{s}_t(x_t, \hat{y}_t, \mathbf{c}_t, w_{x_t})$  to  $M$  ▷ Add samples in prediction-balanced manner
- 10:     **if**  $|M| > N$  **then**
- 11:        $L^* \leftarrow$  class with most samples in  $M$
- 12:       **if**  $\hat{y}_t \notin L^*$  **then** ▷ Removes domain-centroid farthest sample
- 13:          $\mathbf{s}_{\max\_dist} \leftarrow \arg \max_{\mathbf{s}_i \in M \wedge \hat{y}_i \in L^*} w_{x_i}$
- 14:       **else**
- 15:          $\mathbf{s}_{\max\_dist} \leftarrow \arg \max_{\mathbf{s}_i \in M \wedge \hat{y}_i = \hat{y}_t} w_{x_i}$
- 16:       Remove  $\mathbf{s}_{\max\_dist}$  from  $M$
- 17:      $\mathbf{c}_{domain} \leftarrow (1 - \beta)\mathbf{c}_{domain} + \beta\mathbf{c}_t$  ▷ Update domain-centroid
- 18:     Recalculate  $w_{s_i}$  for all  $\mathbf{s}_i$  in  $M$
- 19:     **if**  $b \bmod k == 0$  **then** ▷ Adaptation occurs every  $k$  batches
- 20:       Update model  $\Theta$  using samples in  $M$

---

## 3.2 INFERENCE-ONLY BATCH-AWARE MEMORY NORMALIZATION (IoBMN)

**Motivation.** In Sparse Test-Time Adaptation (STTA) scenarios, models must adapt to domain shifts despite having limited opportunities for updates. In this setting, maintaining robust performance becomes challenging as the stored memory statistics, derived from representative adaptation batches, may not fully align with subsequent inference batches, especially when updates are skipped. This can lead to a potential mismatch between the stored statistics and the current data distribution. Traditional normalization methods, which solely rely on test batches’ statistics, struggle to address these subtle shifts effectively. To tackle this issue, we introduce the Inference-only Batch-aware Memory Normalization (IoBMN) module, which leverages the robustness of class-domain representative statistics while dynamically adjusting for mismatches that arise in skipped batches. By primarily basing normalization on stable, representative memory statistics and selectively adapting with recent inference data, IoBMN efficiently corrects for distributional shifts, ensuring both robustness and adaptability in STTA conditions. This approach significantly enhances model stability in sparse adaptation scenarios, as shown in our ablation study in Section 4.

**Approach.** Given a feature map  $f \in \mathbb{R}^{B \times C \times L}$ , where  $B$  is the batch size,  $C$  is the number of channels, and  $L$  is the number of spatial locations, the batch-wise statistics  $\bar{\mu}_c$  and  $\bar{\sigma}_c^2$  for the  $c$ -th channel are calculated as follows:

$$\bar{\mu}_c = \frac{1}{B \times L} \sum_{b=1}^B \sum_{l=1}^L f_{b,c,l}, \quad \bar{\sigma}_c^2 = \frac{1}{B \times L} \sum_{b=1}^B \sum_{l=1}^L (f_{b,c,l} - \mu_{b,c})^2, \quad (2)$$

where  $\bar{\mu}_m$  and  $\bar{\sigma}_m^2$  are calculated from the most recent adapted CnDRM samples in the same way with Equation 2, using the memory capacity  $M$  with  $m$  representing the memory. We assume that  $\mu_m$  and  $\sigma_m^2$  follow the *sampling distribution* of the feature map size  $L$  and memory capacity  $M$ . The corresponding variances for the memory mean  $\mu_m$  and variance  $\sigma_m^2$  are calculated as:

$$s_{\mu_m}^2 := \frac{\bar{\sigma}_m^2}{C \times M}, \quad s_{\sigma_m^2}^2 := \frac{2\bar{\sigma}_m^4}{C \times M - 1}. \quad (3)$$

For the normalization process to adapt efficiently to the current inference batch statistics, IoBMN corrects  $(\bar{\mu}_m, \bar{\sigma}_m^2)$  only when  $\bar{\mu}_c$  (and  $\bar{\sigma}_c^2$ ) significantly differ from  $\bar{\mu}_m$  (and  $\bar{\sigma}_m^2$ ) through soft shrinkage function:

$$\mu_m^{\text{IoBMN}} = \bar{\mu}_m + S_\lambda(\bar{\mu}_c - \bar{\mu}_m; \alpha s_{\mu_m}), \quad (\sigma_m^{\text{IoBMN}})^2 = \bar{\sigma}_m^2 + S_\lambda(\bar{\sigma}_c^2 - \bar{\sigma}_m^2; \alpha s_{\sigma_m^2}), \quad (4)$$

where  $\alpha \geq 0$  in IoBMN controls the reliance on the normalization layer statistics. A larger  $\alpha$  gives more weight to the last adapted memory normalization statistics, whereas a smaller  $\alpha$  emphasizes the current inference batch normalization statistics. The soft shrinkage function  $S_\lambda(x; \lambda)$  is defined as:

$$S_\lambda(x; \lambda) = \begin{cases} x - \lambda & \text{if } x > \lambda, \\ x + \lambda & \text{if } x < -\lambda, \text{ and} \\ 0 & \text{otherwise,} \end{cases}$$

where  $\lambda$  is the threshold,  $s$  is a scaling factor, and  $x$  is the input. The function allows for proportional adjustments based on the magnitude of the values, where smaller values are adjusted less and larger values more, preserving the critical information inherent in the adapted memory normalization statistics.

Finally, the output of the IoBMN for each feature  $f_{b,c,l}$  is computed as:

$$\text{IoBMN}(f_{b,c,l}; \bar{\mu}_m, \bar{\sigma}_m^2, \mu_m^{\text{IoBMN}}, (\sigma_m^{\text{IoBMN}})^2) := \gamma \cdot \frac{f_{b,c,l} - \mu_m^{\text{IoBMN}}}{\sqrt{(\sigma_m^{\text{IoBMN}})^2 + \epsilon}} + \beta, \quad (5)$$

where  $\gamma$  and  $\beta$  are learnable affine parameters of normalization layer, and  $\epsilon$  is a small constant added for numerical stability. In our experiments, we chose  $\alpha = 4$  to effectively handle various out-of-distribution scenarios. The parameter  $s$  is a hyperparameter that determines the degree of adjustment desired and can be tuned based on specific requirements.

IoBMN utilizes CnDRM’s class-domain representative statistics and adjusts them based on the current inferencing batch statistics. This dual-statistic approach allows IoBMN to correct the outdated and skewed distribution of the memory, ensuring alignment with the data distribution at each inference point. By leveraging the statistics of the data used during model update points, IoBMN adapts effectively without significant computational overhead. Additionally, this method mitigates the performance degradation caused by the prolonged intervals between adaptations so that the model remains well-aligned with the evolving data distribution.

## 4 EXPERIMENTS

This section outlines our experimental setup and presents the results obtained under various STTA settings. Refer to Appendix B for further details.

**Scenario.** We examined how different adaptation rates affect performance to simulate a scenario requiring a certain latency threshold for latency-sensitive applications. We varied the *adaptation rate* to observe its impact on both model accuracy and latency. The main evaluation was run with diverse adaptation rates (0.01, 0.03, 0.05, 0.1, 0.3, and 0.5). We report the average accuracy and standard deviation from three random seeds. Latency measurement was done on our Raspberry Pi 4 (Raspberry Pi Foundation, 2019) testbed.

**Dataset and Model.** We used three standard TTA benchmarks: **CIFAR10-C**, **CIFAR100-C** (Hendrycks & Dietterich, 2019a) and **ImageNet-C** (Hendrycks & Dietterich, 2019b). These datasets include 15 different types of corruption with five levels of severity, and we used the highest one. CIFAR10-C/CIFAR100-C has 10,000 test data with 10/100 classes, and ImageNet-C has 50,000 test data with 1,000 classes for each corruption. We employed **ResNet18** (He et al., 2016) as the backbone network, utilizing models pre-trained on CIFAR10 and CIFAR100 (Krizhevsky & Hinton, 2009). We also use **ResNet50** (He et al., 2016) and **ViT** (Dosovitskiy, 2020) pre-trained on ImageNet (Deng et al., 2009) from the TorchVision (maintainers & contributors, 2016) library.

**Baselines.** SNAP-TTA is designed to integrate with existing TTA algorithms. Therefore, testing existing *TTA algorithms under different adaptation rates* serves as our baseline (implementation details including hyperparameters are in Appendix B.1). We selected five SOTA TTA algorithms: (i) **Tent** (Wang et al., 2021) updates only BN affine parameters, (ii) **CoTTA** (Wang et al., 2022) updates the entire model parameters using a teacher-student framework, (iii) **EATA** (Niu et al., 2022), (iv) **SAR** (Niu et al., 2023), and (v) **RoTTA** (Yuan et al., 2023). For efficiency evaluation, we compared our method against **BN stats** (Nado et al., 2020; Schneider et al., 2020).





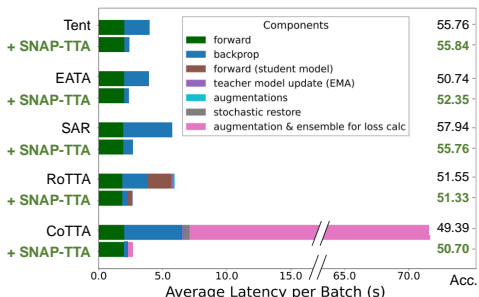


Figure 4: Latency and accuracy comparison of original TTA methods and their SNAP-TTA integration on CIFAR100-C. SNAP-TTA significantly enhances the efficiency.

Table 3: Latency reduction and accuracy gaps of SNAP-TTA (adaptation rate 0.1) compared by original TTA, tested on Raspberry Pi 4. Performance averaged over 15 CIFAR10-C corruptions. Numbers in parentheses represent the performance difference of SNAP-TTA compared to full adaptation.

Methods	Latency per batch (s)		Accuracy (%)	
	Original TTA	SNAP-TTA	naive STTA	SNAP-TTA
Tent	3.97	<b>2.20 (-44.0%)</b>	76.81 (-3.62)	<b>78.95 (-1.48)</b>
CoTTA	71.68	<b>8.96 (-87.5%)</b>	66.42 (-11.58)	<b>78.83 (+0.83)</b>
EATA	3.93	<b>2.18 (-44.6%)</b>	76.29 (-5.27)	<b>78.61 (-2.95)</b>
SAR	5.75	<b>2.30 (-60.1%)</b>	76.01 (-3.04)	<b>78.06 (-0.99)</b>
RoTTA	5.93	<b>2.25 (-62.0%)</b>	74.78 (-2.27)	<b>77.07 (+0.07)</b>

original counterparts, highlighting its adaptability and effectiveness. These results underscore the capability of SNAP-TTA to balance efficiency and performance, providing a significant advantage in sparse adaptation scenarios while maintaining or even enhancing classification accuracy. This validates the effectiveness of utilizing class-domain representative samples in the STTA setting.

Furthermore, Figure 5 shows more computationally complex and latency-intensive methods such as CoTTA tend to have greater performance gain when integrated with SNAP-TTA. This is because methods that update the entire model parameters are more susceptible to the influence of specific adaptation samples, leading to significant performance drops under sparse update conditions, which SNAP-TTA’s CnDRM and IoBMN effectively mitigate. In addition, adaptation rates of 0.5 or 0.3, which represent relatively high adaptation frequencies, sometimes can achieves even better performance with SNAP-TTA than the original TTA, despite in the STTA setting. This is likely because the sampling rate was not critically low but rather comparable to that of existing data-efficient methods such as EATA (Niu et al., 2022), allowing SNAP-TTA to achieve performance gains similar to various sampling-based TTA methods (Niu et al., 2022; 2023; Gong et al., 2022; 2023) using fewer yet effective samples. Overall, SNAP-TTA significantly reduced the average latency per batch while effectively maintaining accuracy, highlighting its benefits for resource-constrained environments. More details on all other adaptation rates are reported in Appendix C.

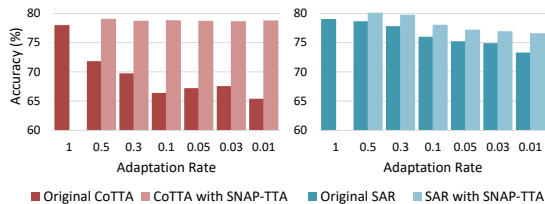


Figure 5: Classification accuracy on CIFAR10-C with varying adaptation rates. SNAP-TTA consistently mitigates accuracy drop across all rates.

Table 4: Classification accuracy (%) comparison of ablative settings on the STTA (adaptation rate 0.1). Performance averaged over 15 CIFAR10-C corruptions.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
Naive	76.81	66.42	76.29	76.01	74.78
Random	77.08	65.61	76.59	76.33	75.01
LowEntropy	75.66	63.19	74.89	74.41	72.60
CRM	77.77	65.71	77.18	74.36	75.27
CnDRM	77.46	77.69	77.17	76.85	75.64
CnDRM+EMA	78.02	72.19	77.05	76.84	76.18
CnDRM+IoBMN	<b>78.95</b>	<b>78.83</b>	<b>78.61</b>	<b>78.06</b>	<b>77.07</b>

**Contribution of individual components of SNAP-TTA.** We conducted an ablative evaluation to understand the effects of the individual components of SNAP-TTA (Table 4; more results on diverse adaptation rates and datasets are on Appendix D). CRM denotes prediction-balanced sampling with a confidence threshold (same as the Class-Representative criteria of CnDRM), and CnDRM denotes both Class and Domain Representative sampling (the first component of SNAP-TTA). For inference, the default uses test batch normalization statistics, EMA uses the exponential moving average of the test batch, and IoBMN uses memory samples’ statistics corrected to match that of the test batch (the second component of SNAP-TTA).

Contrary to the hypothesis that low-entropy samples are beneficial for TTA (Niu et al., 2022; 2023), LowEntropy performed worse than Rand for STTA. This can be attributed to the limited updates of STTA, resulting in poor or longer convergence times due to low entropy minimization loss. CRM, originally used for data-efficient supervised deep learning (Choi et al., 2024; Xia et al., 2022), performed better than Rand. However, as CRM on TTA inevitably relies on uncertain pseudo labels instead of the ground truth, its performance remains lower than utilizing domain representative features (CnDRM) (note that TTA is unsupervised domain adaptation rather than training from scratch (Xia et al., 2022)). The highest accuracy was achieved when inference was performed us-

Table 5: Classification accuracy (%) on ImageNet-C through Adaptation Rate 0.1 using ViT-based model. **Bold** numbers are the highest accuracy.

Methods	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.
Tent	40.56	41.30	41.69	35.76	31.81	42.01	38.02	44.33	<b>53.53</b>	20.69	72.41	30.42	45.87	<b>51.95</b>	<b>56.11</b>	43.10
+ SNAP-TTA	<b>40.98</b>	<b>41.72</b>	<b>42.18</b>	<b>37.16</b>	<b>32.30</b>	<b>42.89</b>	<b>38.44</b>	<b>46.19</b>	52.50	<b>53.11</b>	<b>72.25</b>	<b>39.25</b>	<b>46.77</b>	51.53	55.99	<b>46.22</b>
EATA	20.12	21.52	21.40	20.90	23.42	15.71	18.00	16.12	28.35	22.24	35.97	11.33	19.78	20.22	19.99	21.00
+ SNAP-TTA	<b>40.74</b>	<b>43.22</b>	<b>43.11</b>	<b>40.63</b>	<b>44.59</b>	<b>51.58</b>	<b>50.63</b>	<b>54.77</b>	<b>58.32</b>	<b>61.50</b>	<b>73.91</b>	<b>33.85</b>	<b>60.19</b>	<b>63.35</b>	<b>63.01</b>	<b>52.23</b>
SAR	21.45	23.02	23.17	23.67	24.64	15.98	14.62	7.70	31.49	8.94	41.33	6.82	17.35	22.39	22.49	20.34
+ SNAP-TTA	<b>37.59</b>	<b>38.27</b>	<b>36.78</b>	<b>38.58</b>	<b>39.99</b>	<b>49.00</b>	<b>45.77</b>	<b>43.96</b>	<b>56.61</b>	<b>59.96</b>	<b>73.02</b>	<b>19.69</b>	<b>54.30</b>	<b>61.16</b>	<b>61.85</b>	<b>47.77</b>

ing IoBMN, which primarily utilizes memory statistics and only shifts slightly to the test batch on demand. These results collectively indicate that utilizing CnDRM and IoBMN of SNAP-TTA enhances performance in a low-latency STTA scenario.

**Validation of SNAP-TTA on Vision Transformer (ViT) based Model.** To validate the effectiveness of SNAP-TTA on the Vision Transformer (ViT) (Dosovitskiy, 2020), we conducted experiments on ImageNet-C with adaptation rate of 0.1. Since ViT uses layer normalization (LN), we adjusted CnDRM and IoBMN to use LN from instances, demonstrating that the core concepts of selecting domain-representative samples and mitigating shift in normalization statistics can be applied effectively to a different normalization type (details in Appendix F.3). The results in Table 5 confirm consistent accuracy gains of SNAP-TTA with significant latency decrease, regardless of model and normalization types.

## 5 DISCUSSION AND CONCLUSION

**Limitations and future work.** Our work could be optimized for more realistic data streams, such as continuous domain adaptation scenarios (Appendix F.2). For instance, the adaptation rate can be dynamically altered based on the need for adaptation (i.e., the data distribution just changed). Additionally, while SNAP-TTA employed a fixed confidence threshold in CnDRM as a safeguard to filter noisy samples, its adaptability could be improved. Dynamically adjusting the threshold based on data characteristics presents a promising direction for future research to enhance sampling efficiency and overall performance.

Moreover, while we focused on reducing adaptation latency, memory overhead is another concern. We note that SNAP-TTA introduces negligible additional memory overhead, as detailed in the Appendix E.4, where related analysis and tracking information from real-device experiments are provided. Additionally, we demonstrate in the Appendix E.5 that SNAP-TTA can be effectively used alongside memory-efficient TTA methods such as MECTA (Hong et al., 2023), showcasing its compatibility and practicality. Future works could further explore optimizing SNAP-TTA for both latency and memory.

**Conclusion** We raised the overlooked issue of latency of TTA methods, which is particularly relevant for applications on resource-constrained edge devices. To this end, we propose SNAP-TTA, a Sparse TTA (STTA) framework that could be applied to existing TTA methods to significantly reduce their latency while maintaining competitive accuracy. For effective performance in an STTA setting, we utilize class-domain representative memory of samples for adaptation. Furthermore, we optimize inference by adapting normalization layers using representative samples to account for domain shifts. Extensive experiments and ablative studies demonstrate SNAP-TTA’s effectiveness in latency and adaptation accuracy.

## REPRODUCIBILITY STATEMENT

Details of the experiments, including datasets, scenarios, and hyperparameters for reproducibility, are provided in the Appendix B. Additionally, we share the link (<https://anonymous.4open.science/r/SNAPTTA-DD0E>) of an anonymous repository containing our source code and instructions to validate the reproducibility.

## REFERENCES

- 540  
541  
542 Motasem Alfarrar, Hani Itani, Alejandro Pardo, Shyma Yaser Alhuwaider, Merey Ramazanova,  
543 Juan Camilo Perez, Zhipeng Cai, Matthias Müller, and Bernard Ghanem. Evaluation of test-  
544 time adaptation under computational time constraints. In Ruslan Salakhutdinov, Zico Kolter,  
545 Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.),  
546 *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Pro-  
547 ceedings of Machine Learning Research*, pp. 976–991. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/alfarra24a.html>.  
548
- 549 Hoyong Choi, Nohyun Ki, and Hye Won Chung. Bws: Best window selection based on sample  
550 scores for data pruning across broad ranges. *arXiv preprint arXiv:2406.03057*, 2024.  
551
- 552 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hier-  
553 archical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*,  
554 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 555 Alexey Dosovitskiy. An image is worth 16x16 words: Transformers for image recognition at scale.  
556 *arXiv preprint arXiv:2010.11929*, 2020.  
557
- 558 Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimiza-  
559 tion for efficiently improving generalization. In *International Conference on Learning Representa-*  
560 *tions*, 2021. URL <https://openreview.net/forum?id=6Tm1mposlrM>.
- 561 Taesik Gong, Jongheon Jeong, Taewon Kim, Yewon Kim, Jinwoo Shin, and Sung-Ju Lee. NOTE:  
562 Robust continual test-time adaptation against temporal correlation. In Alice H. Oh, Alekh Agar-  
563 wal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing*  
564 *Systems*, 2022. URL <https://openreview.net/forum?id=E9HNxrCFZPV>.
- 565 Taesik Gong, Yewon Kim, Taekyung Lee, Sorn Chottanaturak, and Sung-Ju Lee. SoTTA: Robust  
566 test-time adaptation on noisy data streams. In *Thirty-seventh Conference on Neural Information*  
567 *Processing Systems*, 2023.  
568
- 569 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
570 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
571 *(CVPR)*, June 2016.  
572
- 573 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-  
574 ruptions and perturbations. In *International Conference on Learning Representations*, 2019a.  
575 URL <https://openreview.net/forum?id=HJz6tiCqYm>.
- 576 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common cor-  
577 ruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019b.  
578
- 579 Junyuan Hong, Lingjuan Lyu, Jiayu Zhou, and Michael Spranger. Mecta: Memory-economic contin-  
580 ual test-time model adaptation. In *International Conference on Learning Representations*, 2023.  
581 URL <https://openreview.net/pdf?id=N92hjsf5NNh>.
- 582 Ziheng Jiang, Chiyuan Zhang, Kunal Talwar, and Michael C Mozer. Characterizing structural  
583 regularities of labeled data in overparameterized models. In Marina Meila and Tong Zhang  
584 (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of  
585 *Proceedings of Machine Learning Research*, pp. 5034–5044. PMLR, 18–24 Jul 2021. URL  
586 <https://proceedings.mlr.press/v139/jiang21k.html>.
- 587 Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*  
588 *Conference on Learning Representations (ICLR)*, 2015.  
589
- 590 Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In  
591 Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on*  
592 *Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 1885–1894.  
593 PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/koh17a.html>.

- 594 A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. *Master's*  
595 *thesis, Department of Computer Science, University of Toronto, 2009.*  
596
- 597 Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting  
598 out-of-distribution samples and adversarial attacks. *Advances in neural information processing*  
599 *systems*, 31, 2018.
- 600 Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch normaliza-  
601 tion for practical domain adaptation, 2016. URL <https://arxiv.org/abs/1603.04779>.  
602
- 603 Ji Lin, Wei-Ming Chen, Yujun Lin, John Cohn, Chuang Gan, and Song Han. Mxnet: Tiny deep  
604 learning on iot devices. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin  
605 (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 11711–11722. Cur-  
606 ran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/](https://proceedings.neurips.cc/paper_files/paper/2020/file/86c51678350f656dcc7f490a43946ee5-Paper.pdf)  
607 [paper/2020/file/86c51678350f656dcc7f490a43946ee5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/86c51678350f656dcc7f490a43946ee5-Paper.pdf).
- 608 Haolan Liu, Zixuan Wang, and Jishen Zhao. Cola: Characterizing and optimizing the tail latency  
609 for safe level-4 autonomous vehicle systems. *arXiv preprint arXiv:2305.07147*, 2023.  
610
- 611 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *Inter-*  
612 *national Conference on Learning Representations (ICLR)*, 2017.  
613
- 614 TorchVision maintainers and contributors. Torchvision: Pytorch's computer vision library. [https://](https://github.com/pytorch/vision)  
615 [github.com/pytorch/vision](https://github.com/pytorch/vision), 2016.
- 616 Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of  
617 machine learning models. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th*  
618 *International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning*  
619 *Research*, pp. 6950–6960. PMLR, 13–18 Jul 2020. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v119/mirzasoleiman20a.html)  
620 [press/v119/mirzasoleiman20a.html](https://proceedings.mlr.press/v119/mirzasoleiman20a.html).  
621
- 622 Zachary Nado, Shreyas Padhy, D Sculley, Alexander D'Amour, Balaji Lakshminarayanan, and  
623 Jasper Snoek. Evaluating prediction-time batch normalization for robustness under covariate  
624 shift. *arXiv preprint arXiv:2006.10963*, 2020.
- 625 Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofu Chen, Shijian Zheng, Peilin Zhao, and Mingkui  
626 Tan. Efficient test-time model adaptation without forgetting. In Kamalika Chaudhuri, Stefanie  
627 Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th*  
628 *International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning*  
629 *Research*, pp. 16888–16905. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v162/niu22a.html)  
630 [press/v162/niu22a.html](https://proceedings.mlr.press/v162/niu22a.html).
- 631 Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Zhiqian Wen, Yaofu Chen, Peilin Zhao, and Mingkui  
632 Tan. Towards stable test-time adaptation in dynamic wild world. In *The Eleventh International*  
633 *Conference on Learning Representations*, 2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=g2YraF75Tj)  
634 [id=g2YraF75Tj](https://openreview.net/forum?id=g2YraF75Tj).  
635
- 636 NVIDIA Corporation. *NVIDIA Jetson Nano*, 2019. URL [https://developer.nvidia.](https://developer.nvidia.com/embedded/jetson-nano)  
637 [com/embedded/jetson-nano](https://developer.nvidia.com/embedded/jetson-nano). Accessed: 2024-11-20.  
638
- 639 Mansheej Paul, Surya Ganguli, and Gintare Karolina Dziugaite. Deep learning on a data diet:  
640 Finding important examples early in training. In A. Beygelzimer, Y. Dauphin, P. Liang, and  
641 J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, 2021. URL  
642 <https://openreview.net/forum?id=Uj7pF-D-YvT>.
- 643 Geoff Pleiss, Tianyi Zhang, Ethan Elenberg, and Kilian Q Weinberger. Identifying  
644 mislabeled data using the area under the margin ranking. In H. Larochelle,  
645 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural In-*  
646 *formation Processing Systems*, volume 33, pp. 17044–17056. Curran Associates, Inc.,  
647 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf)  
[file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/c6102b3727b2a7d8b1bb6981147081ef-Paper.pdf).

- 648 Omead Pooladzandi, David Davini, and Baharan Mirzasoleiman. Adaptive second order core-  
649 sets for data-efficient machine learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,  
650 Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International  
651 Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,  
652 pp. 17848–17869. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/  
653 v162/pooladzandi22a.html](https://proceedings.mlr.press/v162/pooladzandi22a.html).
- 654 Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence.  
655 *Dataset shift in machine learning*. Mit Press, 2008.
- 656
- 657 Raspberry Pi Foundation. *Raspberry Pi 4 Model B*, 2019. URL [https://www.raspberrypi.  
658 com/products/raspberry-pi-4-model-b/](https://www.raspberrypi.com/products/raspberry-pi-4-model-b/). Accessed: 2024-11-20.
- 659
- 660 Raspberry Pi Foundation. *Raspberry Pi Zero 2 W*, 2021. URL [https://www.raspberrypi.  
661 com/products/raspberry-pi-zero-2-w/](https://www.raspberrypi.com/products/raspberry-pi-zero-2-w/). Accessed: 2024-11-20.
- 662 Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias  
663 Bethge. Improving robustness against common corruptions by covariate shift adaptation.  
664 In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances  
665 in Neural Information Processing Systems*, volume 33, pp. 11539–11551. Curran Associ-  
666 ates, Inc., 2020. URL [https://proceedings.neurips.cc/paper/2020/file/  
667 85690f81aadcl749175c187784afc9ee-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/85690f81aadcl749175c187784afc9ee-Paper.pdf).
- 668 Mattia Segu, Alessio Tonioni, and Federico Tombari. Batch normalization embeddings for deep  
669 domain generalization. *Pattern Recognition*, 135:109115, 2023.
- 670
- 671 Junha Song, Jungsoo Lee, In So Kweon, and Sungha Choi. Ecotta: Memory-efficient continual  
672 test-time adaptation via self-distilled regularization. In *Proceedings of the IEEE/CVF Conference  
673 on Computer Vision and Pattern Recognition*, pp. 11920–11929, 2023.
- 674
- 675 Ardi Tampuu, Kristjan Roosild, and Ilmar Uduste. The effects of speed and delays on test-time  
676 performance of end-to-end self-driving. *Sensors*, 24(6):1963, 2024.
- 677 Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio,  
678 and Geoffrey J. Gordon. An empirical study of example forgetting during deep neural network  
679 learning. In *International Conference on Learning Representations*, 2019. URL [https://  
680 openreview.net/forum?id=BJlxm30cKm](https://openreview.net/forum?id=BJlxm30cKm).
- 681
- 682 Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully  
683 test-time adaptation by entropy minimization. In *International Conference on Learning Repre-  
684 sentations*, 2021. URL <https://openreview.net/forum?id=uX13bZLkr3c>.
- 685
- 686 Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In  
687 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
pp. 7201–7211, June 2022.
- 688
- 689 Xiaobo Xia, Jiale Liu, Jun Yu, Xu Shen, Bo Han, and Tongliang Liu. Moderate coreset: A universal  
690 method of data selection for real-world data-efficient deep learning. In *The Eleventh International  
691 Conference on Learning Representations*, 2022.
- 692
- 693 Shuo Yang, Zeke Xie, Hanyu Peng, Min Xu, Mingming Sun, and Ping Li. Dataset pruning: Reduc-  
694 ing training data by examining generalization influence. In *The Eleventh International Confer-  
695 ence on Learning Representations*, 2023. URL [https://openreview.net/forum?id=  
696 4wZiAXD29TQ](https://openreview.net/forum?id=4wZiAXD29TQ).
- 697
- 698 Longhui Yuan, Binhui Xie, and Shuang Li. Robust test-time adaptation in dynamic scenarios. In  
699 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*,  
pp. 15922–15932, June 2023.
- 700
- 701 Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In  
*Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12,  
2014, Proceedings, Part I 13*, pp. 818–833. Springer, 2014.

702 Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via  
703 adaptation and augmentation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Bel-  
704 grave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing*  
705 *Systems*, volume 35, pp. 38629–38642. Curran Associates, Inc., 2022. URL  
706 [https://proceedings.neurips.cc/paper\\_files/paper/2022/file/  
707 fc28053a08f59fccb48b11f2e31e81c7-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/fc28053a08f59fccb48b11f2e31e81c7-Paper-Conference.pdf).  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## A RELATED WORK

**Test-time adaptation.** Test-time adaptation (TTA) aims to improve model performance on Out-of-Distribution (OOD) data by using only the unlabeled test data stream to adapt the model. Test-time normalization (Nado et al., 2020; Schneider et al., 2020) adjusts the batch normalization (BN) statistics using test data to improve performance. Other works mainly involve updating the parameters of the model during test-time. Tent (Wang et al., 2021) adapts the affine parameters of the BN layers to minimize the entropy of its predictions. EATA (Niu et al., 2022) builds upon Tent, sampling reliable and non-redundant samples and utilizing an anti-forgetting regularizer for efficiency. Other works introduce more complex schemes, primarily to improve robustness against more practical test-time scenarios. CoTTA (Wang et al., 2022) addresses a continually changing test-time environment by using weight-averaged and augmentation-averaged predictions with stochastic restoring. SAR (Niu et al., 2023) filters samples with large and noisy gradients to stabilize the model during wilder test-time scenarios. RoTTA (Yuan et al., 2023) targets a practical test-time setting of changing distributions and correlative sampling by introducing a memory bank and a teacher-student model.

**Test-time adaptation on edge devices.** TTA on edge devices primarily inherit the challenges of on-device learning: limited memory and increased latency from general resource constraints (Lin et al., 2020). Several memory-efficient TTA works have been proposed in this regard. MECTA (Hong et al., 2023) aims to reduce the memory consumption of gradient-based TTA, proposing an adaptive normalization layer to reduce the intermediate caches for backpropagation. Another work EcoTTA (Song et al., 2023) proposes memory-efficient continual TTA by adapting lightweight meta networks instead of the originals to reduce the size of intermediate activations. Despite works to promote memory-efficiency, the latency of TTA, especially on resource-constrained edge devices, has been generally overlooked. While many adaptation-based TTA (Wang et al., 2021; Niu et al., 2022; 2023; Yuan et al., 2023) update only the affine parameters for general time and memory concerns, they still involve computationally-heavy operations every batch, which can lead to high latency on edge devices. A recent work (Alfarra et al., 2024) introduces a more realistic TTA evaluation protocol that penalizes slow TTA methods by providing them with fewer samples for adaptation. We build on from this notion, proposing a sparse TTA setting to reduce the latency of existing TTA methods, but at a minimal cost to performance.

**Data-efficient deep learning.** Data-efficient deep learning methods enable deep learning models to achieve competitive performance with less data. Among these methods, data selection, or data sampling, involves utilizing a small subset of the training data in an attempt to match that of full-dataset training. A branch of data-selection is score-based selection, which scores each sample based on some predefined metric, such as a sample’s influence (Koh & Liang, 2017), difficulty (Toneva et al., 2019; Paul et al., 2021), prediction confidence (Pleiss et al., 2020), or consistency (Jiang et al., 2021), and selects samples with scores in a certain range. Another set of data-selection methods involve optimization-based selection, which formulates an optimization problem to find a optimal subset that can best approximate full-dataset training (Mirzasoileiman et al., 2020; Yang et al., 2023; Pooladzandi et al., 2022). While these approaches work well in their preconceived settings, they generally suffer performance drop as their settings change, such as a change in sampling ratio. More recent works like the Moderate Coreset (Xia et al., 2022) proposes a more robust selection approach by using the distance of a sample to the class center as a score criterion, for an effective representation of the dataset. While our proposed sparse TTA setting is more challenging than the conventional data-efficient setting, as we cannot access ground truths labels nor make assumptions regarding the model, we utilize similar ideas of representative sampling as motivation for our method.

## B EXPERIMENT DETAILS

All experiments presented in this paper were conducted using three random seeds (0, 1, 2), and we report the average accuracies along with their corresponding standard deviations. To ensure efficiency in experimentation, accuracy measurements were obtained using NVIDIA GeForce RTX 3090 GPUs, as the performance differences attributable to the random seed are negligible. Latency measurements were conducted on a Raspberry Pi 4 (Raspberry Pi Foundation, 2019), equipped with a Quad-core Cortex-A72 (ARM v8) 64-bit SoC @ 1.8GHz CPU and 4GB RAM.

## B.1 BASELINE IMPLEMENTATION DETAILS

In this study, we utilized the official implementations of the baseline methods. To ensure consistency, we adopted the reported best hyperparameters documented in the respective papers or source code repositories as much as possible. Also, we present information about the implementation specifics of the baseline methods and provide a comprehensive overview of our experimental setup, including detailed descriptions of the employed hyperparameters.

We adopt hyperparameters from the original papers or the official code of the baselines for consistency. To assess the generality of SNAP-TTA, the test batch sizes were set to 16 for all baseline methods to ensure a fair comparison. To minimize overhead and maintain consistency with inference batches, we set the size of CnDRM equal to the batch size. TTA is conducted in an online manner, with adaptation or inference performed per batch. When there was a conflict between the implementation of SNAP-TTA and certain components of the existing baseline methods, we prioritized SNAP-TTA’s features for fair evaluation at the STTA setting.

For Tent (Wang et al., 2021), we update the BN affine parameters using the SGD optimizer (Loshchilov & Hutter, 2017) with a learning rate of  $l = 1e - 3$  for CIFAR10/100C and  $l = 1e - 4$  for ImageNet-C. For separate experimentation on the ViT, we used a learning rate of  $l = 2e - 4$ . For CoTTA (Wang et al., 2022), we update all model parameters using the Adam optimizer (Kingma & Ba, 2015) with a learning rate of  $l = 1e - 4$ . Furthermore, we set CoTTA’s teacher model EMA factor to  $\alpha = 0.99$ , the restoration factor to  $p = 0.1$ , and the anchor probability to  $p_{th} = 0.9$ . For EATA (Niu et al., 2022), we use the SGD optimizer with a learning rate of  $l = 1e - 4$ . We set the entropy threshold as  $E_0 = 0.4 \times \ln |N|$ , where  $N$  is the total number of classes. For SAR (Niu et al., 2023), we use SAM (Foret et al., 2021) with the base optimizer as SGD with a learning rate of  $l = 1e - 3$ . For fair evaluation, we replaced the sample filtering scheme with SNAP-TTA’s CnDRM. For RoTTA (Yuan et al., 2023), we use the SGD optimizer with a learning rate of  $l = 1e - 3$ . For fair evaluation, we replaced RoTTA’s RBN and CSTU with SNAP-TTA’s CnDRM and IoBMN. For the teacher-student structure, we set the teacher model’s exponential moving average update rate as  $v = 1e - 3$ .

Finally, we list the hyperparameters specific to the components of SNAP-TTA. The confidence threshold for CnDRM  $\tau_{conf}$  is set to 0.4 for CIFAR10-C, 0.45 for CIFAR100-C, and 0.5 for ImageNet-C. The entropy threshold for our ablation study  $\tau_{entr}$  is set to  $\log(10) \times 0.40$  for CIFAR10-C and  $\log(100) \times 0.40$  for CIFAR100-C, as referenced in a previous work using entropy-based filtering (Niu et al., 2022). Additionally, the parameters for the soft shrinkage function in IoBMN are fixed with  $\alpha = 4$  for Tent, CoTTA, SAR, RoTTA, and  $\alpha = 2$  for EATA.

## C DETAILED EXPERIMENT RESULTS

In this section, we provide detailed experimental results for the performance comparison of SNAP-TTA across a wide range of adaptation rates. We evaluated the performance on CIFAR10-C, CIFAR100-C, and ImageNet-C datasets with adaptation rates of 0.01, 0.03, 0.05, 0.1, 0.3, and 0.5, and across five state-of-the-art (SOTA) TTA algorithms: Tent, EATA, SAR, CoTTA, and RoTTA. This comprehensive evaluation resulted in a total of 150 combinations (3 datasets, 6 adaptation rates, 5 algorithms).

The results demonstrate that, regardless of the adaptation rate, dataset, or the TTA algorithm, integrating SNAP-TTA consistently outperforms the baseline methods. Specifically, SNAP-TTA achieved the highest accuracy across nearly all of these 150 combinations, effectively demonstrating its robustness in both high and low adaptation settings. For CIFAR10-C and CIFAR100-C, SNAP-TTA showed substantial performance improvements compared to the baseline, even at very low adaptation rates (e.g., 0.01 and 0.05). Similarly, for ImageNet-C, SNAP-TTA maintained superior accuracy across diverse corruption types.

These results highlight that SNAP-TTA effectively balances adaptation and latency, ensuring optimal performance even when the adaptation rate is sparse and regardless of the underlying TTA algorithm. This consistent superiority across all 150 combinations underscores SNAP-TTA’s suitability for practical, real-world applications on resource-constrained devices.





Table 7: STTA classification accuracy (%) comparing with and without SNAP-TTA on ImageNet-C through Adaptation Rates(AR) (0.05, 0.03, and 0.01). **Bold** numbers are the highest accuracy.

AR	Methods	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.
0.05	Tent	23.77	24.65	24.44	20.54	20.27	32.73	43.57	40.82	35.92	52.78	63.82	15.95	49.33	53.46	47.19	36.62
	+ SNAP-TTA	<b>29.12</b>	<b>30.46</b>	<b>30.30</b>	<b>25.77</b>	<b>25.22</b>	<b>38.21</b>	<b>46.14</b>	<b>44.29</b>	<b>39.95</b>	<b>54.65</b>	<b>65.47</b>	<b>33.81</b>	<b>50.83</b>	<b>55.59</b>	<b>49.21</b>	<b>41.27</b>
	CoTTA	11.03	11.91	11.75	11.03	11.20	22.30	34.98	30.87	29.78	43.99	61.87	12.92	40.26	45.23	36.63	27.72
	+ SNAP-TTA	<b>15.22</b>	<b>15.97</b>	<b>15.93</b>	<b>13.91</b>	<b>14.05</b>	<b>24.87</b>	<b>36.48</b>	<b>32.60</b>	<b>31.65</b>	<b>46.09</b>	<b>63.59</b>	<b>15.67</b>	<b>42.00</b>	<b>46.71</b>	<b>37.96</b>	<b>30.18</b>
	EATA	19.53	20.65	20.72	16.74	16.96	29.11	41.22	37.96	34.84	50.75	63.29	19.86	45.92	51.15	44.13	34.19
	+ SNAP-TTA	<b>22.83</b>	<b>23.95</b>	<b>23.62</b>	<b>19.43</b>	<b>19.70</b>	<b>30.34</b>	<b>41.59</b>	<b>38.06</b>	<b>35.06</b>	<b>50.98</b>	<b>63.30</b>	<b>23.72</b>	<b>46.26</b>	<b>51.52</b>	<b>45.46</b>	<b>35.72</b>
	SAR	23.25	24.23	23.66	19.98	20.38	33.05	43.04	40.73	36.06	52.61	64.09	20.17	49.00	53.35	46.73	36.69
	+ SNAP-TTA	<b>27.54</b>	<b>29.03</b>	<b>28.66</b>	<b>24.05</b>	<b>23.42</b>	<b>36.28</b>	<b>44.12</b>	<b>42.89</b>	<b>38.54</b>	<b>53.24</b>	<b>64.25</b>	<b>31.83</b>	<b>48.79</b>	<b>54.04</b>	<b>47.80</b>	<b>39.63</b>
	RoTTA	14.42	15.22	15.02	13.25	13.31	23.79	35.27	32.09	30.43	44.71	62.64	15.24	40.63	45.55	36.75	29.22
	+ SNAP-TTA	<b>14.65</b>	<b>15.48</b>	<b>15.29</b>	<b>13.43</b>	<b>13.45</b>	<b>23.93</b>	<b>35.33</b>	<b>32.18</b>	<b>30.53</b>	<b>44.71</b>	<b>62.58</b>	<b>15.41</b>	<b>40.64</b>	<b>45.55</b>	<b>36.81</b>	<b>29.33</b>
	Tent	21.76	22.76	22.58	19.06	18.90	30.85	42.34	38.94	35.53	51.58	63.42	18.61	47.96	52.41	45.56	35.48
	+ SNAP-TTA	<b>26.42</b>	<b>28.20</b>	<b>27.81</b>	<b>23.79</b>	<b>22.82</b>	<b>35.77</b>	<b>44.80</b>	<b>42.37</b>	<b>38.81</b>	<b>53.34</b>	<b>64.95</b>	<b>30.05</b>	<b>49.28</b>	<b>54.16</b>	<b>47.57</b>	<b>39.34</b>
CoTTA	10.61	12.36	11.78	11.66	11.32	22.25	35.01	30.88	29.84	44.09	61.83	12.92	40.26	45.20	36.58	27.77	
+ SNAP-TTA	<b>15.29</b>	<b>16.02</b>	<b>16.00</b>	<b>13.99</b>	<b>14.06</b>	<b>24.78</b>	<b>36.54</b>	<b>32.62</b>	<b>31.70</b>	<b>46.01</b>	<b>63.49</b>	<b>15.69</b>	<b>42.05</b>	<b>46.75</b>	<b>37.97</b>	<b>30.20</b>	
EATA	17.17	18.34	17.94	14.48	15.04	26.31	39.47	35.51	33.41	49.16	63.06	18.01	44.16	49.90	42.47	32.30	
+ SNAP-TTA	<b>20.75</b>	<b>21.87</b>	<b>21.28</b>	<b>17.34</b>	<b>17.90</b>	<b>28.08</b>	<b>39.84</b>	<b>36.27</b>	<b>33.54</b>	<b>49.50</b>	<b>63.04</b>	<b>20.86</b>	<b>44.68</b>	<b>49.97</b>	<b>43.53</b>	<b>33.90</b>	
SAR	20.38	21.34	21.18	18.24	18.28	30.56	41.63	38.57	35.23	51.19	63.74	20.40	47.32	52.02	44.81	34.99	
+ SNAP-TTA	<b>25.11</b>	<b>26.27</b>	<b>26.00</b>	<b>22.02</b>	<b>21.25</b>	<b>33.51</b>	<b>42.86</b>	<b>40.83</b>	<b>37.09</b>	<b>51.87</b>	<b>63.83</b>	<b>28.36</b>	<b>47.19</b>	<b>52.63</b>	<b>45.80</b>	<b>37.64</b>	
RoTTA	14.36	15.12	14.95	13.30	13.34	23.78	35.23	31.89	30.33	44.52	62.48	15.20	40.50	45.36	36.63	29.13	
+ SNAP-TTA	<b>14.45</b>	<b>15.21</b>	<b>15.06</b>	<b>13.35</b>	<b>13.42</b>	<b>23.83</b>	<b>35.26</b>	<b>31.92</b>	<b>30.36</b>	<b>44.53</b>	<b>62.47</b>	<b>15.27</b>	<b>40.50</b>	<b>45.39</b>	<b>36.65</b>	<b>29.18</b>	
Tent	17.09	17.70	17.69	14.91	15.25	25.23	38.66	34.15	32.28	48.14	62.65	15.76	43.44	49.14	41.18	31.55	
+ SNAP-TTA	<b>20.66</b>	<b>21.73</b>	<b>21.55</b>	<b>18.46</b>	<b>18.28</b>	<b>29.88</b>	<b>40.63</b>	<b>36.97</b>	<b>34.89</b>	<b>49.85</b>	<b>64.29</b>	<b>22.64</b>	<b>45.13</b>	<b>50.77</b>	<b>43.17</b>	<b>34.59</b>	
CoTTA	11.11	13.24	11.86	10.85	10.97	22.18	34.96	30.88	29.63	44.09	61.71	12.81	40.16	45.14	36.73	27.75	
+ SNAP-TTA	<b>15.09</b>	<b>16.00</b>	<b>15.83</b>	<b>13.84</b>	<b>14.06</b>	<b>24.70</b>	<b>36.47</b>	<b>32.59</b>	<b>31.66</b>	<b>46.10</b>	<b>63.62</b>	<b>15.60</b>	<b>42.03</b>	<b>46.74</b>	<b>38.17</b>	<b>30.17</b>	
EATA	14.85	15.61	15.69	13.26	13.37	23.72	36.18	32.57	31.14	46.06	62.35	13.88	41.91	47.00	38.88	29.76	
+ SNAP-TTA	<b>16.73</b>	<b>17.55</b>	<b>17.30</b>	<b>14.35</b>	<b>14.64</b>	<b>24.13</b>	<b>36.83</b>	<b>32.81</b>	<b>31.09</b>	<b>46.63</b>	<b>62.20</b>	<b>15.26</b>	<b>42.34</b>	<b>47.44</b>	<b>39.81</b>	<b>30.61</b>	
SAR	16.08	17.04	16.69	14.72	14.78	25.92	37.85	34.07	32.25	47.66	63.15	17.20	43.05	48.78	40.14	31.29	
+ SNAP-TTA	<b>18.89</b>	<b>19.45</b>	<b>19.70</b>	<b>16.70</b>	<b>16.55</b>	<b>27.69</b>	<b>38.57</b>	<b>35.34</b>	<b>33.09</b>	<b>48.08</b>	<b>63.04</b>	<b>20.39</b>	<b>42.95</b>	<b>48.76</b>	<b>40.99</b>	<b>32.68</b>	
RoTTA	14.30	15.06	14.89	13.30	13.37	23.78	35.22	31.79	30.27	44.40	62.40	15.16	40.42	45.27	36.54	29.08	
+ SNAP-TTA	<b>14.30</b>	<b>15.07</b>	<b>14.92</b>	<b>13.30</b>	<b>13.38</b>	<b>23.78</b>	<b>35.22</b>	<b>31.78</b>	<b>30.26</b>	<b>44.41</b>	<b>62.40</b>	<b>15.15</b>	<b>40.43</b>	<b>45.27</b>	<b>36.54</b>	<b>29.08</b>	



Table 9: STTA classification accuracy (%) comparing with and without SNAP-TTA on CIFAR10-C through Adaptation Rates(AR) (0.05, 0.03, and 0.01). **Bold** numbers are the highest accuracy.

AR	Methods	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.	
0.05	Tent	64.65	67.08	58.48	85.00	62.61	82.76	84.63	79.01	77.66	83.32	88.00	82.34	74.16	77.11	69.40	75.75	
	+ SNAP-TTA	<b>67.71</b>	<b>69.84</b>	<b>59.53</b>	<b>87.10</b>	<b>64.66</b>	<b>85.73</b>	<b>86.35</b>	<b>80.68</b>	<b>78.92</b>	<b>85.60</b>	<b>90.19</b>	<b>86.72</b>	<b>76.16</b>	<b>78.86</b>	<b>70.95</b>	<b>77.93</b>	
	CoTTA	59.27	61.18	56.33	72.22	57.37	74.27	72.61	70.03	68.68	74.82	79.72	65.57	66.92	64.13	65.25	67.22	
	+ SNAP-TTA	<b>60.66</b>	<b>61.12</b>	<b>60.06</b>	<b>81.43</b>	<b>61.10</b>	<b>81.46</b>	<b>81.11</b>	<b>81.02</b>	<b>80.92</b>	<b>81.09</b>	<b>81.07</b>	<b>81.38</b>	<b>81.14</b>	<b>81.27</b>	<b>80.98</b>	<b>81.05</b>	
	EATA	64.68	67.01	58.07	84.90	62.56	82.64	84.57	78.77	77.16	83.09	87.80	81.62	74.05	76.99	69.31	75.55	
	+ SNAP-TTA	<b>67.36</b>	<b>68.73</b>	<b>59.35</b>	<b>87.05</b>	<b>64.36</b>	<b>85.62</b>	<b>86.48</b>	<b>81.31</b>	<b>78.73</b>	<b>85.33</b>	<b>90.03</b>	<b>86.31</b>	<b>76.04</b>	<b>78.79</b>	<b>70.90</b>	<b>77.76</b>	
	SAR	64.79	66.32	57.58	84.66	62.46	81.42	84.13	78.87	77.20	82.62	88.10	82.12	74.04	75.38	69.13	75.25	
	+ SNAP-TTA	<b>66.00</b>	<b>68.85</b>	<b>58.47</b>	<b>86.54</b>	<b>63.06</b>	<b>85.26</b>	<b>86.13</b>	<b>80.38</b>	<b>78.17</b>	<b>85.17</b>	<b>90.93</b>	<b>85.96</b>	<b>75.27</b>	<b>77.37</b>	<b>70.61</b>	<b>77.21</b>	
	RoTTA	63.21	64.87	56.60	84.64	62.16	82.31	84.13	78.16	76.39	82.90	87.44	81.47	73.59	76.02	68.09	74.80	
	+ SNAP-TTA	<b>65.28</b>	<b>66.91</b>	<b>57.88</b>	<b>86.75</b>	<b>63.51</b>	<b>85.48</b>	<b>86.17</b>	<b>80.46</b>	<b>78.38</b>	<b>85.24</b>	<b>89.99</b>	<b>85.82</b>	<b>75.66</b>	<b>77.98</b>	<b>70.15</b>	<b>77.05</b>	
	0.03	Tent	64.36	66.21	57.65	84.73	62.95	83.07	84.50	78.46	76.99	83.00	88.07	82.62	73.93	76.50	68.82	75.46
		+ SNAP-TTA	<b>66.32</b>	<b>68.38</b>	<b>59.00</b>	<b>86.93</b>	<b>64.04</b>	<b>85.58</b>	<b>86.35</b>	<b>80.78</b>	<b>78.68</b>	<b>85.34</b>	<b>90.08</b>	<b>86.19</b>	<b>75.77</b>	<b>78.37</b>	<b>70.49</b>	<b>77.49</b>
CoTTA		60.38	61.26	56.71	72.44	57.58	74.64	72.73	69.68	68.34	74.64	79.52	67.28	67.42	64.89	66.19	67.58	
+ SNAP-TTA		<b>71.12</b>	<b>73.68</b>	<b>66.34</b>	<b>85.30</b>	<b>66.64</b>	<b>84.25</b>	<b>84.55</b>	<b>80.88</b>	<b>80.11</b>	<b>84.06</b>	<b>89.89</b>	<b>81.98</b>	<b>76.27</b>	<b>79.77</b>	<b>75.35</b>	<b>78.68</b>	
EATA		63.99	65.95	57.39	84.71	62.66	83.11	84.44	78.42	76.63	82.97	88.00	82.55	73.85	76.46	68.91	75.34	
+ SNAP-TTA		<b>66.16</b>	<b>67.60</b>	<b>58.81</b>	<b>86.95</b>	<b>64.06</b>	<b>85.49</b>	<b>86.34</b>	<b>80.79</b>	<b>78.65</b>	<b>85.24</b>	<b>90.09</b>	<b>86.23</b>	<b>75.88</b>	<b>78.48</b>	<b>70.56</b>	<b>77.42</b>	
SAR		63.72	65.75	57.89	84.37	62.45	81.47	82.46	78.32	76.79	81.93	88.60	82.72	73.89	74.55	68.79	74.91	
+ SNAP-TTA		<b>65.40</b>	<b>67.68</b>	<b>58.37</b>	<b>86.72</b>	<b>63.11</b>	<b>85.10</b>	<b>86.18</b>	<b>79.93</b>	<b>78.05</b>	<b>84.92</b>	<b>90.93</b>	<b>85.58</b>	<b>75.30</b>	<b>77.22</b>	<b>69.97</b>	<b>76.96</b>	
RoTTA		63.36	65.10	56.64	84.62	62.41	82.96	84.35	78.10	76.42	82.69	87.90	82.34	73.56	76.09	68.39	75.00	
+ SNAP-TTA		<b>65.27</b>	<b>67.05</b>	<b>58.05</b>	<b>86.79</b>	<b>63.48</b>	<b>85.46</b>	<b>86.25</b>	<b>80.39</b>	<b>78.34</b>	<b>85.19</b>	<b>90.10</b>	<b>85.94</b>	<b>75.67</b>	<b>78.04</b>	<b>69.75</b>	<b>77.05</b>	
0.01		Tent	62.43	64.13	55.85	84.03	62.21	82.47	83.87	77.71	76.55	82.75	87.35	81.83	73.24	75.34	67.73	74.50
		+ SNAP-TTA	<b>65.51</b>	<b>67.26</b>	<b>58.05</b>	<b>86.89</b>	<b>63.53</b>	<b>85.44</b>	<b>85.97</b>	<b>80.58</b>	<b>78.35</b>	<b>85.12</b>	<b>90.09</b>	<b>85.86</b>	<b>75.66</b>	<b>78.38</b>	<b>70.12</b>	<b>77.12</b>
	CoTTA	59.75	59.44	54.47	71.12	57.11	72.47	72.83	66.05	65.14	69.75	75.12	64.31	66.22	62.65	64.76	65.41	
	+ SNAP-TTA	<b>71.79</b>	<b>73.61</b>	<b>65.98</b>	<b>85.34</b>	<b>66.76</b>	<b>84.26</b>	<b>84.93</b>	<b>80.64</b>	<b>80.38</b>	<b>83.94</b>	<b>89.98</b>	<b>82.47</b>	<b>76.48</b>	<b>79.61</b>	<b>75.60</b>	<b>78.79</b>	
	EATA	62.36	63.92	55.73	84.05	62.24	82.38	83.90	77.66	76.48	82.67	87.34	81.82	73.30	75.31	67.76	74.46	
	+ SNAP-TTA	<b>65.49</b>	<b>67.19</b>	<b>57.93</b>	<b>86.92</b>	<b>63.65</b>	<b>85.42</b>	<b>85.97</b>	<b>80.46</b>	<b>78.13</b>	<b>85.07</b>	<b>90.03</b>	<b>85.87</b>	<b>75.69</b>	<b>78.20</b>	<b>70.03</b>	<b>77.07</b>	
	SAR	62.50	64.13	55.65	82.30	62.22	77.21	80.11	77.66	76.75	79.12	89.45	81.97	73.39	69.39	67.83	73.31	
	+ SNAP-TTA	<b>65.06</b>	<b>66.93</b>	<b>57.66</b>	<b>86.76</b>	<b>62.78</b>	<b>85.05</b>	<b>85.94</b>	<b>79.95</b>	<b>77.62</b>	<b>84.65</b>	<b>90.72</b>	<b>85.48</b>	<b>75.34</b>	<b>75.72</b>	<b>69.61</b>	<b>76.62</b>	
	RoTTA	62.25	63.71	55.59	84.05	62.17	82.32	83.86	77.56	76.39	82.64	87.27	81.75	73.21	75.15	67.75	74.38	
	+ SNAP-TTA	<b>65.32</b>	<b>66.94</b>	<b>57.85</b>	<b>86.91</b>	<b>63.44</b>	<b>85.32</b>	<b>85.98</b>	<b>80.49</b>	<b>78.22</b>	<b>85.04</b>	<b>90.01</b>	<b>85.77</b>	<b>75.75</b>	<b>78.15</b>	<b>70.06</b>	<b>77.02</b>	



Table 11: STTA classification accuracy (%) comparing with and without SNAP-TTA on CIFAR100-C through Adaptation Rates(AR) (0.05, 0.03, and 0.01). **Bold** numbers are the highest accuracy.

AR	Methods	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.	
0.05	Tent	40.69	41.55	35.14	62.26	40.26	58.92	61.06	51.21	50.00	55.52	64.05	58.45	50.50	54.68	44.36	51.24	
	+ SNAP-TTA	<b>42.87</b>	<b>44.87</b>	<b>37.60</b>	<b>65.01</b>	<b>42.22</b>	<b>62.22</b>	<b>63.72</b>	<b>54.03</b>	<b>53.68</b>	<b>58.03</b>	<b>67.05</b>	<b>63.08</b>	<b>52.97</b>	<b>57.67</b>	<b>46.94</b>	<b>54.13</b>	
	CoTTA	26.15	26.89	25.26	39.48	28.34	41.41	38.77	32.06	30.84	35.56	41.60	28.52	34.99	33.60	34.54	33.20	
	+ SNAP-TTA	<b>42.02</b>	<b>42.70</b>	<b>37.67</b>	<b>58.30</b>	<b>41.57</b>	<b>57.47</b>	<b>58.02</b>	<b>50.55</b>	<b>51.31</b>	<b>52.34</b>	<b>63.63</b>	<b>51.25</b>	<b>49.76</b>	<b>54.94</b>	<b>47.98</b>	<b>50.63</b>	
	EATA	38.46	39.05	33.47	61.07	38.52	58.16	60.59	49.60	49.18	54.41	63.15	57.06	49.09	52.87	42.49	49.81	
	+ SNAP-TTA	<b>40.49</b>	<b>41.64</b>	<b>34.37</b>	<b>64.28</b>	<b>40.38</b>	<b>61.52</b>	<b>63.17</b>	<b>51.66</b>	<b>52.12</b>	<b>56.50</b>	<b>66.03</b>	<b>62.01</b>	<b>51.76</b>	<b>55.66</b>	<b>44.83</b>	<b>52.43</b>	
	SAR	40.28	41.62	35.35	62.84	40.37	59.51	61.68	51.29	50.66	55.60	64.43	58.49	50.90	54.82	44.64	51.50	
	+ SNAP-TTA	<b>41.76</b>	<b>44.24</b>	<b>36.89</b>	<b>64.34</b>	<b>41.54</b>	<b>62.13</b>	<b>63.39</b>	<b>53.24</b>	<b>52.91</b>	<b>57.54</b>	<b>66.89</b>	<b>62.41</b>	<b>52.70</b>	<b>57.23</b>	<b>46.63</b>	<b>53.59</b>	
	RoTTA	36.38	37.38	31.78	61.44	38.26	58.18	60.19	48.98	48.30	53.50	62.73	56.52	49.37	52.19	41.60	49.12	
	+ SNAP-TTA	<b>37.67</b>	<b>38.66</b>	<b>32.47</b>	<b>63.95</b>	<b>40.18</b>	<b>61.33</b>	<b>62.52</b>	<b>51.47</b>	<b>51.32</b>	<b>55.67</b>	<b>65.89</b>	<b>61.24</b>	<b>51.47</b>	<b>54.52</b>	<b>42.84</b>	<b>51.41</b>	
	0.03	Tent	38.55	39.28	33.77	61.64	39.66	58.83	60.89	49.45	49.51	54.64	63.48	57.29	50.34	53.44	43.28	50.27
		+ SNAP-TTA	<b>41.22</b>	<b>42.20</b>	<b>35.31</b>	<b>64.48</b>	<b>40.82</b>	<b>61.96</b>	<b>63.50</b>	<b>52.84</b>	<b>52.36</b>	<b>57.18</b>	<b>66.50</b>	<b>62.17</b>	<b>52.12</b>	<b>56.48</b>	<b>45.72</b>	<b>52.99</b>
CoTTA		27.11	27.73	25.87	40.25	29.52	42.16	39.60	32.74	32.23	36.60	43.33	29.13	36.45	34.51	35.96	34.21	
+ SNAP-TTA		<b>41.77</b>	<b>42.85</b>	<b>37.50</b>	<b>58.61</b>	<b>41.15</b>	<b>57.65</b>	<b>58.05</b>	<b>50.45</b>	<b>51.34</b>	<b>52.72</b>	<b>63.49</b>	<b>51.63</b>	<b>49.87</b>	<b>55.24</b>	<b>48.14</b>	<b>50.70</b>	
EATA		37.94	38.63	32.00	61.02	39.08	58.52	60.28	48.73	49.15	53.89	63.03	56.64	49.45	52.93	42.11	49.56	
+ SNAP-TTA		<b>39.87</b>	<b>41.12</b>	<b>34.48</b>	<b>64.14</b>	<b>40.27</b>	<b>61.91</b>	<b>63.09</b>	<b>52.37</b>	<b>51.93</b>	<b>56.36</b>	<b>66.02</b>	<b>61.88</b>	<b>51.83</b>	<b>55.60</b>	<b>44.59</b>	<b>52.36</b>	
SAR		38.33	39.19	33.15	61.77	39.78	59.09	61.02	49.67	49.86	54.71	63.59	57.45	50.37	53.67	42.88	50.30	
+ SNAP-TTA		<b>39.84</b>	<b>41.83</b>	<b>34.94</b>	<b>63.70</b>	<b>40.49</b>	<b>61.45</b>	<b>63.17</b>	<b>52.27</b>	<b>51.91</b>	<b>56.69</b>	<b>65.91</b>	<b>61.31</b>	<b>51.68</b>	<b>56.06</b>	<b>44.95</b>	<b>52.41</b>	
RoTTA		36.24	36.94	31.15	60.87	38.28	58.25	59.88	48.43	48.17	53.32	62.73	56.18	49.23	52.12	41.28	48.87	
+ SNAP-TTA		<b>37.85</b>	<b>38.68</b>	<b>32.78</b>	<b>63.97</b>	<b>39.75</b>	<b>61.41</b>	<b>62.57</b>	<b>51.53</b>	<b>51.38</b>	<b>55.68</b>	<b>65.56</b>	<b>61.25</b>	<b>51.53</b>	<b>54.84</b>	<b>42.96</b>	<b>51.45</b>	
0.01		Tent	36.08	36.95	31.31	61.03	38.09	57.63	58.76	48.24	48.65	53.45	62.14	55.07	48.59	51.82	40.68	48.57
		+ SNAP-TTA	<b>38.40</b>	<b>39.40</b>	<b>33.26</b>	<b>63.85</b>	<b>40.36</b>	<b>61.23</b>	<b>62.79</b>	<b>51.92</b>	<b>51.73</b>	<b>56.20</b>	<b>65.83</b>	<b>60.95</b>	<b>51.82</b>	<b>54.75</b>	<b>43.53</b>	<b>51.73</b>
	CoTTA	26.59	27.92	24.86	41.34	28.91	43.09	40.11	34.33	33.32	37.99	44.78	28.80	36.26	34.70	35.67	34.58	
	+ SNAP-TTA	<b>42.05</b>	<b>42.91</b>	<b>37.50</b>	<b>58.70</b>	<b>41.22</b>	<b>57.38</b>	<b>58.14</b>	<b>50.39</b>	<b>51.13</b>	<b>52.23</b>	<b>63.42</b>	<b>51.74</b>	<b>49.87</b>	<b>54.84</b>	<b>47.72</b>	<b>50.62</b>	
	EATA	36.10	37.05	31.03	60.86	37.83	57.64	58.77	48.02	48.75	53.37	62.18	54.95	48.55	51.89	40.75	48.51	
	+ SNAP-TTA	<b>38.54</b>	<b>39.78</b>	<b>33.11</b>	<b>63.82</b>	<b>39.98</b>	<b>61.33</b>	<b>62.53</b>	<b>51.76</b>	<b>51.50</b>	<b>56.03</b>	<b>65.94</b>	<b>61.16</b>	<b>51.47</b>	<b>54.52</b>	<b>43.67</b>	<b>51.68</b>	
	SAR	36.04	37.02	31.38	61.13	38.07	58.00	59.08	48.44	48.84	53.52	62.57	55.19	48.87	52.01	40.71	48.72	
	+ SNAP-TTA	<b>37.91</b>	<b>38.85</b>	<b>32.92</b>	<b>63.17</b>	<b>39.35</b>	<b>60.51</b>	<b>62.01</b>	<b>51.11</b>	<b>50.48</b>	<b>55.47</b>	<b>65.07</b>	<b>59.69</b>	<b>51.24</b>	<b>54.10</b>	<b>42.80</b>	<b>50.98</b>	
	RoTTA	35.55	36.34	30.55	60.76	37.42	57.50	58.57	47.87	48.31	53.11	61.90	54.70	48.25	51.37	40.29	48.16	
	+ SNAP-TTA	<b>37.82</b>	<b>38.72</b>	<b>32.60</b>	<b>63.53</b>	<b>39.80</b>	<b>61.00</b>	<b>62.27</b>	<b>51.42</b>	<b>51.33</b>	<b>55.71</b>	<b>65.64</b>	<b>60.89</b>	<b>51.50</b>	<b>54.27</b>	<b>42.92</b>	<b>51.30</b>	

## D ADDITIONAL RESULTS ON ABLATION STUDY

In this section, we provide additional details on the ablation study to evaluate the contributions of the CnDRM and IoBMN components in SNAP-TTA. Specifically, we measured the average accuracy across 15 corruption types on CIFAR10-C and CIFAR100-C datasets under varying adaptation rates (0.3, 0.1, 0.05) to thoroughly assess the effectiveness of each component.

Tables 12, 13, 14, 15, and 16 summarize the results for different combinations of CnDRM and IoBMN across these adaptation rates. The results indicate that the combination of CnDRM (Class and Domain Representative sampling) and IoBMN (inference using memory statistics corrected to match the test batch) consistently yields the highest accuracy. This trend is observed across all evaluated adaptation rates, suggesting that both components contribute significantly to enhancing adaptation performance.

Moreover, individual evaluations show that each component has a distinct positive effect, as evidenced by consistently higher accuracy compared to using no adaptation or only a single component. This emphasizes the complementary nature of CnDRM and IoBMN, which together provide

robust adaptation capabilities for domain-shifted scenarios. These tables provide further insight into the benefits of each configuration and how the synergy of CnDRM and IoBMN results in improved robustness against various corruptions.

Table 12: STTA classification accuracy (%) of ablative settings on the CIFAR10-C, adaptation rate 0.5. Averaged over all 15 corruptions. **Bold** numbers are the highest accuracy.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
naïve	78.86	69.75	79.02	77.83	75.39
Random	78.90	66.04	78.97	77.77	75.06
LowEntropy	78.68	63.74	78.42	76.21	72.83
CRM	80.32	66.50	80.14	75.78	75.49
CnDRM	79.62	77.68	79.63	78.22	75.85
CnDRM+EMA	80.96	72.42	80.27	78.19	76.73
<b>CnDRM+IoDMN</b>	<b>81.23</b>	<b>78.75</b>	<b>81.30</b>	<b>79.77</b>	<b>77.41</b>

Table 13: STTA classification accuracy (%) of ablative settings on the CIFAR10-C, adaptation rate 0.05. Averaged over all 15 corruptions. **Bold** numbers are the highest accuracy.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
naïve	75.75	67.22	75.55	75.25	74.80
Random	75.82	65.90	75.56	75.27	74.91
LowEntropy	74.07	64.08	73.73	73.58	72.83
CRM	76.55	66.14	76.06	74.02	75.23
CnDRM	76.53	77.67	76.29	76.18	75.61
CnDRM+EMA	76.86	71.69	75.98	75.43	75.95
<b>CnDRM+IoDMN</b>	<b>77.93</b>	<b>78.73</b>	<b>77.76</b>	<b>77.21</b>	<b>77.05</b>

Table 14: STTA classification accuracy (%) of ablative settings on the CIFAR100-C, adaptation rate 0.3. Averaged over all 15 corruptions. **Bold** numbers are the highest accuracy.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
naïve	53.36	39.11	49.97	56.65	49.84
Random	53.00	33.49	49.24	56.06	49.00
LowEntropy	53.53	32.29	45.51	55.84	44.77
CRM	54.21	32.86	47.42	56.40	46.68
CnDRM	55.15	50.02	51.36	57.72	50.74
CnDRM+EMA	55.39	41.34	50.11	57.68	49.88
<b>CnDRM+IoDMN</b>	<b>57.27</b>	<b>50.32</b>	<b>52.19</b>	<b>58.44</b>	<b>51.55</b>

Table 15: STTA classification accuracy (%) of ablative settings on the CIFAR100-C, adaptation rate 0.1. Averaged over all 15 corruptions. **Bold** numbers are the highest accuracy.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
naïve	52.84	35.86	49.70	53.49	49.11
Random	52.68	33.18	49.39	53.42	48.84
LowEntropy	51.76	32.30	46.03	52.15	45.18
CRM	52.43	32.54	47.68	53.12	47.01
CnDRM	54.46	50.06	51.41	55.24	50.47
CnDRM+EMA	54.36	41.63	50.21	54.84	49.95
<b>CnDRM+IoDMN</b>	<b>55.84</b>	<b>50.52</b>	<b>52.35</b>	<b>55.76</b>	<b>51.33</b>

Table 16: STTA classification accuracy (%) of ablative settings on the CIFAR100-C, adaptation rate 0.05. Averaged over all 15 corruptions. **Bold** numbers are the highest accuracy.

Methods	Tent	CoTTA	EATA	SAR	RoTTA
naïve	51.24	33.20	49.81	51.50	49.12
Random	51.35	33.71	49.57	51.48	48.98
LowEntropy	49.79	32.36	46.65	49.51	45.41
CRM	50.17	32.74	47.47	50.49	46.58
CnDRM	52.86	50.08	51.47	53.09	50.44
CnDRM+EMA	52.68	41.43	50.32	52.80	50.04
<b>CnDRM+IoDMN</b>	<b>54.13</b>	<b>50.63</b>	<b>52.43</b>	<b>53.59</b>	<b>51.41</b>

## E ADDITIONAL ABLATE ANALYSIS

### E.1 DOMAIN INFLUENCE IN EARLY LAYER REPRESENTATIONS

In deep learning models, early layers capture low-level features such as textures, edges, and frequency components (Zeiler & Fergus, 2014). These features are inherently domain-specific, making these layers more sensitive to shifts in input data distribution—a critical challenge for tasks requiring domain adaptation and generalization (Lee et al., 2018; Segu et al., 2023). This sensitivity arises because early layers encapsulate domain-specific patterns that may not generalize to new distributions. Under the covariate shift assumption (Quiñonero-Candela et al., 2008), while input distributions differ between source and target domains, the conditional distribution of labels remains the same. This discrepancy between input distributions makes early layers particularly vulnerable to domain shifts.

Visualizing early layer feature embeddings using 2D PCA on CIFAR-10C domains reveals distinct domain-specific patterns, highlighting the significant influence of domain information in these representations (Figure 6). Our preliminary experiments further confirm that sparse TTA, using the Wasserstein distance between moving batch normalization statistics and instance-specific statistics derived from early layer hidden features, can significantly improve performance. Selecting instances closer to the target domain distribution center using this distance metric yields better adaptation results, as demonstrated by performance comparisons between the top 20% and bottom 20% of samples (Figure 3). These findings emphasize the crucial role of domain-sensitive early layers in achieving effective adaptation.

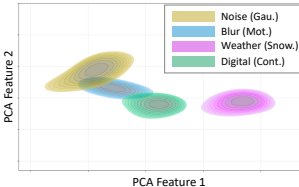


Figure 6: PCA embedding of early layer features for one domain from each of the four main CIFAR10-C corruption categories, showing clear separation between domains.

### E.2 ANALYSIS ON CONFIDENCE THRESHOLD ON PSEUDO-LABEL ACCURACY

We analyzed the impact of using a confidence threshold for pseudo-label selection by comparing random sampling with high-confidence sampling across three benchmarks: CIFAR10-C, CIFAR100-C, and ImageNet-C. Table 17 shows that high-confidence sampling consistently outperformed random sampling, achieving significantly higher pseudo-label accuracy in all datasets. This result demonstrates the effectiveness of selecting high-confidence samples to improve the quality of pseudo-labels, thereby enhancing model adaptation under domain shift conditions.

Table 17: Pseudo-label accuracy comparison between random and high-confidence sampling on three benchmarks: CIFAR10-C, CIFAR100-C, and ImageNet-C. **Bold** numbers are the highest accuracy.

	CIFAR10-C	CIFAR100-C	ImageNet-C
Random	69.91	45.30	23.90
<b>HighConf</b>	<b>74.80</b>	<b>59.38</b>	<b>59.40</b>

### E.3 LATENCY TRACKING OF SNAP-TTA ON DIVERSE EDGE-DEVICES

To evaluate the latency efficiency of SNAP-TTA on resource-constrained edge devices, we measured the adaptation latency across three devices: NVIDIA Jetson Nano (NVIDIA Corporation, 2019), Raspberry Pi 4 (Raspberry Pi Foundation, 2019), and Raspberry Pi Zero 2 W (Raspberry Pi Foundation, 2021). These experiments compared the latency of SNAP-TTA with the Original TTA framework, specifically focusing on five state-of-the-art TTA algorithms: Tent (Wang et al., 2021), EATA (Niu et al., 2022), SAR (Niu et al., 2023), RoTTA (Yuan et al., 2023), and CoTTA (Wang et al., 2022). The experiments were conducted at an adaptation rate of 0.1, demonstrating the effectiveness of SNAP-TTA in reducing adaptation latency while maintaining competitive accuracy.



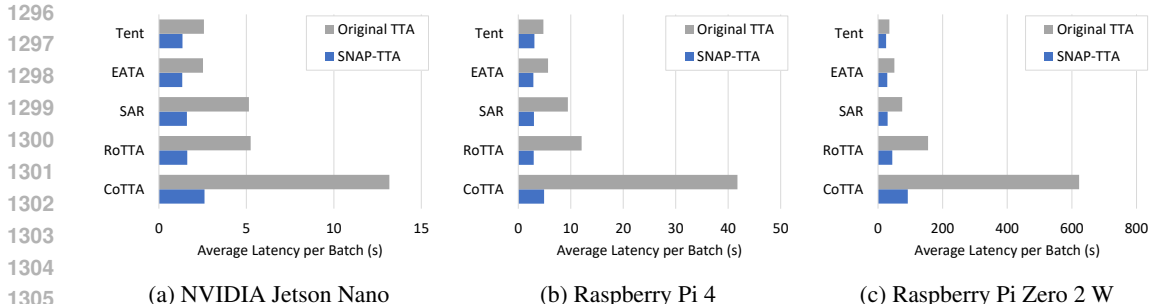


Figure 7: Latency comparison between SNAP-TTA and Original TTA across five state-of-the-art TTA algorithms (Tent, EATA, SAR, RoTTA, CoTTA) on three edge devices: (a) NVIDIA Jetson Nano, (b) Raspberry Pi 4, and (c) Raspberry Pi Zero 2 W. SNAP-TTA demonstrates significant latency reductions while maintaining competitive adaptation performance. The experiments were conducted at an adaptation rate of 0.1.

Figure 7 illustrates the latency performance for each device. It is evident that SNAP-TTA achieves a significant reduction in adaptation latency compared to the Original TTA framework. Notably, the latency reduction was proportional to the adaptation rate, validating the efficiency of SNAP-TTA in sparse adaptation scenarios. For instance, the latency for CoTTA was reduced by up to 87.5% on the Raspberry Pi 4, emphasizing the practical benefits of SNAP-TTA in latency-sensitive environments. Additionally, similar trends were observed across other devices, including the resource-limited Raspberry Pi Zero 2 W.

The results confirm that SNAP-TTA not only ensures substantial latency reductions but also adapts effectively to real-world conditions on diverse edge devices, proving its suitability for deployment in latency-sensitive applications.

#### E.4 MEMORY OVERHEAD OF SNAP-TTA

The SNAP-TTA framework achieves substantial latency reduction and accuracy improvements with minimal memory overhead, even under resource-constrained scenarios like edge devices. In this section, we present both a theoretical analysis of the memory requirements and empirical results obtained from evaluations on a Raspberry Pi 4 (Raspberry Pi Foundation, 2019) (CPU-only edge device).

The memory overhead of SNAP-TTA arises from two main components: (1) the memory buffer in Class and Domain Representative Memory (CnDRM) for storing representative samples, including both feature statistics (mean and variance) and the raw image samples, and (2) the statistics required for Inference-only Batch-aware Memory Normalization (IoBMN). For a batch size  $B$ , the total theoretical memory overhead can be expressed as:  $\text{Memory Overhead} = B \times (\text{Image Size} + 2 \times \text{Feature Dimension} \times \text{Bytes per Value}) + \text{Feature Dimension} \times \text{Bytes per Value} \times 2$ . The last term accounts for the storage of IoBMN statistics (mean and variance for each feature channel). The image size is calculated based on the dataset resolution and data type.

For ResNet18 on CIFAR10-C, CIFAR10 images have a resolution of  $32 \times 32 \times 3$  with each value stored as 1 byte. For a feature dimension of 512 and batch size  $B = 16$ , the total overhead is: Image Overhead =  $16 \times (32 \times 32 \times 3 \times 1) = 49,152$  bytes (48 KB), Feature Overhead (CnDRM) =  $16 \times (512 \times 2 \times 4) = 65,536$  bytes (64 KB), Feature Overhead (IoBMN) =  $512 \times 2 \times 4 = 4,096$  bytes (4 KB). Thus, the total memory overhead is: Total Overhead =  $48 \text{ KB} + 64 \text{ KB} + 4 \text{ KB} = 116 \text{ KB}$ .

For ResNet50 on ImageNet-C, ImageNet images have a resolution of  $224 \times 224 \times 3$ , stored as 1 byte per value. For a feature dimension of 2048 and batch size  $B = 16$ , the total overhead is: Image Overhead =  $16 \times (224 \times 224 \times 3 \times 1) = 12,044,928$  bytes (11.5 MB), Feature Overhead (CnDRM) =  $16 \times (2048 \times 2 \times 4) = 262,144$  bytes (256 KB), Feature Overhead (IoBMN) =  $2048 \times 2 \times 4 = 16,384$  bytes (16 KB). Thus, the total memory overhead is: Total Overhead =  $11.5 \text{ MB} + 256 \text{ KB} + 16 \text{ KB} \approx 11.77 \text{ MB}$ .

Table 18 shows the empirical memory usage of SNAP-TTA compared to Original TTA methods (Tent, EATA, CoTTA, SAR, and RoTTA). The results were averaged across three seeds of experiments and represent the memory footprint observed in a CPU-only edge device, Raspberry Pi 4. While minor variations in measurements are expected due to the nature of CPU memory footprint tracking, the results robustly indicate that the actual memory overhead of SNAP-TTA on edge devices is extremely low across all algorithms, ranging from 0.02% to 1.74%. Furthermore, while peak memory usage is either slightly increased or remains comparable to Original TTA methods, the average memory usage of SNAP-TTA is consistently lower. This is because SNAP-TTA performs backpropagation infrequently, which is the most memory-intensive operation in TTA.

Table 18: Comparison of memory usage (Average Memory, Peak Memory, and Memory Overhead) between Original TTA and SNAP-TTA (adaptation rate 0.3) across various methods (Tent, EATA, CoTTA, SAR, and RoTTA) tested on Raspberry Pi 4. **Bold** numbers are the lowest memory usage.

Methods	Average Mem (MB)		Peak Mem (MB)		Mem Overhead (MB)
	Original TTA	SNAP-TTA	Original TTA	SNAP-TTA	SNAP - Original
Tent	764.24	<b>751.35</b>	<b>822.93</b>	828.46	5.52 (0.67%)
CoTTA	1133.52	<b>1099.64</b>	<b>1211.21</b>	1227.99	16.78 (1.13%)
EATA	816.69	<b>749.95</b>	<b>847.73</b>	862.51	14.78 (1.74%)
SAR	786.65	<b>753.69</b>	<b>863.77</b>	865.18	1.41 (0.02%)
RoTTA	933.23	<b>871.64</b>	<b>972.23</b>	983.94	11.71 (1.20%)

These findings demonstrate that **SNAP-TTA’s memory overhead is negligible compared to its benefits in latency reduction and accuracy improvements**. By leveraging a small memory buffer for representative samples and minimizing backpropagation operations, SNAP-TTA not only achieves a lightweight memory profile but also becomes more efficient in terms of average memory usage compared to Original TTA. This lightweight design, combined with its advantages in latency and accuracy, underscores the practicality of SNAP-TTA for deployment in latency-sensitive applications on edge devices.

#### E.5 INTEGRATION OF SNAP-TTA WITH MEMORY-EFFICIENT TTA ALGORITHM: MECTA (HONG ET AL., 2023)

This section evaluates the integration of SNAP-TTA with MECTA, a memory-efficient TTA algorithm, to demonstrate its applicability for resource-constrained edge devices. The experimental setup follows the evaluation settings presented in the MECTA paper to ensure a fair and consistent comparison. Specifically, we analyze the performance of Tent and EATA, enhanced with MECTA and further integrated with SNAP-TTA, using the ResNet50 model with a batch size of 64 on the ImageNet-C dataset.

Table 19 presents the classification accuracy and peak memory usage for Tent+MECTA and EATA+MECTA configurations with and without SNAP-TTA. Integrating SNAP-TTA with Tent+MECTA improves accuracy from 35.21% to 39.52%, while reducing peak memory usage by approximately 30% compared to the Tent baseline. Similarly, SNAP-TTA boosts the accuracy of EATA+MECTA from 35.55% to 42.86% while maintaining an efficient memory footprint.

Table 19: Comparison of classification (%) and memory peak (MB) in STTA with an adaptation rate of 0.1. MECTA significantly reduces memory consumption, and SNAP-TTA is applied alongside it to boost the performance of sparse adaptation. The accuracy is the average over 15 corruptions in ImageNet-C. **Bold** numbers indicate either the lowest memory usage or the highest accuracy.

Methods	Accuracy (%)	Max Memory (MB)
Tent	35.21	6805.26
+MECTA	37.62	<b>4620.25 (-32.10%)</b>
<b>+ MECTA + SNAP-TTA</b>	<b>39.52</b>	4622.12 (-32.08%)
EATA	35.55	6541.02
+MECTA	41.41	<b>4512.38 (-31.01%)</b>
<b>+ MECTA + SNAP-TTA</b>	<b>42.86</b>	4535.44 (-30.66%)

Further details are provided in Table 20, which evaluates the combination of SNAP-TTA with MECTA across various corruption types and adaptation rates (AR = 0.3, 0.1, and 0.05). These results show that SNAP-TTA consistently outperforms baseline configurations across all adaptation rates and corruption types. This demonstrates the robustness of SNAP-TTA when integrated with MECTA and its suitability for real-world applications.

By adhering to the evaluation settings of the MECTA paper, this study ensures high reliability and comparability of results. The findings confirm that SNAP-TTA is highly compatible with MECTA, significantly improving both accuracy and memory efficiency. This synergy highlights the potential of combining SNAP-TTA and MECTA for deployment in resource-constrained environments such as edge devices.

Table 20: Evaluation of SNAP-TTA with MECTA on ImageNet-C through Adaptation Rates(AR) (0.3, 0.1, and 0.05). **Bold** numbers are the highest accuracy.

AR	Methods	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.
0.3	Tent + MECTA	28.20	30.13	29.58	23.07	23.35	34.49	45.95	40.97	35.68	55.66	66.56	14.72	53.09	57.16	50.74	39.29
	+ SNAP-TTA	<b>30.49</b>	<b>31.98</b>	<b>31.66</b>	<b>26.29</b>	<b>26.19</b>	<b>38.47</b>	<b>47.38</b>	<b>43.79</b>	<b>40.12</b>	<b>56.38</b>	<b>66.81</b>	<b>28.87</b>	<b>53.53</b>	<b>57.61</b>	<b>50.86</b>	<b>42.03</b>
	EATA + MECTA	32.18	34.85	33.06	28.80	29.18	41.02	49.24	47.10	41.56	57.35	66.27	34.56	55.38	58.19	52.87	44.11
	+ SNAP-TTA	<b>33.67</b>	<b>35.76</b>	<b>34.86</b>	<b>30.35</b>	<b>30.29</b>	<b>42.78</b>	<b>49.55</b>	<b>47.46</b>	<b>42.32</b>	<b>57.50</b>	66.18	<b>39.08</b>	<b>55.38</b>	<b>58.35</b>	52.72	<b>45.08</b>
		$\pm 0.19$	$\pm 0.24$	$\pm 0.10$	$\pm 0.11$	$\pm 0.04$	$\pm 0.06$	$\pm 0.10$	$\pm 0.10$	$\pm 0.05$	$\pm 0.15$	$\pm 0.06$	$\pm 0.81$	$\pm 0.16$	$\pm 0.12$	$\pm 0.02$	$\pm 0.15$
0.1	Tent + MECTA	24.94	26.73	25.63	21.11	21.46	32.11	44.05	38.22	36.36	53.92	66.48	18.50	50.80	55.67	48.33	37.62
	+ SNAP-TTA	<b>27.49</b>	<b>28.90</b>	<b>28.26</b>	<b>23.49</b>	<b>23.76</b>	<b>34.92</b>	<b>45.18</b>	<b>40.21</b>	<b>38.40</b>	53.78	<b>66.54</b>	<b>27.72</b>	<b>51.00</b>	55.48	47.61	<b>39.52</b>
	EATA + MECTA	29.42	31.72	29.44	24.41	25.48	37.04	47.10	43.60	39.43	55.95	66.42	28.85	53.70	57.34	51.20	41.41
	+ SNAP-TTA	<b>31.26</b>	<b>32.71</b>	<b>32.22</b>	<b>27.31</b>	<b>27.61</b>	<b>38.88</b>	<b>47.83</b>	<b>44.52</b>	<b>40.58</b>	<b>56.42</b>	66.24	<b>35.38</b>	53.67	<b>57.39</b>	50.83	<b>42.86</b>
		$\pm 0.11$	$\pm 0.17$	$\pm 0.17$	$\pm 0.46$	$\pm 0.28$	$\pm 0.28$	$\pm 0.09$	$\pm 0.14$	$\pm 0.05$	$\pm 0.06$	$\pm 0.21$	$\pm 0.63$	$\pm 0.17$	$\pm 0.13$	$\pm 0.12$	$\pm 0.20$
0.05	Tent + MECTA	21.22	23.19	21.90	18.69	19.39	29.89	42.02	36.53	35.23	51.75	66.23	19.64	48.43	53.54	45.43	35.54
	+ SNAP-TTA	<b>23.93</b>	<b>25.37</b>	<b>24.10</b>	<b>20.42</b>	<b>21.14</b>	<b>31.83</b>	<b>42.68</b>	<b>37.53</b>	<b>36.31</b>	51.42	66.19	<b>23.84</b>	<b>48.62</b>	53.20	44.57	<b>36.74</b>
	EATA + MECTA	24.97	26.95	21.87	21.19	21.94	33.61	45.11	40.92	37.73	54.64	66.60	23.03	51.87	56.60	49.15	38.41
	+ SNAP-TTA	<b>28.39</b>	<b>30.10</b>	<b>29.45</b>	<b>24.32</b>	<b>25.12</b>	<b>35.54</b>	<b>46.04</b>	<b>41.87</b>	<b>39.16</b>	<b>55.12</b>	<b>66.61</b>	<b>30.34</b>	<b>52.06</b>	56.42	49.11	<b>40.64</b>
		$\pm 0.57$	$\pm 0.38$	$\pm 0.22$	$\pm 0.20$	$\pm 0.07$	$\pm 0.20$	$\pm 0.27$	$\pm 0.07$	$\pm 0.15$	$\pm 0.01$	$\pm 0.09$	$\pm 0.34$	$\pm 0.24$	$\pm 0.11$	$\pm 0.07$	$\pm 0.20$

## F ADDITIONAL DISCUSSIONS

### F.1 EFFICIENT STRATEGY FOR RE-CALCULATION OF SAMPLE'S DISTANCE

The domain centroid in our framework is updated using a momentum-based approach to effectively capture recent shifts in the target domain. This ensures that the centroid remains adaptive to evolving distributions without being overly influenced by temporary fluctuations. However, during sparse adaptation (SA), where model updates occur at extended intervals, the data distribution can shift substantially between updates. Consequently, distances calculated for older samples may become outdated, leading to inconsistencies when comparing them to more recently added samples that are evaluated based on the updated centroid.

To address this issue efficiently, our Class and Domain Representative Memory (CnDRM) recalculates the distance of samples only when the shift in the domain centroid exceeds a predefined significance threshold. Specifically, if the change in the domain centroid  $\Delta c_{\text{domain}}$  surpasses a threshold  $\tau_{\Delta}$ , the distances of all samples in memory are updated to reflect the new domain conditions. This threshold-based approach ensures that recalculations occur only when necessary, thereby minimizing computational costs while maintaining the representativeness of the memory.

In practice, we observed that the performance was not significantly affected as long as the threshold  $\tau_{\Delta}$  was not set too high, indicating robustness to the choice of threshold. Based on these observations, we set  $\tau_{\Delta} = 0.1$  and used this value consistently for all evaluations. By focusing recalculations on significant shifts, this strategy preserves consistency in sample selection, ensuring that both older and newer samples are compared fairly in the context of the current domain characteristics without excessive computational overhead.

## F.2 STRATEGY FOR CONTINUOUS DOMAIN SHIFT SETTING

In our proposed framework, the centroid used for selecting domain-representative samples naturally adapts to changes in the domain as new data is encountered. This mechanism inherently ensures that the centroid evolves to reflect the characteristics of the current domain, allowing for effective performance even under continual Test-Time Adaptation (TTA) scenarios, where the domain may gradually or abruptly shift during adaptation.

Instead of employing additional mechanisms like z-score evaluation to detect domain shifts, we rely on the natural adaptability of the centroid to adjust to the incoming data. This simplifies the design and avoids unnecessary overhead while maintaining robustness. As the domain characteristics evolve, the centroid continuously aligns with the new domain without requiring explicit detection of changes or manual intervention.

To validate the effectiveness of SNAP-TTA under continual domain shift scenarios, we conducted experiments across various benchmark datasets with incremental and abrupt domain shifts. Table 21 summarizes the results, demonstrating that SNAP-TTA maintains strong performance across evolving domains without requiring additional computational overhead for explicit domain shift detection.

Table 21: Performance of SNAP-TTA under continual domain shift scenarios. The table reports the accuracy (%) for different datasets with incremental and abrupt shifts. **Bold** numbers are the highest accuracy.

AR	Method	Gau.	Shot	Imp.	Def.	Gla.	Mot.	Zoom	Snow	Fro.	Fog	Brit.	Cont.	Elas.	Pix.	JPEG	Avg.
0.1	Tent	24.68	19.65	5.12	0.63	0.43	0.40	0.44	0.41	0.30	0.33	0.42	0.24	0.32	0.31	0.31	3.60
	+ SNAP-TTA	<b>28.71</b>	<b>30.60</b>	<b>22.91</b>	<b>6.13</b>	<b>1.62</b>	<b>0.87</b>	<b>0.88</b>	<b>0.64</b>	<b>0.64</b>	<b>0.66</b>	<b>0.75</b>	<b>0.44</b>	<b>0.60</b>	<b>0.63</b>	<b>0.61</b>	<b>6.45</b>
	CoTTA	10.99	12.21	11.54	11.28	11.13	22.08	34.80	30.69	29.45	43.87	61.92	12.76	40.03	44.99	36.43	27.61
	+ SNAP-TTA	<b>15.19</b>	<b>15.97</b>	<b>15.91</b>	<b>13.94</b>	<b>14.18</b>	<b>24.76</b>	<b>36.50</b>	<b>32.61</b>	<b>31.76</b>	<b>46.14</b>	<b>63.60</b>	<b>15.60</b>	<b>42.17</b>	<b>46.77</b>	<b>38.08</b>	<b>30.21</b>
		$\pm 0.17$	$\pm 0.11$	$\pm 0.02$	$\pm 0.04$	$\pm 0.03$	$\pm 0.07$	$\pm 0.23$	$\pm 0.04$	$\pm 0.06$	$\pm 0.10$	$\pm 0.14$	$\pm 0.04$	$\pm 0.02$	$\pm 0.06$	$\pm 0.12$	$\pm 0.08$
0.05	Tent	23.31	27.08	22.71	9.72	4.14	2.03	1.16	0.66	0.45	0.47	0.61	0.33	0.47	0.47	0.46	6.27
	+ SNAP-TTA	<b>27.10</b>	<b>33.41</b>	<b>31.78</b>	<b>19.85</b>	<b>16.94</b>	<b>14.75</b>	<b>12.46</b>	<b>5.53</b>	<b>2.69</b>	<b>1.47</b>	<b>1.52</b>	<b>0.67</b>	<b>0.88</b>	<b>0.89</b>	<b>0.84</b>	<b>11.39</b>
	CoTTA	11.04	12.25	11.73	11.62	11.25	22.05	34.89	30.73	29.50	44.09	61.87	12.87	40.15	45.06	36.53	27.71
	+ SNAP-TTA	<b>15.20</b>	<b>15.89</b>	<b>15.93</b>	<b>13.81</b>	<b>14.15</b>	<b>24.74</b>	<b>36.68</b>	<b>32.51</b>	<b>31.71</b>	<b>46.11</b>	<b>63.48</b>	<b>15.73</b>	<b>42.20</b>	<b>46.69</b>	<b>38.05</b>	<b>30.19</b>
		$\pm 0.15$	$\pm 0.02$	$\pm 0.10$	$\pm 0.04$	$\pm 0.03$	$\pm 0.16$	$\pm 0.27$	$\pm 0.04$	$\pm 0.20$	$\pm 0.05$	$\pm 0.09$	$\pm 0.19$	$\pm 0.12$	$\pm 0.10$	$\pm 0.04$	$\pm 0.10$

These results indicate that SNAP-TTA effectively handles both incremental and abrupt domain shifts, consistently outperforming baseline methods. By leveraging the natural adaptability of the centroid, SNAP-TTA provides a robust solution for continual domain adaptation in real-world scenarios. Notably, SNAP-TTA mitigates catastrophic forgetting not only through its sparse adaptation strategy but also by leveraging domain centroid-based sampling, allowing performance to be sustained longer in continual shift scenarios. Unlike Tent, CoTTA is specifically designed for continual domain shift environments, which highlights its superior performance under such conditions.

Future work could explore augmenting this adaptive mechanism by incorporating techniques like z-score evaluation to enable even more responsive adjustments. For instance, a z-score-based approach could further refine the centroid’s responsiveness to subtle, gradual domain shifts by monitoring discrepancies between incoming data statistics and the current centroid. Such enhancements could make the system even more effective at handling continual domain evolution, particularly in scenarios with complex or noisy data streams.

## F.3 MODIFICATION FOR LAYER NORMALIZATION OF ViT

The main text describes the use of Batch Normalization (BN) statistics for calculating domain centroids and centroid-instance distances, with subsequent adjustment of memory statistics to match the target test batch using the Inference-only Batch-aware Memory Normalization (IoBMN) method. Specifically, these calculations leverage the mean and variance across batches as follows:

$$\bar{\mu}_c = \frac{1}{B \times L} \sum_{b=1}^B \sum_{l=1}^L f_{b,c,l}, \quad \bar{\sigma}_c^2 = \frac{1}{B \times L} \sum_{b=1}^B \sum_{l=1}^L (f_{b,c,l} - \mu_{b,c})^2, \quad (6)$$

where  $B$  represents the batch size,  $L$  the number of spatial locations, and  $c$  the channel index.

1512 However, modern models like Vision Transformer (ViT) utilize Layer Normalization (LN) instead  
 1513 of BN. Unlike BN, which calculates statistics across the entire batch, LN normalizes each instance  
 1514 independently by using the statistics calculated over individual feature dimensions. Specifically, for  
 1515 a feature vector  $\mathbf{f}_b$  belonging to the  $b$ -th instance, LN computes:

$$1516 \mu_b = \frac{1}{C} \sum_{c=1}^C f_{b,c}, \quad \sigma_b^2 = \frac{1}{C} \sum_{c=1}^C (f_{b,c} - \mu_b)^2, \quad (7)$$

1519 where  $C$  is the number of channels. This difference implies that LN operates without batch-level  
 1520 interactions, focusing solely on within-instance normalization, which makes the method inherently  
 1521 more suitable for handling variable batch sizes, particularly in latency-sensitive applications like  
 1522 those considered in our Test-Time Adaptation (TTA) setting.

1523 Despite the differences between BN and LN, the fundamental mechanism of using feature statistics  
 1524 to capture domain information remains valid. The key domain characteristics in early layer features  
 1525 are preserved in both normalization types, enabling the construction of a domain centroid that re-  
 1526 flects the distributional characteristics of the test data. For LN, this centroid can be computed by  
 1527 aggregating across instances instead of across batches:

$$1528 \bar{\mu}_c^{\text{LN}} = \frac{1}{M} \sum_{b=1}^M \mu_b, \quad \bar{\sigma}_c^{2\text{LN}} = \frac{1}{M} \sum_{b=1}^M \sigma_b^2, \quad (8)$$

1531 where  $M$  is memory capacity. This modified approach allows the domain centroid to still represent  
 1532 the overall domain-specific characteristics effectively, despite the lack of direct batch-level statistics.

1533 Furthermore, this methodology extends seamlessly to other normalization layers, such as Group  
 1534 Normalization (GN). In GN, the statistics are computed across smaller groups of channels within  
 1535 each instance, but the procedure for aggregating these statistics to form a domain centroid remains  
 1536 the same—by averaging the group-level statistics across instances.

1537 To maintain the core concept of selecting domain-representative samples with minimal modifica-  
 1538 tions, we continue to use the memory of high-confidence domain-representative samples in the  
 1539 Inference-only Batch-aware Memory Normalization (IoBMN) strategy. The adjustment for LN re-  
 1540 quires: 1. Calculating LN-specific centroids as described in Equation 8. 2. Replacing BN statistics  
 1541 with LN statistics in the IoBMN module, thereby aligning the feature normalization during inference  
 1542 with the domain-representative information derived from memory.

1543 The effectiveness of this modification was validated experimentally, as shown in Table 5, where ViT  
 1544 models using LN showed improved performance even under sparse TTA conditions. This indicates  
 1545 that, with minimal adjustments, SNAP-TTA remains effective for ViT with LN. The core principle  
 1546 of utilizing domain-representative statistics for aligning test-time feature distributions continues to  
 1547 provide significant benefits, ensuring robust adaptation in shifting domains with limited latency and  
 1548 computational overhead.

#### 1550 F.4 IMPACT OF MEMORY SIZE ON SNAP-TTA PERFORMANCE

1551 The memory size of the Class and Domain Representative Memory (CnDRM) in SNAP-TTA has  
 1552 implications for both performance and privacy. Increasing memory size allows storing more sam-  
 1553 ples, which intuitively could improve adaptation. However, such an approach raises privacy con-  
 1554 cerns and needs additional memory and latency when storing sensitive samples. To evaluate the  
 1555 trade-off, we conducted experiments on ImageNet-C under Gaussian noise corruption, using Tent +  
 1556 SNAP-TTA (adaptation rate 0.3) with a batch size of 16 and varying the memory size.

1558 As shown in Table 22, increasing the memory size beyond  
 1559 the base configuration of 16 does not lead to significant per-  
 1560 formance gains. This observation highlights the efficiency of  
 1561 SNAP-TTA’s representative sampling strategy, which priori-  
 1562 tizes storing samples based on proximity to class and domain  
 1563 centroids. The saturation in accuracy suggests that a carefully  
 1564 aligned memory size to the batch size is sufficient to balance  
 1565 computational efficiency, performance, and privacy considera-  
 tions.

Table 22: Performance comparison with varying memory sizes on ImageNet-C (Gaussian noise).

Memory Size	Accuracy (%)
16 (Base)	26.60
32	28.44
64	28.89
128	28.60

In conclusion, to minimize computational overhead while ensuring robust test-time adaptation, the memory size in SNAP-TTA is designed to align with the batch size. This configuration addresses privacy and memory overhead risks by limiting the number of stored samples without compromising adaptation effectiveness.

## F.5 EFFECT OF LEARNING RATE ON SPARSE AND FULL ADAPTATION

To investigate the impact of learning rates on the performance of SNAP-TTA and baseline methods, we conducted experiments under sparse adaptation settings. Initially, the same learning rate was applied for each SOTA TTA algorithms across all adaptation rates to ensure fair comparisons (Table 6, 7, 8, 9, 10, and 11). However, as sparse adaptation inherently limits the number of updates, the updates might be insufficient at lower adaptation rates and explored the effect of increasing the learning rate.

The results, summarized in Table 23, reveal that higher learning rates improve the accuracy of both the naive baseline and SNAP-TTA under sparse settings. Notably, while the naive TTA baseline benefits from a higher learning rate, its performance still falls short of that achieved with full adaptation. In contrast, SNAP-TTA surpasses the performance of full adaptation at optimal learning rates, demonstrating its ability to leverage sparse adaptation effectively. At the same time, applying these higher learning rates to full adaptation results in model instability and collapse, underscoring the need to carefully tune learning rates based on adaptation frequency. Therefore, we selected a stable learning rate of  $1 \times 10^{-4}$  for the evaluations in our work that balances model convergence and performance across all adaptation rates. These findings suggest that SNAP-TTA not only adapts effectively under sparse settings but also maintains robustness under optimized learning rates.

Table 23: Accuracy (%) with varying Learning Rates (LR) on ImageNet-C Gaussian noise adaptation rate 0.3.

LR	Tent(Full)	Tent(STTA)	Tent+SNAP	CoTTA(Full)	CoTTA(STTA)	CoTTA+SNAP	EATA(Full)	EATA(STTA)	EATA+SNAP
$2 \times 10^{-3}$	2.31	7.04	13.69	13.31	11.88	14.67	0.36	0.59	0.75
$1 \times 10^{-3}$	4.54	16.13	27.63	13.18	11.86	14.68	1.31	0.95	24.35
$5 \times 10^{-4}$	10.22	<b>24.96</b>	<b>29.95</b>	13.15	11.85	15.11	21.96	20.96	27.72
$1 \times 10^{-4}$	<b>27.03</b>	23.63	26.60	13.12	11.74	<b>15.26</b>	<b>29.42</b>	<b>27.35</b>	<b>29.48</b>
$5 \times 10^{-5}$	26.34	20.94	24.87	<b>13.34</b>	<b>11.92</b>	14.85	29.37	26.07	27.9

In conclusion, selecting an appropriately high learning rate for sparse adaptation significantly enhances performance while ensuring model stability. This strategy is particularly useful for real-world deployment of SNAP-TTA, where computational efficiency and robust performance are paramount.

## G LICENSE OF ASSETS

**Datasets** CIFAR10/CIFAR100 (MIT License), CIFAR10-C/CIFAR100-C (Creative Commons Attribution 4.0 International), and ImageNet-C (Apache 2.0).

**Codes** Torchvision for ResNet18, ResNet50, and ViTBase-LN (Apache 2.0), the official repository of CoTTA (MIT License), the official repository of Tent (MIT License), the official repository of EATA (MIT License), the official repository of SAR (BSD 3-Clause License), the official repository of RoTTA (MIT License), and the official repository of MECTA (Sony AI).