# ADVMEM: Adversarial Memory Initialization for Realistic Test-Time Adaptation via Tracklet-Based Benchmarking

Shyma Alhuwaider*    Motasem Alfarra    Juan C. Pérez    Merey Ramazanova    Bernard Ghanem
Center of Excellence in Generative AI, KAUST, Saudi Arabia

## Abstract

*We introduce a novel tracklet-based dataset for benchmarking test-time adaptation (TTA) methods. The aim of this dataset is to mimic the intricate challenges encountered in real-world environments such as images captured by hand-held cameras, self-driving cars, etc. The current benchmarks for TTA focus on how models face distribution shifts, when deployed, and on violations to the customary independent-and-identically-distributed (i.i.d.) assumption in machine learning. Yet, these benchmarks fail to faithfully represent realistic scenarios that naturally display temporal dependencies, such as how consecutive frames from a video stream likely show the same object across time. We address this shortcoming of current datasets by proposing a novel TTA benchmark we call the "Inherent Temporal Dependencies" (ITD) dataset. We ensure the instances in ITD naturally embody temporal dependencies by collecting them from tracklets—sequences of object-centric images we compile from the bounding boxes of an object-tracking dataset. We use ITD to conduct a thorough experimental analysis of current TTA methods, and shed light on the limitations of these methods when faced with the challenges of temporal dependencies. Moreover, we build upon these insights and propose a novel adversarial memory initialization strategy to improve memory-based TTA methods. We find this strategy substantially boosts the performance of various methods on our challenging benchmark. [1].*

## 1. Introduction

Deep neural networks (DNNs) have demonstrated impressive performance across various domains [15]. However, their reliability often diminishes in real-world scenarios due to natural corruptions and distribution shifts [16, 17, 21]. These shifts can manifest as unforeseen distortions that cause the input data to deviate from the model's training distribution. Additionally, the distribution of image classes may differ from what the model has learned, fur-

ther compounding the challenge. Consider images captured by hand-held cameras—these introduce two major difficulties: (1) the visual distribution may differ significantly from training data, such as in foggy or rainy conditions, and (2) images arrive sequentially as part of a continuous video stream. The first issue represents a distribution shift in the visual space, while the second introduces temporal dependencies that break the *independence* assumption inherent in conventional training (i.i.d.). Addressing both aspects is crucial for ensuring the robustness of DNNs in practical deployments.

**Test-Time Adaptation and its Challenges.** Test-Time Adaptation (TTA) seeks to mitigate performance degradation by adapting a pre-trained model on-the-fly using an incoming data stream [29]. At inference, TTA methods perform online unsupervised learning to adjust model parameters in response to new data [19, 27, 46, 51]. While TTA has shown promise, current benchmarks oversimplify the problem, primarily simulating distribution shifts without accounting for temporal dependencies that violate the i.i.d. assumption. For instance, many TTA approaches [27] assume that distribution shifts are purely covariate shifts [6], as seen in datasets like CIFAR10-C and ImageNet-C [16] (Fig. 1, left). Meanwhile, other methods [3] address non-i.i.d. scenarios by modifying label distributions while neglecting the visual continuity inherent in sequential data. A more recent effort by Yuan *et al.* [53] considers both distribution shifts and non-i.i.d. labels, but still overlooks the critical role of temporal dependencies. We argue that the lack of benchmarks that jointly capture distribution shifts and temporal dependencies has limited the development of deployable TTA methods. To bridge this gap, we take inspiration from object tracking to introduce a benchmark that inherently accounts for both challenges.

**Introducing ITD.** Our benchmark, Inherent Temporal Dependencies (ITD), is built using tracklets—short sequences of images tracking the same object across consecutive frames. By leveraging TrackingNet [34], we create a realistic test-time adaptation setting where temporal dependencies naturally emerge, leading to i.i.d. violations. To introduce controlled distribution shifts, we apply standard trans-
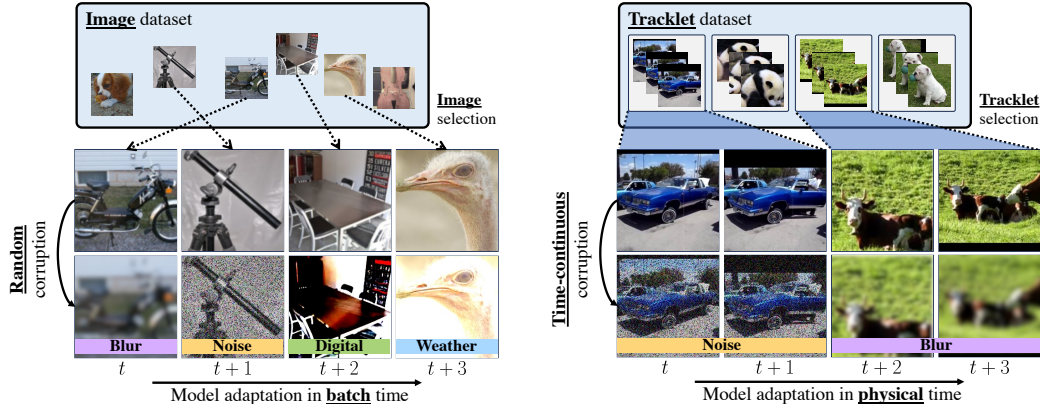
Figure 1. **A tracklet-based benchmark for realistic evaluation of Test-Time Adaptation (TTA) methods (Inherent Temporal Dependencies).** *(Left)* Existing benchmarks evaluate TTA methods using streams of images depicting different objects across batches, with random corruptions applied independently to each image. *(Right)* Our proposed ITD benchmark addresses these limitations by (i) presenting images of the same object in a sequence, preserving temporal dependencies from real-world tracklets, and (ii) applying consistent corruptions whose intensity may evolve over time. This framework offers a more realistic setting for evaluating the adaptability of TTA methods.

formations and corruptions [16, 21] (*e.g.* Gaussian noise, glass blur) consistently across tracklets rather than as independent perturbations (Fig. 1, right). This setup better reflects real-world challenges by aligning the temporal structure of the dataset with the sequential nature of data streams, following a protocol inspired by RoTTA [53]. Using ITD, we rigorously analyze how temporal dependencies interact with distribution shifts, revealing significant weaknesses in existing TTA methods.

**Advancing TTA with ADVMEM.** Our investigation highlights that most TTA methods struggle under the compounded effects of distribution shifts and temporal dependencies, leading to severe performance drops. We examine memory-bank-based methods, which should, in principle, handle temporal challenges well due to their ability to adapt to selectively stored samples. However, our results show that these methods suffer from poor initialization of the memory bank, significantly impacting performance.

To address this, we propose ADVMEM, a novel adversarial memory initialization strategy that enhances stability in adaptation. By leveraging synthetic noise generated in a class-diverse manner, ADVMEM operates as a plug-and-play enhancement for memory-based TTA methods. Experiments demonstrate its effectiveness—equipping SHOT-IM [28], with ADVMEM reduces error rates by 44% in Tracklet-Wise i.i.d. settings (see Table 2).

**Our Contributions.**

- **ITD Benchmark.** We introduce Inherent Temporal Dependencies, a novel benchmark for TTA that integrates object tracklets, capturing real-world distribution shifts and temporal dependencies.

- **Comprehensive Evaluation of TTA Methods.** We systematically assess existing TTA approaches on ITD, highlighting their limitations under realistic non-i.i.d. conditions.

- **ADVMEM.** We equip existing TTA methods with memory and benchmark their memory-adapted versions, demonstrating the impact of incorporating memory mechanisms on adaptation performance. Additionally, we propose an adversarial memory initialization strategy that significantly improves the performance, particularly under severe non-i.i.d. scenarios.

Our work advances the study of test-time adaptation by providing a more realistic evaluation framework and a novel solution to enhance the stability of model adaptation in dynamic environments.

## 2. Related Work

**Test-Time Adaptation.** TTA leverages the unlabeled data that arrives at test time to adapt the forward pass of pre-trained DNNs according to some proxy task [28, 50]. Many existing TTA methods focus on covariate distribution shifts [26, 28, 35, 36, 52]. Several TTA methods tackle this challenge by updating the statistics of the Batch Normalization layers at test time [26, 43]. For example, AdaBN [26] introduces Adaptive Batch Normalization, an algorithm to adapt to the target domain. Another group of methods uses an entropy minimization strategy. For instance, TENT [50] minimizes the entropy of the model's predictions. ETA and EATA [35] extend TENT by selecting reliable and non-redundant samples to update the model weights. More recently, RoTTA [53] attempts to combat non-i.i.d. streams at test time by leveraging a memory bank for adapting to an
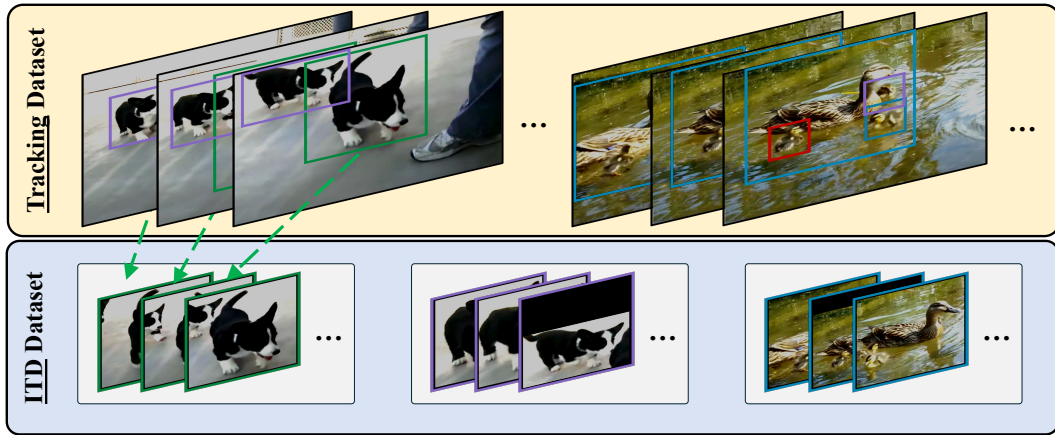
Figure 2. **We build ITD with realistic TTA instances by constructing them from a tracking dataset.** We extract object-centric sequential video frames, encapsulating the small variations of the same entity over time. We source the frames and bounding boxes from TrackingNet, a well-established tracking dataset, such that the instances focus on particular objects of interest. As such, these sequences naturally exhibit the temporal dependencies inherent in real-world scenarios.

incoming stream of data. In this work, we introduce AD-VMEM, a novel adversarial memory initialization strategy to significantly enhance the adaptability of TTA methods under complex, non-i.i.d. scenarios.

**Benchmarking TTA Methods.** The fundamental premise of TTA involves deploying a pre-trained model onto edge devices like self-driving cars or surveillance cameras, where it faces potential changes in data distribution [19, 27, 47, 51]. This scenario unfolds as the model encounters a continuous stream of data, with each input potentially coming from a distribution different than the one the model was originally trained on. To emulate such scenarios, the TTA literature commonly creates a stream of data with samples from the test set of well-established image classification datasets, such as ImageNet [9] and CIFAR [24]. Setups then systematically simulate covariate distribution shifts by inducing corruptions on individual images, such as those from Common Corruptions [16] and 3D Common Corruptions [21]. In this work, we present a comprehensive benchmark for simulating more realistic and complex scenarios.

## 3. Dataset and Methodology

**Motivation.** Recent advancements in the field of Test Time Adaptation (TTA) have departed from traditional independent and identically distributed (i.i.d.) setups towards more nuanced non-i.i.d. configurations. RoTTA [53] introduced correlation sampling to enforce non-i.i.d. distributions in labels, motivated by real-world deployment scenarios—e.g., in edge devices, where objects in a scene often appear with correlated labels (e.g., frequent occurrences of "pedestrian" in a crowded area). This shift revealed a critical limitation: existing TTA methods struggle when faced with such non-i.i.d. streams. In real-world streaming data (e.g., videos from surveillance cameras), consecutive frames often depict the same object with minor variations, leading to visual redundancy. To highlight the impact of temporal dependen-

Table 1. **Average Error Rates on CIFAR-10-C:** Comparison between non-i.i.d. episodic evaluation and tracklet mimic evaluation averaged across corruptions. The tracklet mimic setting simulates real-world temporal dependencies.

| TTA Method | non i.i.d. (%) | Tracklet Mimic (%) | Δ (%) |
|---|---|---|---|
| Source | 44.1 | **44.1** | **0** |
| AdaBN | 75.4 | 78.7 | 3.3 |
| CoTTA | 75.5 | 89.1 | 13.6 |
| SAR | 75.2 | 82.4 | 7.2 |
| ETA | 75.4 | 78.6 | <u>3.2</u> |
| TENT | 75.3 | 85.1 | 9.8 |
| RoTTA | <u>27.6</u> | <u>66.8</u> | 39.2 |

cies on TTA performance, we conducted a simple experiment on CIFAR-10-C. In this experiment, we compared the error rates of several TTA methods under PTTA [53] evaluation against a modified setup where each batch contains images duplicated to mimic a video clip or tracklet. As shown in Table 1, all methods experience a noticeable performance drop in the tracklet mimic setting. This finding clearly illustrates the challenges posed by temporal dependencies and motivates the need for our new ITD benchmark, which better reflects real-world data streams.

These observations highlight the need for adaptation strategies robust to *visual* non-i.i.d. shifts and a benchmark that captures these complexities. To this end, we introduce ITD, a benchmark built explicitly for evaluating methods under complex non-i.i.d. scenarios, and propose ADVMEM, a plugin designed to mitigate over-adaptation in such conditions.

**Tracklets: A Natural Source of Visual Non-i.i.d. Data.** To construct a realistic non-i.i.d. benchmark, we leverage the field of Object Tracking. Unlike artificially simulated correlation sampling used in prior works [53], object-tracking datasets inherently capture realistic temporal de-

pendencies by tracking specific objects across video frames. We propose using *tracklets*—sequences of object-centric images extracted from tracking datasets to model the gradual variations encountered in real-world image streams. This approach not only enhances the realism of our benchmark but also faithfully captures intrinsic characteristics of natural image sequences, establishing a robust foundation for evaluating TTA in real-world scenarios.

## 3.1. Dataset Construction

We construct Inherent Temporal Dependencies (ITD) from a large-scale object-tracking dataset. Specifically, we utilize TrackingNet [34], which originates from the YouTube Bounding-Boxes dataset [40].

For each video, we extract tracklets by iterating through frames where a given object appears, cropping bounding boxes around it. This process results in ITD, a dataset composed of object tracklets—sequences of images capturing realistic temporal variations of objects. Figure 2 illustrates the tracklet extraction process.

In practice, we apply the following preprocessing steps:

- **Frame Selection:** Extract crops at a 5-frame interval to balance dataset size and temporal redundancy.
- **Crop Resizing:** Extract square crops 10% larger than the largest side of the bounding box to retain contextual information.
- **Standardization:** Resize all crops to 224×224 for consistency and compatibility with batch-based training.

These design choices ensure a realistic yet tractable dataset, allowing for systematic evaluation of TTA methods.

**Dataset Properties.** Unlike conventional datasets where samples are independent, ITD is composed of *tracklets*, preserving the temporal continuity of objects. Each tracklet consists of images depicting the same object in different frames, naturally encoding non-i.i.d. dependencies. Additionally, our dataset supports temporally consistent corruptions (detailed in Section 4.4), further enhancing its relevance for evaluating TTA under realistic conditions.

**Statistics.** The ITD dataset contains over 23K objects that span 21 classes. The dataset is divided into training (50%), validation (30%), and test (20%) sets. In total, it comprises over 220K images—more than four times the size of ImageNet-C (50K)—while also providing object-instance relationships via tracklets. Further statistics, including class distributions, are provided in the appendix underscoring the scale and diversity of ITD, making it a valuable resource for advancing TTA research.

## 4. Benchmarking on ITD

**Overview.** Unlike previous benchmarks that assume independent samples, ITD introduces a tracklet-based evaluation to reflect real-world challenges, such as sequential dependencies and non-i.i.d. distributions. We systematically evaluate TTA methods across three levels of complexity to assess their adaptability in streaming environments:

- **Frame-wise i.i.d. (Section 4.5):** Frames are sampled independently and identically distributed (i.i.d.), without considering sequential dependencies.
- **Tracklet-wise i.i.d. (Section 4.6):** Entire tracklets are sampled i.i.d., preserving intra-tracklet dependencies while maintaining inter-tracklet randomness.
- **Tracklet-wise non-i.i.d. (Section 4.7):** Tracklets are sampled following a Dirichlet distribution [53], enforcing stronger non-i.i.d. properties across the dataset.

These setups enable a systematic evaluation of TTA methods under increasingly complex real-world conditions. Each scenario is tested under a single domain shift (e.g., fog) to isolate the effect of adaptation techniques.

## 4.1. Tracklet-Based Adaptation Strategies

Unlike frame-wise adaptation, tracklet-based setups introduce sequential dependencies, making naive entropy minimization unreliable due to biased batch statistics. To ensure fair comparisons, we extend TENT and SHOT-IM by incorporating RoTTA's memory bank $\mathcal{M}$. This allows them to store and utilize previously seen samples, ensuring more stable adaptation in non-i.i.d. streams. In these adaptations, we replace RoTTA's original objective with entropy minimization for TENT and information maximization for SHOT-IM.

## 4.2. Experimental Setup

Having established the dataset structure, corruption strategies, and adaptation mechanisms, we now outline our standardized experimental setup for evaluating TTA methods. Unless stated otherwise, we use ResNet-18 as the base model, apply corruptions at the highest severity level (5), and set a batch size of 64 for streaming data. We evaluate eight TTA methods and compare them against a pre-trained model $f_\theta$(Source) as the baseline. The details of these methods are provided in Table 2. Each method is assessed using its optimal hyperparameters, determined through an extensive search.

Table 2. Overview of TTA Methods Evaluated in ITD.

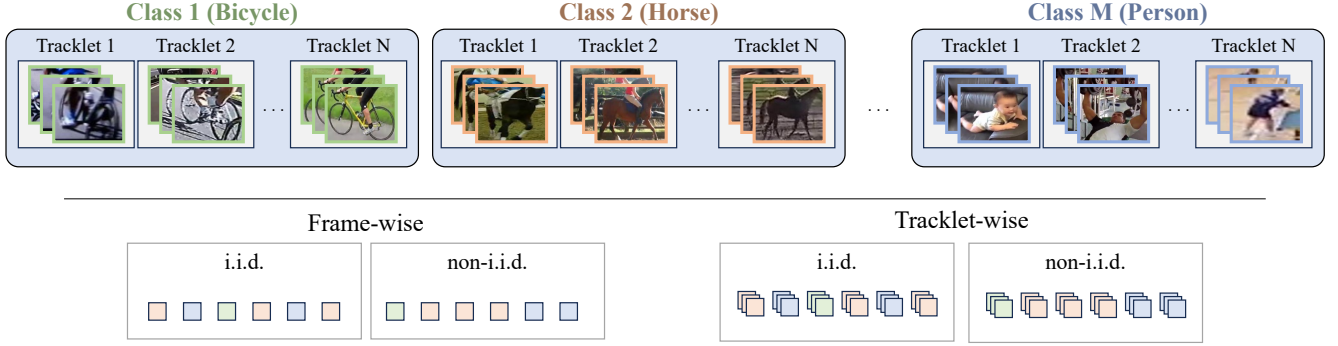| Method | Adaptation Strategy |
|---|---|
| **AdaBN** [26] | Updates batch normalization statistics. |
| **SHOT-IM** [28] | Maximizes mutual information. |
| **TENT** [50] | Utilizes entropy minimization for adaptation. |
| **SAR** [36] | Employs sharpness-aware optimization. |
| **EATA** [35] | Entropy-based sample selection. |
| **CoTTA** [52] | Uses consistency-based distillation for continual adaptation. |
| **RoTTA** [53] | Maintains a memory bank to stabilize adaptation. |

Figure 3. **Frame-wise and Tracklet-wise Experiment Setup:** We illustrate the construction of the frame-wise and tracklet-wise experiments. In the frame-wise setup, one frame is sampled from each tracklet to ensure each object is observed once. In the tracklet-wise setup, the frames within each tracklet are sequentially processed. Both i.i.d. and non-i.i.d. settings are depicted for each setup.

Table 3. **Types of Corruptions Applied in ITD.** To evaluate TTA robustness, we introduce common distribution shifts from ImageNet-C.

| Category | Corruptions |
|---|---|
| **Noise** | Gaussian, Shot, Impulse |
| **Blur** | Defocus, Glass, Zoom |
| **Weather** | Snow, Frost, Fog, Brightness |
| **Digital Artifacts** | Contrast, Elastic Transform, Pixelate, JPEG Compression |

### 4.3. Considered Corruptions

Corruptions, detailed in Table 3, are applied dynamically as the data stream unfolds. Given the temporal nature of ITD, we also explore scenarios where corruption severity varies within a tracklet, simulating transient environmental fluctuations (e.g., changing weather conditions). Additional extended evaluations can be found in the appendix.

### 4.4. Preparation for TTA

To evaluate the effectiveness of TTA methods on our ITD benchmark, we require models that are either pre-trained or fine-tuned specifically on this dataset. While ITD shares some class overlap with ImageNet, its distribution differs significantly, limiting the effectiveness of direct transfer. Our experiments confirm that ImageNet-pretrained models struggle to generalize well to ITD, highlighting the need for dataset-specific fine-tuning. Therefore, we fine-tune two ImageNet-pretrained models, ResNet-18 and ViT-B-16, on ITD's training set (see Table 4).

Table 4. **Error rates on all splits of ITD.** Detailed training and test loss/accuracy results are provided in the appendix.

| Error Rate | ResNet-18 | ViT |
|---|---|---|
| Train | 4.1 | 2.0 |
| Validation | 8.2 | 3.8 |
| Test | 9.4 | 4.0 |

### 4.5. Frame-Wise i.i.d. Scenario

In this scenario, we evaluate TTA methods under a conventional setting where each frame is treated as an independent and identically distributed (i.i.d.) sample, ignoring temporal dependencies. This setup assumes a uniform label distribution across time, thereby oversimplifying real-world conditions where label distributions may be highly imbalanced.

To construct the test stream, we shuffle tracklets and sample one frame per tracklet, eliminating the notion of temporal continuity. Following standard practice [16, 21], we assess performance degradation under 15 different corruptions. Table 5 presents the average error rates of the eight evaluated TTA methods. As expected, distribution shifts significantly degrade the performance of pre-trained models. For example, the error rate of the Source model (ResNet-18 without TTA) rises from below 10% (Table 4) to over 60% (Table 5). Notably, TENT [50] reduces the error rate to below 50% through entropy minimization, while SHOT-IM achieves the lowest error, averaging below 40%.

### 4.6. Tracklet-Wise i.i.d. Scenario

To move towards a more realistic evaluation, we introduce a setup where the model processes entire tracklets at test time. Within each tracklet, consecutive frames share labels and contextual consistency, though tracklets are sampled in an i.i.d. manner. This setup emulates real-world scenarios where the model observes a single object over multiple frames before switching to a new one.

We evaluate SHOT-IM (the best performer from Section 4.5), TENT, and RoTTA. Although RoTTA performed poorly in the frame-wise i.i.d. scenario, its memory-based approach is specifically designed for non-i.i.d. streams, making it relevant for this setting. To ensure a fair comparison, we extend TENT and SHOT-IM by incorporating RoTTA's memory bank, allowing them to only adapt to informative samples selected and retained in memory based

Table 5. **Performance Comparison under Frame-wise i.i.d. Assumption.** Average error rates reported for TTA methods on images from the ITD dataset. In this setup, all images/frames are shuffled and then independently subjected to corruptions. Notably, SHOT-IM outperforms other methods across all corruptions, showcasing its robustness against these domain shifts.

| Method | Source | AdaBN | SHOT-IM | TENT | SAR | CoTTA | ETA | EATA | RoTTA |
|---|---|---|---|---|---|---|---|---|---|
| Avg. Err. ↓ | 62.4 | 46.9 | **39.3** | 46.8 | 46.5 | 46.7 | 46.7 | 46.8 | 53.4 |

on RoTTA Category-balanced sampling heuristics.

Table 6 summarizes our results. We find that SHOT-IM and TENT struggle under tracklet-based evaluation, with SHOT-IM's error rate increasing from under 40% (Table 5) to nearly 95% (Table 6). This decline stems from the biased statistics computed within tracklets, which skew entropy-based adaptation. In contrast, RoTTA demonstrates superior stability, reducing the average error to around 50%, due to its distillation-based approach.

Table 6. **Tracklet-wise i.i.d. and Tracklet-wise *non*-i.i.d. Evaluation with and without Memory.** We report average error rates for TTA methods on our ITD dataset under both the tracklet-wise i.i.d. scenario (*i.e.* entire tracklets are sampled i.i.d.) and when tracklets are non-i.i.d., *i.e.* tracklets are sampled such that their labels display correlation in time (by following a Dirichlet distribution). The results are grouped to reflect the presence (✓) or absence (✗) of memory during adaptation. RoTTA shows notable robustness by using memory, significantly outperforming the non-memory variants across various corruption types. Under *non*-i.i.d. RoTTA showcases it's proficiency against non-i.i.d. data streams.

| Method | Memory | Error Rate↓ Tracklet-wise i.i.d. | Error Rate ↓ Tracklet-wise *non*-i.i.d. |
|---|---|---|---|
| TENT | ✗ | 94.0 | 94.0 |
|  | ✓ | 93.8 | 93.8 |
| SHOT-IM | ✗ | 94.7 | 95.1 |
|  | ✓ | 93.4 | 93.6 |
| RoTTA | ✓ | **51.3** | **79.3** |

### 4.7. Tracklet-Wise non-i.i.d. Scenario

In this final and most challenging setup, tracklets are sampled non-i.i.d. to simulate real-world streaming conditions such as autonomous driving or surveillance, where object categories appear in bursts. To model this, we follow Yuan *et al.* [53] and sample tracklets using a Dirichlet distribution $\text{Dir}(\gamma)$. As $\gamma \to 0$, label correlation within the stream increases, deviating from the i.i.d. assumption.

We evaluate RoTTA, SHOT-IM, and TENT under $\gamma = 10^{-4}$ to enforce strong non-i.i.d. conditions (additional $\gamma$ values are analyzed in Section 5). Table 6 reports the results. Even RoTTA, designed for non-i.i.d. streams, experiences a 28% performance drop compared to the tracklet-wise i.i.d. setup (Table 6). We hypothesize that this decline results from imbalanced memory due to empty memory initialization, where certain classes are observed late in the

stream, causing forgetting effects. These findings suggest that memory-based TTA methods can be further improved by introducing class-balancing initialization mechanisms to stabilize adaptation over long, non-i.i.d. streams

## 5. Experiments: Enhancing Performance with Memory

In Section 4.6, we observed that equipping TENT and SHOT-IM with a memory bank, while seemingly beneficial, does not yield significant performance improvements. Furthermore, in Section 4.7, we demonstrated how a tracklet-wise non-i.i.d. stream severely impacts RoTTA's performance. These findings indicate that while memory banks are essential for adapting to non-i.i.d. streams [53], even strong TTA methods experience significant degradation when evaluated on ITD.

We hypothesize that this degradation stems from the empty initialization of the memory bank. Consider the extreme case where $\gamma \to 0$, resulting in the memory bank lacking any examples for labels revealed later in the stream until those labels actually appear. Additionally, methods such as TENT rely on statistical measures like entropy for updates. When key classes are absent from the memory bank, model updates become skewed, leading to catastrophic forgetting. This motivates the need for a carefully designed memory initialization strategy that ensures stable adaptation steps.

### 5.1. Adversarial Memory Initialization

To address these challenges, we propose a novel approach for initializing the memory bank to enhance adaptation. Our goal is to populate the memory bank with class-representative samples to: **(i)** prevent forgetting by ensuring all classes are accounted for, and **(ii)** balance computed statistics in the output space.

To that end, we propose initializing the memory bank with synthetic data generated by adversarial algorithms [49]. Specifically, each memory bank entry is initialized as Gaussian noise, assigned a random label, and subjected to a targeted adversarial attack that maximizes the network's confidence in classifying it correctly, following [1, 14]. Formally, let $\mathcal{M}$ be an initially empty memory bank with a maximum capacity of $N$. We populate $\mathcal{M}$ iteratively with synthetic examples $x^*$, where:

$$x^* = \underset{x}{\arg\min}\, \mathcal{L}_{\text{ce}}(f_\theta(x), y), \quad (1)$$

6

where $y$ is a randomly assigned label, and $\mathcal{L}_{ce}$ represents the cross-entropy loss. We solve this optimization problem by applying gradient descent, starting from Gaussian noise. This process is repeated $N$ times to fully initialize $\mathcal{M}$. We term this procedure "ADVMEM" and present it in Algorithm 1. Our adversarial memory initialization offers two key advantages: **(i)** The memory bank remains populated with class-representative samples throughout adaptation, mitigating forgetting under strong non-i.i.d. streams. **(ii)** The initialization method is independent of how $\mathcal{M}$ is updated or utilized. For example, when applied to RoTTA, the adaptation and update mechanisms remain unchanged. A straightforward alternative would be initializing $\mathcal{M}$ with uniformly sampled, non-corrupted training examples. However, even in cases where privacy concerns are not an issue, our experiments (detailed in the appendix) indicate that this approach does not improve performance.

**Performance with ADVMEM.** We implement ADVMEM initialization strategy within RoTTA, TENT, and SHOT-IM, replacing their default empty memory initialization with adversarially generated samples. We then evaluate on the experimental setups from Sections 4.6 and 4.7 to assess its impact on adaptation performance. **Tracklet-wise**

---

**Algorithm 1** ADVMEM

**function** INITIALIZEMEMORY($f_\theta, K, N$)
    Initialize $\mathcal{M} = \{\}$
    **while** $|\mathcal{M}| < N$ **do**
        $x \sim \mathcal{N}(0, I), y \sim \mathcal{U}\{1, 2, \ldots, K\}$
        **while** $f_\theta(x) \neq y$ **do**
            $x \leftarrow x - \alpha \cdot (\nabla_x \mathcal{L}_{ce}(f_\theta(x), y))$
        **end while**
        $\mathcal{M} \leftarrow \mathcal{M} \cup x$
    **end while**
    **return** $\mathcal{M}$
**end function**

---

**i.i.d. Setup:** We first analyze the Tracklet-wise i.i.d. setting from Section 4.6 (i.e., $\gamma \to \infty$). Table 2 reports the results, showing significant performance improvements across all baselines. Notably, ADVMEM reduces the average error rate of TENT by $\sim$15% and that of SHOT-IM by over 40%, making SHOT-IM the top-performing method. These improvements highlight the effectiveness of ADVMEM in stabilizing adaptation.

**Tracklet-wise non-i.i.d. Setup:** To further examine its impact, we evaluate ADVMEM in the extreme non-i.i.d. setting from Section 4.7 ($\gamma \to 0$). Table 3 presents the results, demonstrating substantial performance gains. In particular, ADVMEM enhances RoTTA's performance by an average of 4%, with specific improvements of 16% and 10% against pixelate and zoom corruptions, respectively. This trend holds for other methods as well, with TENT improv-

ing by over 10% on JPEG corruption and achieving a 5% average improvement across all corruptions.

# 6. Ablation Studies and Analysis

This section presents an in-depth analysis of the ITD dataset and our proposed ADVMEM. We specifically evaluate both the ResNet-18 and Vision Transformer (ViT) architectures to investigate the impact of ADVMEM on performance across different experimental conditions.

## 6.1. Controlling Label Distribution

Using a Dirichlet distribution $\text{Dir}(\gamma)$ for sampling labels, as introduced in [53], allows control over the distributional shift in the label space through the parameter $\gamma$. Specifically, as $\gamma \to 0$, we approach a class-incremental non-i.i.d. setup, while as $\gamma \to \infty$, we transition towards a uniform i.i.d. setup.

We extend our experiments by analyzing intermediate stages with $\gamma \in \{10^{-4}, 10^{-1}, 10^3\}$ and report the results in Figure 4a.

For low values of $\gamma$, the model encounters a challenging class-incremental non-i.i.d. scenario. In this case, ADVMEM proves instrumental in mitigating class bias within the incoming image stream, reducing degradation in performance. Conversely, as $\gamma$ increases, the scenario becomes easier, as the model has higher chances of encountering uniformly sampled classes. In this context, the adversarially initialized memory samples are rapidly replaced by reliable examples from the stream, making the initialization inconsequential. Consequently, ADVMEM neither enhances nor degrades performance in the uniform i.i.d. setup.

## 6.2. Sensitivity to Batch Size

In this section, we analyze the impact of batch size on performance by varying it across $\{8, 16, 32\}$. The results, reported in Figure 4b, confirm that increasing the batch size improves performance, as larger batches expose models to more diverse examples, facilitating better adaptation.

When comparing i.i.d. and non-i.i.d. setups ($\gamma \to \infty$ vs. $\gamma \to 0$), we observe that methods incorporating ADVMEM consistently achieve improved or maintained performance across batch sizes. However, as in Section 8.3, the impact of ADVMEM is less pronounced in the i.i.d. setup, indicating that memory initialization has a limited effect in this scenario.

These insights highlight the robustness of ADVMEM in highly non-i.i.d. environments while demonstrating that its advantages diminish in simpler, uniform settings. Overall, the findings reinforce the necessity of well-designed memory initialization strategies in real-world, dynamically shifting data distributions.

Table 2. **Effect of Adversarial Memory Initialization on TTA performance in the Tracklet-Wise i.i.d. scenario.** We evaluate TTA methods on ITD in an i.i.d. tracklet context, contrasting standard (✗) and adversarial (✓) memory initialization (ADVMEM). When we equip SHOT-IM with ADVMEM, it outperforms all other methods, demonstrating its capacity to enhance adaptability.

| Method | ADVMEM | Noise | | | Blur | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | defoc. | glass | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| TENT | ✗ | 94.2 | 94.2 | 94.3 | 94.5 | 94.5 | 93.9 | 93.7 | 93.7 | 93.2 | 93.0 | 92.8 | 93.8 | 93.4 | 93.5 | 93.8 |
| | ✓ | 85.1 | 83.2 | 88.2 | 82.9 | 84.8 | 74.0 | 85.5 | 84.2 | 83.9 | 68.9 | 71.6 | 79.9 | 62.6 | 69.5 | 78.9 |
| SHOT-IM | ✗ | 93.8 | 93.9 | 93.9 | 94.2 | 94.3 | 93.7 | 93.3 | 93.0 | 93.0 | 92.8 | 92.1 | 93.5 | 93.4 | 93.2 | 93.4 |
| | ✓ | **61.7** | **58.0** | **62.3** | **52.5** | **51.0** | **39.7** | 57.2 | 60.6 | 51.4 | 32.8 | **55.0** | **38.6** | **27.6** | 35.6 | **48.9** |
| RoTTA | ✗ | 68.0 | 62.4 | 68.6 | 58.5 | 56.7 | 40.9 | **56.6** | 60.7 | 52.2 | 30.2 | 60.7 | 39.7 | 27.8 | **35.1** | 51.3 |
| | ✓ | 69.0 | 62.8 | 68.9 | 58.1 | 58.0 | 40.9 | 57.2 | 61.0 | 53.0 | **30.0** | 62.4 | 41.2 | 27.7 | 35.8 | 51.9 |

Table 3. **Effect of Adversarial Memory Initialization on TTA performance in the Tracklet-Wise *non*-i.i.d. scenario.** We evaluate methods on ITD in a non-i.i.d. tracklet setting, comparing standard (✗) and adversarial (✓) memory initialization (ADVMEM). RoTTA with ADVMEM significantly outperforms the alternatives.

| Method | ADVMEM | Noise | | | Blur | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | defoc. | glass | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| TENT | ✗ | 94.2 | 94.3 | 94.3 | 94.5 | 94.5 | 94.0 | 93.7 | 93.6 | 93.4 | 93.0 | 92.9 | 93.8 | 93.4 | 93.6 | 93.8 |
| | ✓ | 89.2 | 89.8 | 92.4 | 92.1 | 92.4 | 86.6 | 91.4 | 90.2 | 90.7 | 83.7 | 83.9 | 89.0 | 79.5 | 81.5 | 88.0 |
| SHOT-IM | ✗ | 93.9 | 93.8 | 93.9 | 94.3 | 94.3 | 93.7 | 93.4 | 93.6 | 93.2 | 92.7 | 92.5 | 93.6 | 93.5 | 93.2 | 93.6 |
| | ✓ | 92.3 | 92.2 | 92.7 | 92.4 | 93.1 | 91.6 | 92.5 | 91.8 | 91.5 | 90.5 | 92.0 | 91.6 | 89.1 | 90.9 | 91.7 |
| RoTTA | ✗ | 85.6 | **81.6** | 86.5 | 86.9 | 87.3 | 78.9 | 81.7 | 83.5 | 75.0 | 67.2 | **71.8** | 79.7 | 77.1 | 67.4 | 79.3 |
| | ✓ | **84.4** | 82.2 | **84.1** | **83.0** | **83.3** | **68.6** | **80.2** | **80.5** | **74.7** | **62.8** | 73.6 | **75.2** | **61.3** | **62.9** | **75.5** |



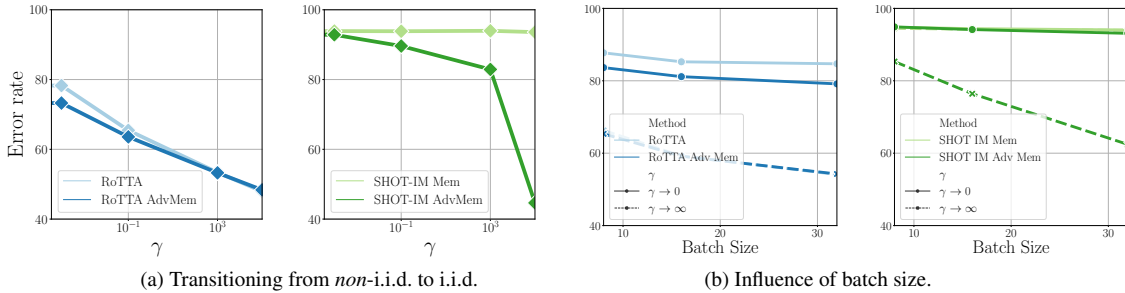(a) Transitioning from *non*-i.i.d. to i.i.d.   (b) Influence of batch size.

Figure 4. **Error rate as (a) we transition the evaluation scenario from non-i.i.d. to i.i.d., and (b) as we vary the batch size.** In **(a)**, we control the i.i.d. nature of the label distribution by varying the $\gamma$ parameter in the Dir($\gamma$) distribution ($\gamma$-axis in log-scale). Adding our proposed ADVMEM consistently enhances or maintains performance across different i.i.d.-ness regimes and methods. In **(b)**, we examine how batch size influences method performance, showing that larger batch sizes improve performance, with ADVMEM further enhancing results. All results presented here are for ResNet-18.

## 7. Conclusion

This work introduces ITD, a novel benchmark designed to challenge existing TTA methods. Unlike previous benchmarks, ITD captures the temporal dependencies inherent in real-world data streams, an aspect often overlooked in traditional evaluations. Additionally, we propose ADVMEM, an adversarial memory initialization strategy that enhances the adaptability of TTA methods. By preloading the memory bank with adversarially crafted samples, ADVMEM effectively mitigates model forgetfulness and leads to significant performance improvements, particularly in non-i.i.d. scenarios. Our findings advocate for a paradigm shift towards benchmarks that more accurately reflect real-world complexities. We emphasize the necessity of robust adaptation strategies capable of handling evolving data distributions, paving the way for the next generation of resilient machine learning models. Future research should explore further enhancements to memory-based adaptation techniques and extend the scope of realistic TTA benchmarks.

## Acknowledgements

## References

[1] Motasem Alfarra, Juan C Pérez, Ali Thabet, Adel Bibi, Philip HS Torr, and Bernard Ghanem. Combating adversaries with anti-adversaries. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5992–6000, 2022. 6

[2] Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. *arXiv preprint arXiv:2002.06470*, 2020.

[3] Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8344–8353, 2022. 1

[4] Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8344–8353, 2022.

[5] Zhipeng Cai, Ozan Sener, and Vladlen Koltun. Online continual learning with natural distribution shifts: An empirical study with visual data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8281–8290, 2021.

[6] J Quinonero Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. Dataset shift in machine learning. *The MIT Press*, 1:5, 2009. 1

[7] Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. Contrastive test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 295–305, 2022.

[8] Francesco Croce, Sven Gowal, Thomas Brunner, Evan Shelhamer, Matthias Hein, and Taylan Cemgil. Evaluating the adversarial robustness of adaptive test-time defenses. In *International Conference on Machine Learning*, pages 4421–4435. PMLR, 2022.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 3

[10] Jin Gao, Jialing Zhang, Xihui Liu, Trevor Darrell, Evan Shelhamer, and Dequan Wang. Back to the source: Diffusion-driven test-time adaptation. *arXiv preprint arXiv:2207.03442*, 2022.

[11] Yasir Ghunaim, Adel Bibi, Kumail Alhamoud, Motasem Alfarra, Hasan Abed Al Kader Hammoud, Ameya Prabhu, Philip HS Torr, and Bernard Ghanem. Real-time evaluation in online continual learning: A new paradigm. *arXiv preprint arXiv:2302.01047*, 2023.

[12] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. *arXiv preprint arXiv:1803.07728*, 2018.

[13] Taesik Gong, Jongheon Jeong, Taewon Kim, Yewon Kim, Jinwoo Shin, and Sung-Ju Lee. Note: Robust continual test-time adaptation against temporal correlation. In *Advances in Neural Information Processing Systems*.

[14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. 6

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1

[16] Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019. 1, 2, 3, 5

[17] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, Dawn Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *ICCV*, 2021. 1

[18] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

[19] Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic domain generalization. *Advances in Neural Information Processing Systems*, 34:2427–2440, 2021. 1, 3

[20] Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic domain generalization. *Advances in Neural Information Processing Systems*, 34:2427–2440, 2021.

[21] Oğuzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18963–18974, 2022. 1, 2, 3, 5

[22] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.

[23] Takeshi Kojima, Yutaka Matsuo, and Yusuke Iwasawa. Robustifying vision transformer without retraining from scratch by test-time class-conditional feature alignment. *arXiv preprint arXiv:2206.13951*, 2022.

[24] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 3

[25] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, page 896, 2013.

[26] Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch normalization for practical do-

main adaptation. *arXiv preprint arXiv:1603.04779*, 2016. 2, 4

[27] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. *CoRR*, abs/2002.08546, 2020. 1, 3

[28] Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source hypothesis transfer for unsupervised domain adaptation. In *International Conference on Machine Learning*, pages 6028–6039. PMLR, 2020. 2, 4

[29] Jian Liang, Ran He, and Tieniu Tan. A comprehensive survey on test-time adaptation under distribution shifts, 2023. 1

[30] Yuejiang Liu, Parth Kothari, Bastien Van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre Alahi. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural Information Processing Systems*, 34: 21808–21820, 2021.

[31] Massimiliano Mancini, Hakan Karaoguz, Elisa Ricci, Patric Jensfelt, and Barbara Caputo. Kitting in the wild through online domain adaptation. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1103–1109. IEEE, 2018.

[32] M Jehanzeb Mirza, Jakub Micorek, Horst Possegger, and Horst Bischof. The norm must go on: dynamic unsupervised domain adaptation by normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14765–14775, 2022.

[33] Muhammad Jehanzeb Mirza, Pol Jané Soneira, Wei Lin, Mateusz Kozinski, Horst Possegger, and Horst Bischof. Actmad: Activation matching to align distributions for test-time-training, 2022.

[34] Matthias Muller, Adel Bibi, Silvio Giancola, Salman Alsubaihi, and Bernard Ghanem. Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In *Proceedings of the European conference on computer vision (ECCV)*, pages 300–317, 2018. 1, 4

[35] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui Tan. Efficient test-time model adaptation without forgetting. In *International conference on machine learning*, pages 16888–16905. PMLR, 2022. 2, 4

[36] Shuaicheng Niu14, Jiaxiang Wu, Yifan Zhang, Zhiquan Wen, Yaofo Chen, Peilin Zhao, and Mingkui Tan15. Towards stable test-time adaptation in dynamic wild world. 2, 4, 13

[37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[38] Juan C Pérez, Motasem Alfarra, Guillaume Jeanneret, Laura Rueda, Ali Thabet, Bernard Ghanem, and Pablo Arbeláez. Enhancing adversarial robustness via test-time transforma-

tion ensembling. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 81–91, 2021.

[39] Ori Press, Steffen Schneider, Matthias Kümmerer, and Matthias Bethge. Rdumb: A simple approach that questions our progress in continual test-time adaptation. *Advances in Neural Information Processing Systems*, 36, 2024.

[40] Esteban Real, Jonathon Shlens, Stefano Mazzocchi, Xin Pan, and Vincent Vanhoucke. Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. *CoRR*, abs/1702.00824, 2017. 4

[41] Kate Saenko, Brian Kulis, Mario Fritz, and Trevor Darrell. Adapting visual category models to new domains. In *Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part IV 11*, pages 213–226. Springer, 2010.

[42] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Acdc: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10765–10775, 2021.

[43] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. *Advances in Neural Information Processing Systems*, 2020. 2

[44] Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends® in Machine Learning*, 4(2):107–194, 2012.

[45] Yongyi Su, Xun Xu, and Kui Jia. Revisiting realistic test-time training: Sequential inference and adaptation by anchored clustering. *arXiv preprint arXiv:2206.02721*, 2022.

[46] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pages 9229–9248. PMLR, 2020. 1

[47] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning*, pages 9229–9248. PMLR, 2020. 3

[48] Antonio Torralba and Alexei A Efros. Unbiased look at dataset bias. In *CVPR 2011*, pages 1521–1528. IEEE, 2011.

[49] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7167–7176, 2017. 6

[50] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020. 2, 4, 5

[51] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020. 1, 3

[52] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In *Proceedings of*

the *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7201–7211, 2022. 2, 4

[53] Longhui Yuan, Binhui Xie, and Shuang Li. Robust test-time adaptation in dynamic scenarios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15922–15932, 2023. 1, 2, 3, 4, 6, 7, 12

[54] Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and augmentation. *arXiv preprint arXiv:2110.09506*, 2021.

# ADVMEM: Adversarial Memory Initialization for Realistic Test-Time Adaptation via Tracklet-Based Benchmarking

## Supplementary Material

## 8. ITD Construction Analysis & Experimental Setup Details

### 8.1. Dataset Distribution Overview

The dataset distribution of the 21 object classes, as shown in Figure 5, presents a varied representation in terms of the number of objects and instances across different classes. The distribution is not uniform, with some classes having a higher number of instances compared to others. This variation provides a diverse range of object occurrences, which can be reflective of different scenarios where certain objects appear more frequently while others are less common. Such a distribution can be valuable in assessing model performance across a broad spectrum of object categories, ensuring that both commonly and less commonly occurring objects are adequately represented in the dataset.
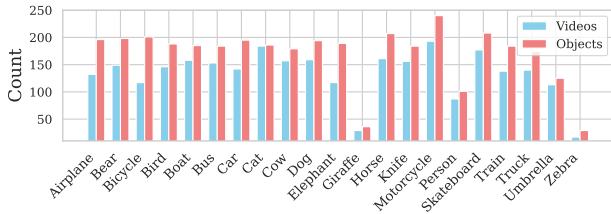


Figure 5. **Class distribution of the test set.** Here we report a detailed breakdown of the distribution of the 21 object classes in terms of the number of objects and instances.

### 8.2. Experiment Setup Overview

Figure 3 provides a visual representation of the two distinct experimental setups used in our study: the frame-wise and tracklet-wise configurations. In the frame-wise setup, a single frame is sampled from each tracklet, ensuring that each object is observed only once. This setup allows for a broad, non-redundant sampling of the dataset, suitable for scenarios where each object instance is considered independently. On the other hand, the tracklet-wise setup processes the frames within each tracklet sequentially, capturing the temporal continuity and variations that occur as the object is observed over time. Both setups are further divided into i.i.d. (independent and identically distributed) and non-i.i.d. settings, providing a comprehensive evaluation framework. The i.i.d. setting assumes that the frames or tracklets are sampled without any dependency, simulating a random observation scenario. Conversely, the non-i.i.d. setting intro-

duces dependencies between samples, reflecting more realistic conditions where observations are not entirely independent, such as in continuous video streams. This dual approach allows for a thorough assessment of model performance under varying assumptions of data distribution and sampling.

### 8.3. Controlling the Label Distribution in Streaming Scenarios

Using a Dirichlet distribution $\text{Dir}(\gamma)$ for sampling labels, as introduced in [53], provides a mechanism for controlling the distribution shift in the label space through the distribution's $\gamma$ parameter.

Specifically, as $\gamma \rightarrow 0$, we converge to a class-incremental non-i.i.d. setup, while, as $\gamma \rightarrow \infty$, we transition towards a uniform i.i.d. setup. We thus extend our experiments by analyzing intermediate stages with $\gamma \in \{10^{-4}, 10^{-1}, 10^3\}$, and report the results in Figure 4a.

For low values of $\gamma$, the model observes a challenging class-incremental non-i.i.d. scenario. In this evolving context, our proposed memory initialization technique proves instrumental in enhancing model performance. The adversarially initialized memory addresses class bias within the incoming image stream and mitigates degradation in performance.

Conversely, as $\gamma$ increases, the scenario becomes easier, as the model has higher chances of encountering images from uniformly-sampled classes. In this scenario, the memory samples that were initialized by ADVMEM are likely to be quickly and completely replaced with reliable samples from the stream. This quick replacement is a result of the stream's extreme uniformity, which causes any initialization to be inconsequential. In this context, ADVMEM neither augments nor diminishes performance.

### 8.4. Sensitivity to Batch Size in Streaming Adaptation

In this section, we explore the impact of batch size on performance. In particular, we vary the batch size in $\{8, 16, 32\}$, and report our results in Figure 4b. As expected, increasing the batch size improves performance, since larger batches provide models with more (and more diverse) examples on which to adapt. When comparing between the i.i.d. and non-i.i.d. setups (*i.e.* $\gamma \rightarrow \infty$ *vs.* $\gamma \rightarrow 0$), we find that methods equipped with ADVMEM exhibit improved or at least sustained performance across batch sizes. However, similar to our observations in previ-

ous sections, the impact of ADVMEM is less pronounced in the i.i.d. setting ($\gamma \to \infty$), indicating that memory initialization has limited effect in such scenarios.

As mentioned in previous section, we use a default batch size of 64. This batch size ensures that exactly one tracket (64 frames) fits in one forward pass. In previous sections, we experiment with different batch sizes, which are smaller than the default. When the batch size is smaller than the number of frames in a tracklet, we adopt a sequential processing approach, where a tracklet is consecutively processed until all frames are observed. This sequential approach ensures that the entire content of a tracklet is leveraged, potentially enabling the model to adapt effectively to the inherent temporal dependencies and patterns within the data.

However, counterintuitively, we can see from Fig. 4 that smaller batch sizes do not lead to a lower error rate in our experiments. While sequential processing allows for a detailed examination of temporal intricacies within a tracklet, the models, influenced by label distribution, exhibit a degradation in performance with smaller batch sizes. This observation underscores the intricate interplay between batch size, sequential information utilization, and the model's robustness to label distribution.

Importantly, our proposed memory initialization (ADVMEM) consistently improves or maintains performance compared to standard memory initialization, regardless of batch size. This fact highlights the robustness and effectiveness of ADVMEM in diverse experimental conditions.

## 9. Dynamic Corruption Incorporation: Analysis and Results

Incorporating dynamic corruptions, as illustrated in Figure 6, into our experiments involves the continuous application of corruptions within the tracklet, where the severity level is defined as a function of time. This dynamic approach enables us to precisely control the severity level of each corruption, closely mimicking real-world scenarios. For example, when simulating defocus blur, the dynamic corruption setting introduces fluctuations in focus, alternating between in and out of focus. Similarly, for weather-related corruptions, such as rain, the dynamic application varies the intensity over time, simulating realistic variations in weather conditions within the video clips. Dynamic corruptions are controlled by the severity function $S(t) = s \cdot |\text{sign}(t)|$, where $s$ represents the severity level and $t$ is the index of the frame in a given tracklet. In contrast to static setups, each frame $k_t$ of the $k$-th tracklet experiences variable severity $S(t)$ at time $t$. The severity function can be customized for each corruption type by adjusting the function's parameters (*i.e.* frequency) or additional factors such as random noise. As depicted in Table 4, the deploy-

ment of our proposed ADVMEM in the context of dynamic corruptions exhibits a consistent trend, akin to our findings from the experiments on static corruptions. This observation underscores that the utilization of ADVMEM consistently again enhances or maintains performance across diverse scenarios.

## 10. Additional Ablations: Vision Transformer (ViT) Experiments

We extend our experiments to include the Vision Transformer (ViT) architecture. ViT outperforms ResNet-18, even at lower batch sizes, due to its reduced sensitivity to batch size [36]. Our ViT experiments focus on the impact of varying batch sizes on method performance (Figure 7). Larger batch sizes do consistently boost performance, with ADVMEM adding further improvements. This analysis highlights the adaptability and effectiveness of ADVMEM.

## 11. Evaluation of Memory Initialization using Training Samples

---
**Algorithm 2** TrainMem
---
    **function** INITIALIZEMEMORY($K, N$)
        Initialize $\mathcal{M} = \{\}$
        **while** $|\mathcal{M}| < N$ **do**
            $y \sim \mathcal{U}\{1, 2, \ldots, K\}$
            $x \sim \mathcal{U}(\mathcal{D}\{x|y\})$
            $\mathcal{M} \leftarrow \mathcal{M} \cup x$
            $\mathcal{D} \leftarrow \mathcal{D} \setminus \{x\}$
        **end while**
        **return** $\mathcal{M}$
    **end function**
---

In contrast to the approach outlined in Algorithm 1, where adversarial samples are employed, Algorithm 2 initializes the memory with images from the training set, denoted by $\mathcal{D}$. Here, $\mathcal{D}$ consists of images and their corresponding labels. For any label $y$, $\mathcal{D}\{x|y\}$ represents the subset of $\mathcal{D}$ corresponding to images $x$ labeled as $y$. The initialization process involves uniformly selecting images from $\mathcal{D}$ without replacement until the memory is full. This procedure, which we refer to as TrainMem, results in a class-wise balanced memory initialization, similar to our adversarial one.

Table 5 summarizes the results for the tracklet-wise non-i.i.d. setup. We observe that while initializing the memory with TrainMem, *i.e.* via Algorithm 2, has positive impact in reducing the error rate of RoTTA, it significantly underperforms our novel ADVMEM. For example, TrainMem reduces the error rate against glass blur by 1.7% when com-

Figure 6. **Dynamic Severity:** We consider temporal variations in the severity of corruptions, reflecting the realistic scenarios where the impact of corruptions may change over time within a single video clip.



Table 4. **Effect of Adversarial Memory Initialization on TTA performance under Dynamic Severity:** We assess the impact of memory initialization techniques on TTA methods in the dynamic severity setting, where the intensity of corruptions varies within the tracklet. The table presents results for standard (✗) and adversarial (✓) memory initialization (ADVMEM).

(a) Tracklet wise *i.i.d.*

| Method | ADVMEM | Noise | | | Weather | | | Avg. |
|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | frost | fog | brigh. | |
| TENT | ✗ | 89.4 | 89.6 | 91.3 | 92.3 | 91.8 | 92.3 | 91.1 |
| | ✓ | 68.5 | 69.0 | 70.6 | 81.2 | 73.4 | 56.1 | 69.8 |
| SHOT-IM | ✗ | 89.7 | 89.2 | 91.2 | 92.1 | 91.9 | 92.2 | 91.0 |
| | ✓ | **38.4** | **38.1** | **45.5** | 53.3 | 40.3 | 25.0 | 40.1 |
| RoTTA | ✗ | 41.1 | 38.1 | 50.4 | **49.4** | **38.3** | **23.0** | **40.0** |
| | ✓ | 40.7 | 38.3 | 50.0 | 49.6 | 39.1 | 23.4 | 40.2 |

(b) Tracklet wise *non*-i.i.d.

| Method | ADVMEM | Noise | | | Weather | | | Avg. |
|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | frost | fog | brigh. | |
| TENT | ✗ | 89.5 | 89.8 | 91.3 | 92.3 | 91.9 | 92.3 | 91.2 |
| | ✓ | 78.1 | 78.1 | 82.7 | 87.3 | 84.7 | 78.8 | 81.6 |
| SHOT-IM | ✗ | 89.8 | 88.9 | 91.2 | 92.3 | 91.9 | 92.1 | 91.0 |
| | ✓ | 89.8 | 89.7 | 91.4 | 91.1 | 90.9 | 89.8 | 90.4 |
| RoTTA | ✗ | 61.2 | 62.1 | 68.8 | 75.1 | 66.9 | 59.4 | 65.6 |
| | ✓ | **58.6** | **60.3** | **67.2** | **74.7** | **66.1** | **51.1** | **63.0** |

pared to RoTTA. However, ADVMEM improves over this naive initialization by over 2% against the same corruption.

## 12. Visualizing Adversarial Examples for AD-VMEM Initialization

In this section, we present visualizations of adversarial examples used for initializing ADVMEM. These examples are generated during the memory bank initialization process. For details on the creation of adversarial examples, please refer to ADVMEM sections. Figure 8 showcases a selection of these adversarial examples.



Figure 7. **Error rate as we vary the batch size.** We study the influence of batch size on method performance. Larger batch sizes enhance performance across the board, with our proposed AD-VMEM consistently contributing to further improvements. All results presented here are for ViT-B-16.
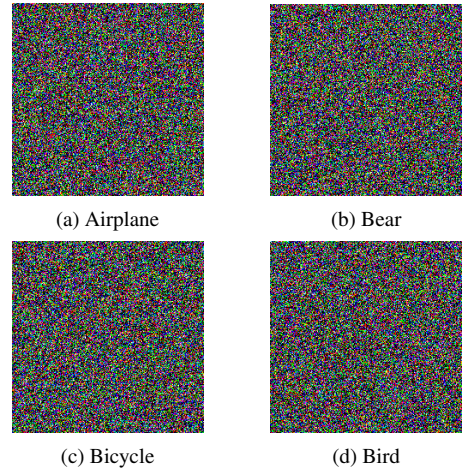


(a) Airplane      (b) Bear

(c) Bicycle      (d) Bird

Figure 8. Visualizations of selected adversarial examples used for initializing the memory bank in ADVMEM.

Table 5. **Tracklet-wise *non*-i.i.d.:** Assessment of TTA methods on ITD in a *non*-i.i.d. tracklet context. We contrast standard memory initialization (✗), memory initialized with training samples (✗), and adversarial memory initialization (✓).

| Method | ADVMEM | Noise | | | Blur | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | defoc. | glass | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| RoTTA | ✗ | 85.6 | **81.6** | 86.5 | 86.9 | 87.3 | 78.9 | 81.7 | 83.5 | 75.0 | 67.2 | 71.8 | 79.7 | 77.1 | 67.4 | 79.3 |
| | ✗ | **83.0** | **81.6** | **83.1** | 86.3 | 85.6 | 77.3 | 84.4 | **79.7** | 75.4 | 71.0 | **69.5** | 78.3 | 71.8 | 67.7 | 78.2 |
| | ✓ | 84.4 | 82.2 | 84.1 | **83.0** | 83.3 | 68.6 | 80.2 | 80.5 | **74.7** | 62.8 | 73.6 | **75.2** | 61.3 | 62.9 | **75.5** |

Table 6. **Performance Comparison under Frame-wise i.i.d. Assumption.** We report average error rates for TTA methods on images from our ITD dataset. In this setup, all images/frames are shuffled and then independently subjected to corruptions. Notably, SHOT-IM outperforms all other methods and across all corruptions, showcasing its robustness against these domain shifts.

| Method | Noise | | | Blur | | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | gauss. | shot | impul. | defoc. | glass | motion | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| Source | 93.8 | 90.8 | 94.0 | 50.8 | 45.1 | 51.2 | 45.2 | 61.7 | 70.4 | 70.6 | 31.4 | 86.6 | 64.0 | 39.1 | 41.1 | 62.4 |
| AdaBN | 57.9 | 56.0 | 58.2 | 50.8 | 51.7 | 44.2 | 37.9 | 53.0 | 55.7 | 50.6 | 32.5 | 50.8 | 37.8 | 29.0 | 36.8 | 46.9 |
| SHOT-IM | **47.5** | **45.5** | **47.1** | **43.6** | **43.2** | **36.4** | **31.2** | **47.1** | **50.2** | **41.7** | **27.8** | **44.5** | **29.6** | **24.6** | **30.0** | **39.3** |
| TENT | 57.9 | 55.8 | 58.4 | 50.6 | 51.7 | 44.2 | 37.8 | 52.9 | 55.6 | 50.4 | 32.5 | 50.4 | 37.7 | 28.9 | 36.8 | 46.8 |
| SAR | 57.6 | 55.3 | 59.1 | 50.4 | 51.3 | 43.7 | 37.6 | 52.8 | 55.4 | 50.1 | 32.2 | 49.4 | 37.6 | 28.6 | 36.3 | 46.5 |
| CoTTA | 57.9 | 55.8 | 58.6 | 50.3 | 51.3 | 43.8 | 37.6 | 53.0 | 55.8 | 50.4 | 32.3 | 50.4 | 37.8 | 28.9 | 36.7 | 46.7 |
| ETA | 57.9 | 54.6 | 59.0 | 50.5 | 51.6 | 44.0 | 37.5 | 53.1 | 55.7 | 50.7 | 32.4 | 50.9 | 37.4 | 28.8 | 36.7 | 46.7 |
| EATA | 57.9 | 55.1 | 58.1 | 50.2 | 51.9 | 43.4 | 38.0 | 54.4 | 56.2 | 50.9 | 31.9 | 50.3 | 38.8 | 28.5 | 36.7 | 46.8 |
| RoTTA | 69.4 | 64.0 | 71.2 | 54.4 | 55.0 | 50.6 | 42.4 | 60.0 | 62.0 | 59.6 | 34.3 | 64.3 | 43.3 | 31.1 | 38.7 | 53.4 |

Table 7. **Tracklet-wise i.i.d. Evaluation with and without Memory.** We report average error rates for TTA methods on our ITD dataset under the tracklet-wise i.i.d. scenario (*i.e.* entire tracklets are sampled i.i.d.). The results are grouped to reflect the presence (✓) or absence (✗) of memory during adaptation. RoTTA shows notable robustness by using memory, significantly outperforming the non-memory variants across various corruption types, as demonstrated by the marked difference in error rate.

| Method | Memory | Noise | | | Blur | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | defoc. | glass | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| TENT | ✗ | 94.3 | 94.3 | 94.4 | 94.5 | 94.6 | 94.1 | 94.1 | 94.5 | 93.4 | 93.1 | 93.7 | 94.0 | 93.5 | 93.7 | 94.0 |
| | ✓ | 94.2 | 94.2 | 94.3 | 94.5 | 94.5 | 93.9 | 93.7 | 93.7 | 93.2 | 93.0 | 92.8 | 93.8 | 93.4 | 93.5 | 93.8 |
| SHOT-IM | ✗ | 94.8 | 94.7 | 94.8 | 94.8 | 94.7 | 94.8 | 94.7 | 94.8 | 94.6 | 94.4 | 94.7 | 94.9 | 94.7 | 94.8 | 94.7 |
| | ✓ | 93.8 | 93.9 | 93.9 | 94.2 | 94.3 | 93.7 | 93.3 | 93.0 | 93.0 | 92.8 | 92.1 | 93.5 | 93.4 | 93.2 | 93.4 |
| RoTTA | ✓ | **68.0** | **62.4** | **68.6** | **58.5** | **56.7** | **40.9** | **56.6** | **60.7** | **52.2** | **30.2** | **60.7** | **39.7** | **27.8** | **35.1** | **51.3** |

Table 8. **Tracklet-wise *non*-i.i.d. Evaluation with and without Memory Banks.** We report the performance of TTA methods when tracklets are non-i.i.d., *i.e.* tracklets are sampled such that their labels display correlation in time (by following a Dirichlet distribution). We further examine whether using a memory bank (✓) influences outcomes. RoTTA with memory exhibits the lowest error rates, indicating proficiency against non-i.i.d. data. Without memory (✗), all methods experience worse error rates, underscoring the impact of memory in adapting to complex data streams.

| Method | Memory | Noise | | | Blur | | | Weather | | | | Digital | | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | gauss. | shot | impul. | defoc. | glass | zoom | snow | frost | fog | brigh. | contr. | elast. | pixel. | jpeg | |
| TENT | ✗ | 94.3 | 94.3 | 94.4 | 94.5 | 94.6 | 94.0 | 94.0 | 95.0 | 93.5 | 93.2 | 93.5 | 94.0 | 93.5 | 93.7 | 94.0 |
| | ✓ | 94.2 | 94.3 | 94.3 | 94.5 | 94.5 | 94.0 | 93.7 | 93.6 | 93.4 | 93.0 | 92.9 | 93.8 | 93.4 | 93.6 | 93.8 |
| SHOT-IM | ✗ | 95.1 | 94.9 | 95.1 | 95.1 | 95.1 | 95.7 | 94.9 | 95.2 | 94.8 | 95.0 | 95.0 | 94.8 | 95.3 | 95.2 | 95.1 |
| | ✓ | 93.9 | 93.8 | 93.9 | 94.3 | 94.3 | 93.7 | 93.4 | 93.6 | 93.2 | 92.7 | 92.5 | 93.6 | 93.5 | 93.2 | 93.6 |
| RoTTA | ✓ | **85.6** | **81.6** | **86.5** | **86.9** | **87.3** | **78.9** | **81.7** | **83.5** | **75.0** | **67.2** | **71.8** | **79.7** | **77.1** | **67.4** | **79.3** |

# 13. Supplementary Table Details

This appendix we present the expanded versions of the tables from the main paper, maintaining the same titles for consistency. These tables contain additional data and detailed results. The supplementary information includes detailed breakdowns and results. These expanded tables are intended to offer a comprehensive reference for readers seeking further insights and details related to the study.