

# Appendix

## NOTE: Robust Continual Test-time Adaptation Against Temporal Correlation

### A Experimental details

For all the experiments in the paper, we used three different random seeds (0, 1, 2) and reported the average errors (and standard deviations). We ran our experiments on NVIDIA GeForce RTX 3090 GPUs.

#### A.1 Baseline details

We referred to the official implementations of the baselines. We use the reported best hyperparameters from their paper or code. We further tuned hyperparameters if there exists a hyperparameter selection guideline. Here, we provide additional details of the baseline implementations, including hyperparameters.

**PL.** Following the previous studies [41, 44], we update the BN layers only in PL. We set the learning rate as  $LR = 0.001$  as the same as [41].

**ONDA.** ONDA [27] has two hyperparameters, the update frequency  $N$  and the decay of the moving average  $m$ . The authors set  $N = 10$  and  $m = 0.1$  as the default values throughout the experiments, and we follow this choice unless specified.

**TENT.** TENT [41] set the learning rate as  $LR = 0.001$  for all datasets except for ImageNet, and we follow this choice. We referred to the official code<sup>1</sup> for implementing TENT.

**LAME.** LAME [4] needs an affinity matrix and has hyperparameters related to it. We follow the authors' hyperparameter selection specified in the paper and their official code. Namely, we use the kNN affinity matrix with the value of  $k$  set as 5. We referred to the official code<sup>2</sup> for implementing LAME.

**CoTTA.** CoTTA [44] has three hyperparameters, augmentation confidence threshold  $p_{th}$ , restoration factor  $p$ , and exponential moving average (EMA) factor  $m$ . We follow the authors' choice for restoration factor ( $p = 0.01$ ) and EMA factor ( $\alpha = 0.999$ ). For the augmentation confidence threshold, the authors provide a guideline to choose it, using 5% quantile for the softmax predictions' confidence on the source domains. We follow this guideline, which results in  $p_{th} = 0.92$  for MNIST-C and CIFAR10-C,  $p_{th} = 0.72$  for CIFAR100-C, and  $p_{th} = 0.55$  for KITTI. For 1D time-series datasets (HARTH and ExtraSensory), the authors do not provide augmentations, and it is non-trivial to select appropriate augmentations for them. We thus do not use augmentations for these datasets. We referred to the official code<sup>3</sup> for implementing CoTTA.

#### A.2 Dataset details

##### A.2.1 Robustness to corruptions

**MNIST-C.** MNIST-C [28] applies 15 corruptions to the MNIST [21] dataset. Specifically, the corruptions include Shot Noise, Impulse Noise, Glass Blur, Motion Blur, Shear, Scale, Rotate, Brightness, Translate, Stripe, Fog, Spatter, Dotted Line, Zigzag, and Canny Edges, as illustrated in Figure 1. Note that the result of this dataset is included only in the supplementary material. In total, MNIST-C has 60,000 clean training data and 150,000 corrupted test data (10,000 for each corruption type). We use ResNet18 [12] as the backbone network. We train it on the clean training

---

<sup>1</sup><https://github.com/DequanWang/tent>

<sup>2</sup><https://github.com/fiveai/LAME>

<sup>3</sup><https://github.com/qinenergy/cotta>

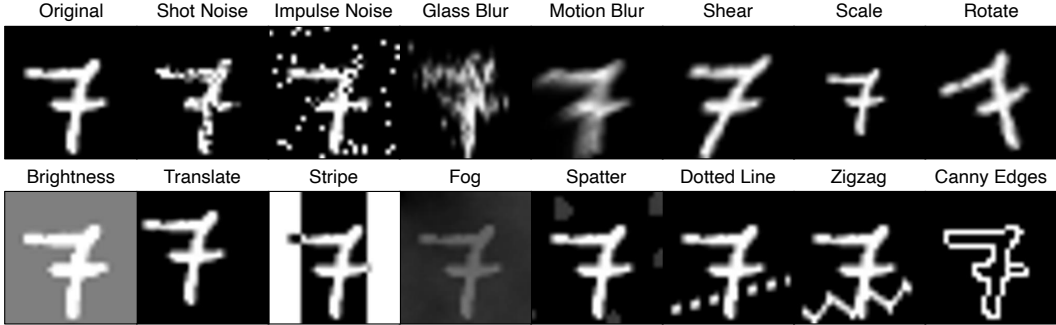


Figure 1: Illustration of the 15 corruption types in the MNIST-C dataset.

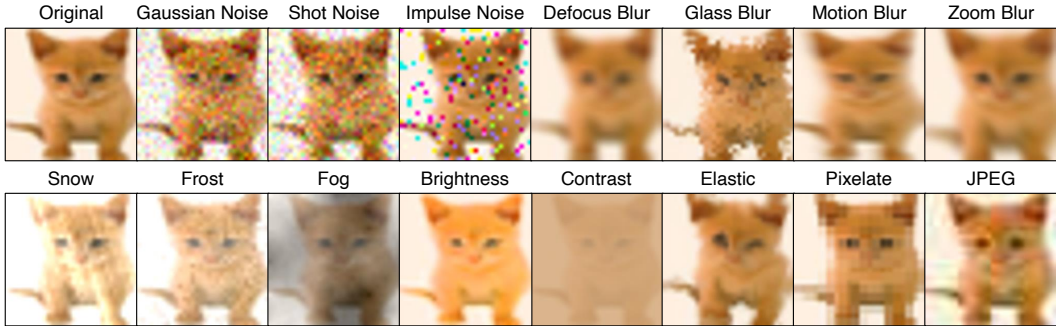


Figure 2: Illustration of the 15 corruption types in the CIFAR10-C/CIFAR100-C/ImageNet-C dataset.

data to generate source models, using stochastic gradient descent with momentum=0.9 and cosine annealing learning rate scheduling [26] for 100 epochs with an initial learning rate of 0.1.

**CIFAR10-C/CIFAR100-C.** CIFAR10-C/CIFAR100-C [13] are common TTA benchmarks for evaluating the robustness to corruptions [29, 33, 41, 44]. Both CIFAR10/CIFAR100 [19] have 50,000/10,000 training/test data. CIFAR10/CIFAR100 have 10/100 classes, respectively. CIFAR10-C/CIFAR100-C apply 15 types of corruptions to CIFAR10/CIFAR100 test data: Gaussian Noise, Shot Noise, Impulse Noise, Defocus Blur, Frosted Glass Blur, Motion Blur, Zoom Blur, Snow, Frost, Fog, Brightness, Contrast, Elastic Transformation, Pixelate, and JPEG Compression, as illustrated in Figure 2. We use the most severe corruption level of 5, similar to previous studies [29, 33, 41, 44]. This results in a total of 150,000 test data for CIFAR10-C/CIFAR100-C, respectively. We use ResNet18 [12] as the backbone network. We train it on the clean training data to generate source models, using stochastic gradient descent with momentum=0.9 and cosine annealing learning rate scheduling [26] for 200 epochs with an initial learning rate of 0.1 and a batch size of 128.

**ImageNet-C.** ImageNet-C is another common TTA benchmark for evaluating the robustness to corruptions [29, 33, 41, 44, 4]. ImageNet [7] has 1,281,167/50,000 training/test data. ImageNet-C applies the same 15 types of corruption used in CIFAR10-C and CIFAR100-C. We use a pre-trained ResNet18 [12] on ImageNet training data and fine-tune it by replacing BN layers with IABN layers on the clean ImageNet training data. For fine-tuning, we use stochastic gradient descent with momentum=0.9 for 30 epochs with a fixed learning rate of 0.001 and a batch size of 256.

**Temporally correlated streams via Dirichlet distribution.** Note that most public vision datasets are not time-series data, and existing TTA studies usually shuffled the order of these datasets resulting in i.i.d. streams, which might be unrealistic in real-world scenarios. To simulate non-i.i.d. streams from these “static” datasets, we utilize Dirichlet distribution that is widely used to simulate non-i.i.d. settings. [23, 15, 43, 42] Specifically, we simulate a non-i.i.d partition for  $T$  tokens on  $C$  classes. For each class  $c$ , we draw a  $T$ -dimensional vector  $\mathbf{q}_c \sim \text{Dir}(\delta \mathbf{p})$ , where  $\text{Dir}(\cdot)$  denotes the Dirichlet distribution,  $\mathbf{p}$  is a prior class distribution over  $T$  classes, and  $\delta > 0$  is a concentration parameter. We assign data from each class to each token  $t$ , following proportion  $\mathbf{q}_c[n]$ . To simulate the nature of

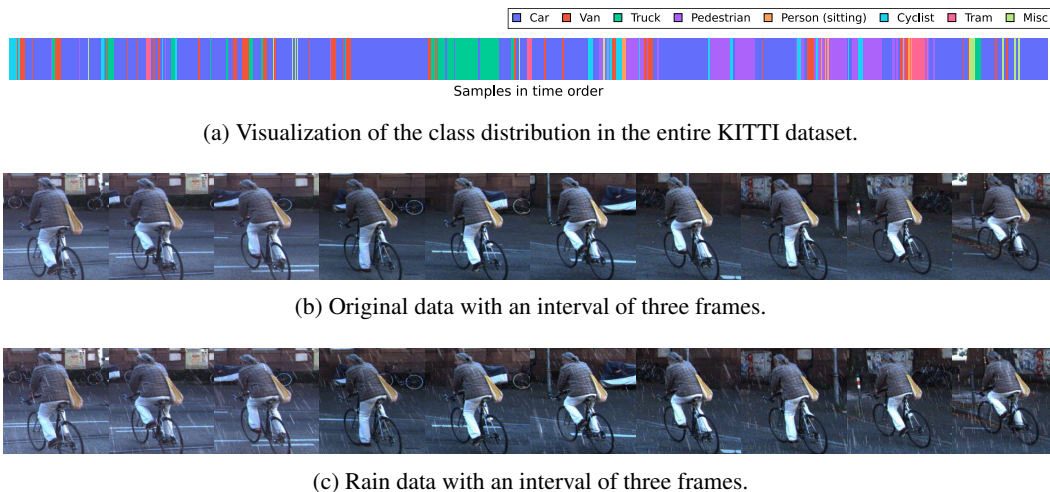


Figure 3: Illustration of the test stream of the KITTI dataset. We apply a 200mm/hr rain intensity to the original data.

real-world online data where sequences are temporally correlated, and data from the same classes appear multiple times (e.g., walking, jogging, and then walking, see Figure 4 and 5 for illustrations), we concatenate the generated  $T$  tokens to create a synthetic non-i.i.d. sequential data. We use  $\delta = 0.1$  as the default value if not specified.

### A.2.2 Real-distributions with domain shift

The following illustrates the summary and preprocessing steps of datasets collected in the real world or have a resemblance to class distributions in the real world.

**KITTI, KITTI-Rain.** KITTI [9] is a well-known dataset used in numerous tasks such as object detection, object tracking, depth estimation, etc. It must be emphasized that the dataset was collected by driving around the city, in rural areas and on highways, which captures the real-world distribution. From the available tasks, we select the object tracking task; to utilize its temporal correlation. In order to reduce the task to a single image classification task, we crop each frame with respect to the largest bounding box. Domain gap is introduced through synthetic generation of corresponding “rainy” frames, hereby denoted as KITTI-Rain [11]. KITTI-Rain is generated via a two-step procedure: (1) generation of a depth-map estimation of each frame, and (2) generation of rainy images from the vanilla frame and its corresponding depth map, as described in [11]. For the depth map generation, we used Monodepth [10], and for rainy image generation, we used the source code available in [11]. The rain intensity is set to 200mm/hr for training and testing. The final source domain consists of 7,481 samples, and each of the target domains consists of 7,800 samples. We use ResNet50 [12] pre-trained on ImageNet [8] as the backbone network. We fine-tune it on the KITTI training data to generate source models, using the Adam optimizer [18] and cosine annealing learning rate scheduling [26] for 50 epochs with an initial learning rate of 0.1 and a batch size of 64.

**HARTH.** Human Activity Recognition Trondheim dataset [25] was collected from 22 users, with two three-axial Axivity AX3 accelerometers, each attached to the subject’s thigh and lower back. HARTH was also collected in a free-living environment and labeled through recorded video. We set the source domain as the accelerometer data collected from the back (15 users), and set the target domain as one collected from the thigh (from the remaining seven users). We deem such a setting to be natural, for one of the most dominant forms of domain shift in wearable sensory data is by the positioning of sensors on the human body [20]. We use a window size of 50 and min-max scaled (0-1) the data, following the original paper [25]. The final source domain consists of 82,544 samples, and each of the seven target domains consists of {S008: 8,140, S018: 6,241, S019: 5,846, S021: 5,910, S022: 6,448, S028: 3,271, S029: 3,521} samples. We use four one-dimensional convolutional layers followed by one fully-connected layer as the backbone network. We train it on the source data to generate source models, using stochastic gradient descent with momentum=0.9 for 100 epochs and

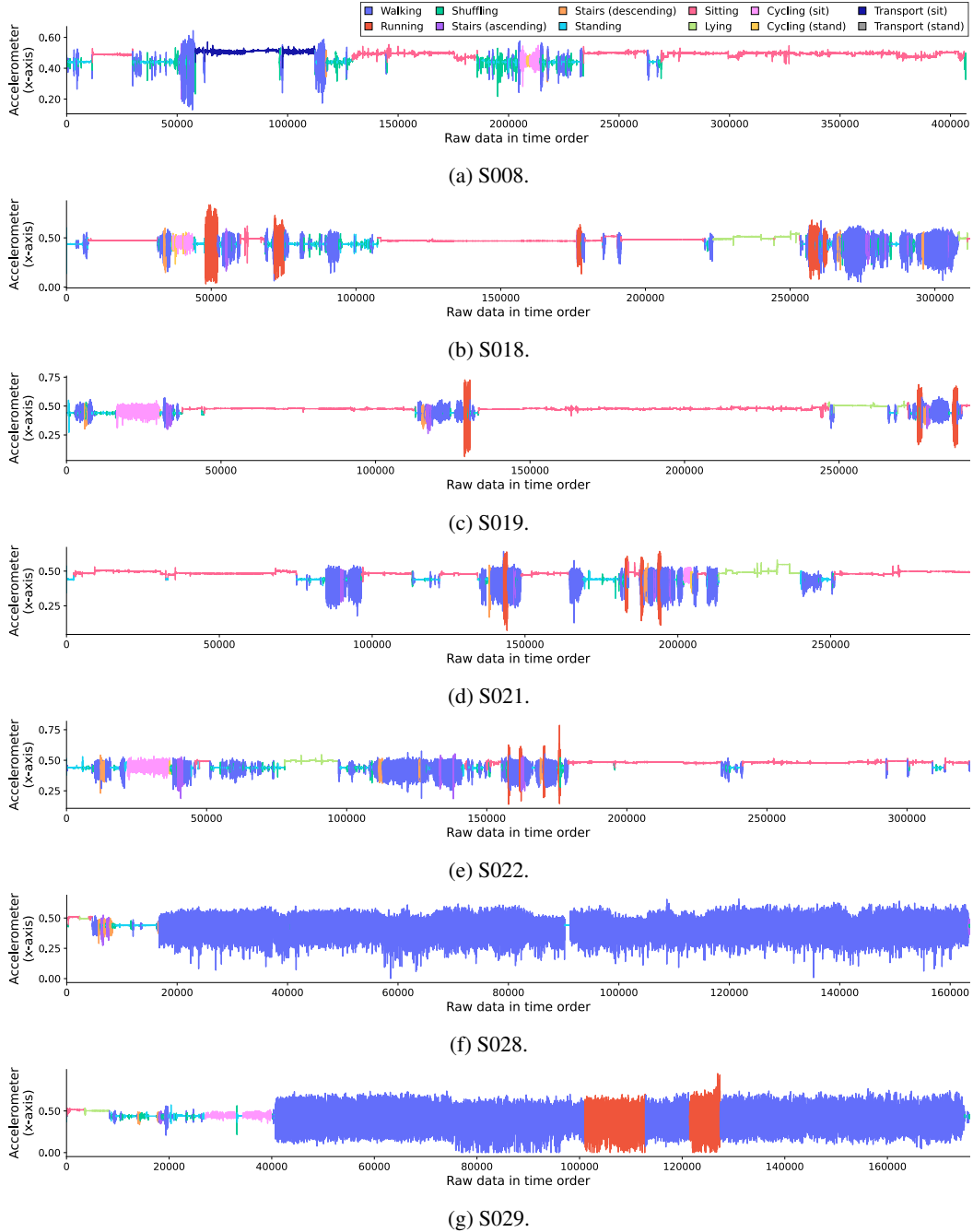


Figure 4: Illustration of the target streams of the HARTH dataset. We specify x-axis accelerometer values only.

cosine annealing learning rate scheduling [26] with an initial learning rate of 0.1 and a batch size of 64.

**ExtraSensory.** Extrasensory dataset [38] was collected from 60 users with the user’s own smartphones over a seven-day period in the wild, i.e., data was collected from users who engaged in their regular natural behavior. As there were no constraints on the subject’s activity, the distribution varied from user to user. We select the five most frequently occurred, mutually exclusive activities (lying down, sitting, walking, standing, running) and omit other labels. We further process the data to only those consisting of the following sensor modalities - accelerometer, gyroscope, magnetometer, and audio. We used a window size of five, with no overlap, and standardly scaled

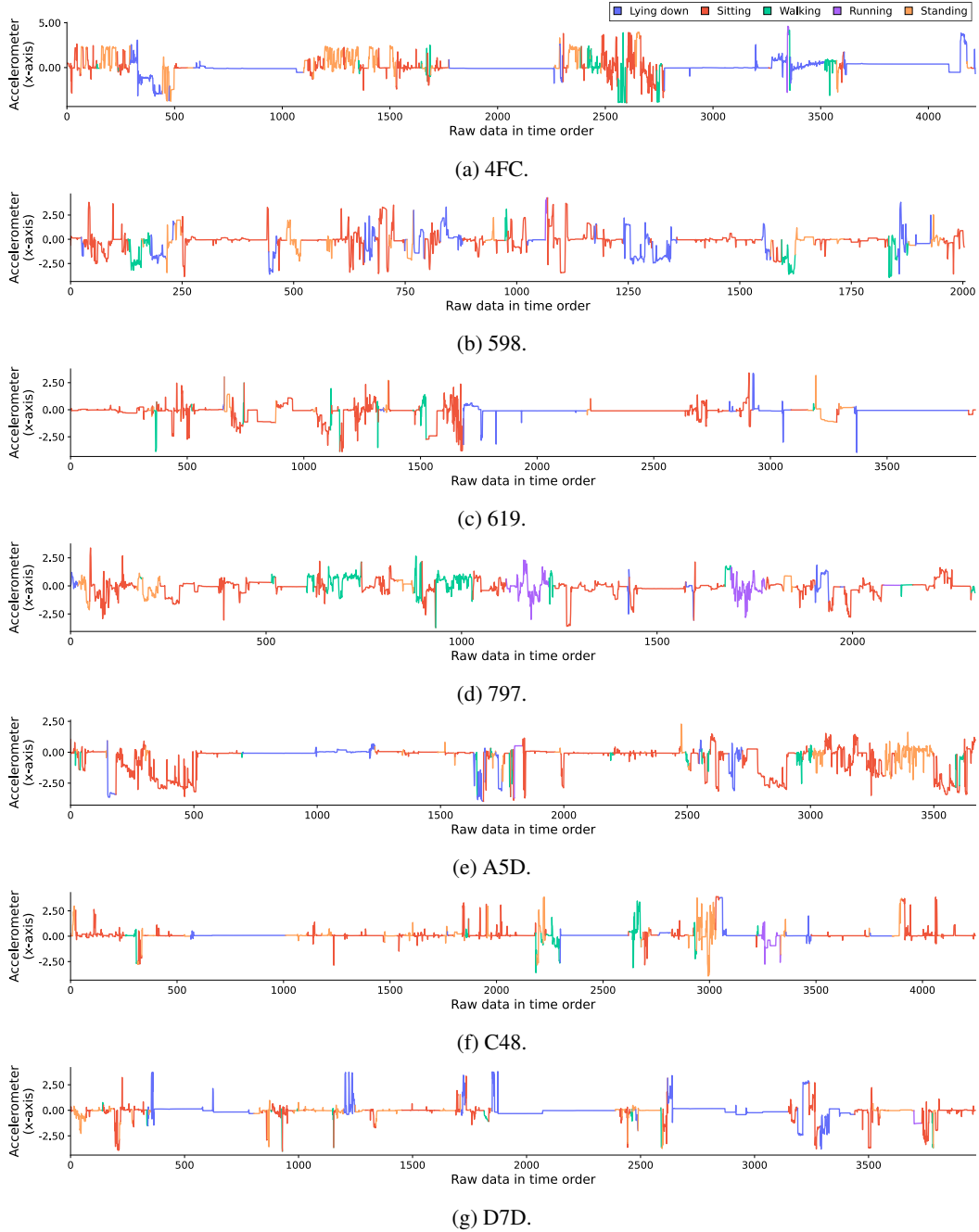
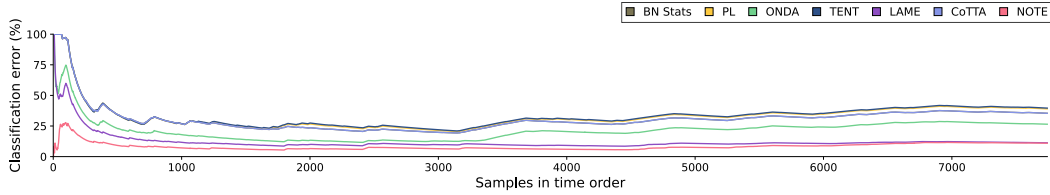


Figure 5: Illustration of the target streams of the Extrasensory dataset. We specify x-axis accelerometer values only. Due to the length of the name of each domain, denoted here with the first three characters.

the datasets. After the pre-processing step, 23 users were left, 16 of them were used as source domains, and seven of them were used as target domains. The final source domain consists of 17,777 samples, and each of the seven target domains consists of {4FC32141-E888-4BFF-8804-12559A491D8C: 844, 59818CD2-24D7-4D32-B133-24C2FE3801E5: 401, 61976C24-1C50-4355-9C49-AAE44A7D09F6: 776, 797D145F-3858-4A7F-A7C2-A4EB721E133C: 463, A5CDF89D-02A2-4EC1-89F8-F534FDABDD96 : 734, C48CE857-A0DD-4DDB-BEA5-3A25449B2153 : 850, D7D20E2E-FC78-405D-B346-DBD3FD8FC92B: 794} samples. We use two one-dimensional convolutional layers followed by one fully-connected layer as the backbone network. We train it on the source data to generate source models, using stochastic gradient descent with momentum=0.9 for



(a) Rain-200.

Figure 6: Illustration of the real-time cumulative classification error change of different methods on the KITTI dataset. The x-axis denotes the samples in order, whereas the y-axis denotes the error rate in percentage. Note that some lines are not clearly visible due to overlap.

100 epochs and cosine annealing learning rate scheduling [26] with an initial learning rate of 0.1 and a batch size of 64.

**Error on the source domain.** We also measure the domain gap between the source and the targets in the three real-distribution datasets: Table 1 for KITTI, Table 2 for HARTH, and Table 3 for Extrasensory. As shown, there is a clear performance degradation from the source domain to the target domain. For HARTH and ExtraSensory, the performance degradation was severe (30~40%p increased error rates compared with Source), indicating the importance of overcoming the domain shift problem in sensory applications.

Table 1: Average classification error (%) and their corresponding standard deviations on the KITTI dataset of the source model. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	Src domain	Rain	Avg
Source	<b>7.4 ± 1.0</b>	12.3 ± 2.3	9.9

Table 2: Average classification error (%) and their corresponding standard deviations on the HARTH dataset of the source model. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	Src domain	S008	S018	S019	S021	S022	S028	S029	Avg
Source	<b>11.7 ± 0.7</b>	86.2 ± 1.3	44.7 ± 2.1	50.4 ± 9.5	74.8 ± 3.8	72.0 ± 2.6	53.0 ± 24.0	57.0 ± 16.7	56.2

Table 3: Average classification error (%) and their corresponding standard deviations on the ExtraSensory dataset of the source model. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	Src domain	4FC	598	619	797	A5C	C48	D7D	Avg
Source	<b>8.3 ± 0.7</b>	34.6 ± 2.5	40.1 ± 0.7	63.8 ± 5.7	45.3 ± 2.4	64.6 ± 3.7	39.6 ± 6.8	63.0 ± 3.9	44.9



## B Domain-wise results

### B.1 Robustness to corruptions

Table 4: Average classification error (%) and their corresponding standard deviations on CIFAR10-C with **temporally correlated test streams**, shown per corruption. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	74.0	66.8	75.3	43.3	48.0	32.6	35.2	22.0	33.0	25.9	8.5	66.1	23.4	53.6	26.8	42.3
BN Stats [29]	$\pm 3.3$	$\pm 3.5$	$\pm 4.2$	$\pm 2.7$	$\pm 2.7$	$\pm 1.2$	$\pm 2.6$	$\pm 0.4$	$\pm 2.5$	$\pm 0.9$	$\pm 0.3$	$\pm 1.8$	$\pm 0.7$	$\pm 0.7$	$\pm 0.7$	73.4
ONDA [27]	<b>77.2</b>	<b>76.7</b>	<b>78.9</b>	<b>70.0</b>	<b>78.6</b>	<b>70.5</b>	<b>71.1</b>	<b>72.5</b>	<b>71.9</b>	<b>70.6</b>	<b>68.7</b>	<b>69.1</b>	<b>75.1</b>	<b>73.6</b>	<b>76.8</b>	63.6
PL [22]	$\pm 0.7$	$\pm 1.0$	$\pm 0.8$	$\pm 1.7$	$\pm 0.6$	$\pm 1.5$	$\pm 1.4$	$\pm 1.4$	$\pm 1.1$	$\pm 1.6$	$\pm 1.9$	$\pm 1.9$	$\pm 1.5$	$\pm 1.4$	$\pm 1.4$	75.4
TENT [41]	69.3	<b>68.5</b>	71.8	<b>58.5</b>	71.0	59.9	59.5	62.4	62.1	59.6	55.6	58.4	65.6	63.9	67.6	76.4
LAME [4]	$\pm 1.0$	$\pm 1.0$	$\pm 0.6$	$\pm 1.4$	$\pm 0.2$	$\pm 1.0$	$\pm 1.0$	$\pm 1.4$	$\pm 1.0$	$\pm 1.3$	$\pm 1.4$	$\pm 1.4$	$\pm 1.0$	$\pm 1.4$	$\pm 1.1$	36.2
CoTTA [44]	<b>78.3</b>	<b>78.0</b>	<b>80.4</b>	<b>72.2</b>	80.1	72.4	73.1	74.5	73.9	73.4	71.5	71.7	77.3	75.7	78.6	75.5
NOTE	$\pm 1.0$	$\pm 1.5$	$\pm 1.0$	$\pm 1.6$	$\pm 1.2$	$\pm 2.2$	$\pm 1.4$	$\pm 2.5$	$\pm 1.8$	$\pm 1.7$	$\pm 2.7$	$\pm 2.5$	$\pm 2.1$	$\pm 1.5$	$\pm 2.7$	21.1
	79.0	<b>78.8</b>	<b>80.6</b>	<b>73.3</b>	<b>80.5</b>	74.4	74.5	74.8	75.0	74.0	72.3	74.9	78.2	76.5	79.0	
	$\pm 2.9$	$\pm 2.8$	$\pm 2.2$	$\pm 1.7$	$\pm 2.9$	$\pm 2.4$	$\pm 3.3$	$\pm 2.2$	$\pm 2.3$	$\pm 2.2$	$\pm 3.4$	$\pm 3.2$	$\pm 2.8$	$\pm 2.9$	$\pm 2.9$	
	73.6	64.8	74.8	36.2	37.7	24.9	27.9	<b>12.4</b>	22.4	19.4	<b>3.6</b>	65.1	<b>12.6</b>	50.3	<b>16.4</b>	
	$\pm 5.2$	$\pm 4.6$	$\pm 6.4$	$\pm 4.4$	$\pm 5.3$	$\pm 1.6$	$\pm 3.4$	$\pm 1.0$	$\pm 3.9$	$\pm 0.9$	$\pm 0.3$	$\pm 1.5$	$\pm 0.8$	$\pm 0.9$	$\pm 1.2$	
	<b>77.0</b>	<b>76.8</b>	<b>79.0</b>	<b>74.1</b>	<b>79.6</b>	<b>74.3</b>	<b>74.0</b>	<b>74.8</b>	<b>73.3</b>	<b>72.9</b>	<b>72.2</b>	<b>76.5</b>	<b>76.5</b>	<b>75.1</b>	<b>76.6</b>	
	$\pm 0.7$	$\pm 0.6$	$\pm 0.7$	$\pm 0.9$	$\pm 0.6$	$\pm 0.5$	$\pm 0.8$	$\pm 1.1$	$\pm 0.9$	$\pm 0.5$	$\pm 0.9$	$\pm 0.8$	$\pm 0.9$	$\pm 0.8$	$\pm 0.6$	
	<b>34.9</b>	<b>32.3</b>	<b>39.6</b>	<b>13.6</b>	<b>35.8</b>	<b>11.8</b>	<b>14.5</b>	14.1	<b>15.2</b>	<b>14.2</b>	7.7	<b>7.6</b>	20.8	<b>27.7</b>	26.4	
	$\pm 1.6$	$\pm 3.1$	$\pm 2.5$	$\pm 0.5$	$\pm 1.9$	$\pm 0.8$	$\pm 0.5$	$\pm 0.6$	$\pm 1.3$	$\pm 0.6$	$\pm 0.3$	$\pm 0.6$	$\pm 0.7$	$\pm 2.6$	$\pm 0.5$	

Table 5: Average classification error (%) and their corresponding standard deviations on CIFAR100-C with **temporally correlated test streams**, shown per corruption. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	88.1	86.8	93.7	64.9	79.7	55.5	57.7	53.8	66.3	59.3	33.0	81.4	49.2	73.6	55.5	66.6
BN Stats [29]	$\pm 0.2$	$\pm 0.6$	$\pm 0.6$	$\pm 0.4$	$\pm 0.9$	$\pm 0.3$	$\pm 0.2$	$\pm 0.4$	$\pm 0.8$	$\pm 0.4$	$\pm 0.3$	$\pm 0.4$	$\pm 0.4$	$\pm 1.1$	$\pm 0.3$	65.0
ONDA [27]	73.9	73.5	77.2	56.9	72.3	<b>58.8</b>	<b>57.9</b>	<b>65.3</b>	65.0	<b>62.4</b>	<b>55.6</b>	57.6	<b>64.6</b>	63.6	<b>71.0</b>	49.6
PL [22]	$\pm 0.5$	$\pm 0.4$	$\pm 0.7$	$\pm 0.2$	$\pm 0.5$	$\pm 0.3$	$\pm 0.4$	$\pm 0.4$	$\pm 0.4$	$\pm 0.6$	$\pm 0.2$	$\pm 0.4$	$\pm 0.5$	$\pm 0.3$	$\pm 0.4$	66.4
TENT [41]	<b>63.0</b>	<b>62.5</b>	<b>68.0</b>	37.3	<b>60.0</b>	40.0	38.3	49.6	50.0	45.2	<b>35.7</b>	40.9	48.6	<b>46.9</b>	<b>57.5</b>	66.9
LAME [4]	$\pm 0.7$	$\pm 0.4$	$\pm 0.5$	$\pm 0.2$	$\pm 0.2$	$\pm 0.3$	$\pm 0.1$	$\pm 0.3$	$\pm 0.6$	$\pm 0.6$	$\pm 0.2$	$\pm 0.5$	$\pm 0.5$	$\pm 0.3$	$\pm 0.2$	63.3
CoTTA [44]	71.9	72.0	76.3	59.3	73.8	<b>61.5</b>	<b>59.9</b>	<b>67.1</b>	<b>66.7</b>	<b>63.0</b>	<b>57.9</b>	62.2	<b>67.6</b>	65.2	71.1	64.2
NOTE	$\pm 1.4$	$\pm 0.5$	$\pm 0.7$	$\pm 0.8$	$\pm 0.9$	$\pm 0.9$	$\pm 0.5$	$\pm 0.9$	$\pm 1.4$	$\pm 1.0$	$\pm 0.5$	$\pm 1.5$	$\pm 1.0$	$\pm 0.3$	$\pm 0.5$	47.0
	71.8	71.0	76.4	60.2	75.0	<b>61.9</b>	<b>60.2</b>	<b>67.8</b>	<b>67.8</b>	<b>63.3</b>	<b>58.4</b>	65.0	68.4	<b>65.0</b>	<b>71.8</b>	
	$\pm 0.9$	$\pm 0.4$	$\pm 1.2$	$\pm 0.6$	$\pm 1.0$	$\pm 0.9$	$\pm 0.7$	$\pm 0.5$	$\pm 0.7$	$\pm 1.1$	$\pm 0.7$	$\pm 1.8$	$\pm 0.9$	$\pm 0.2$	$\pm 0.1$	
	<b>89.0</b>	<b>87.1</b>	<b>94.5</b>	62.3	79.7	49.4	52.8	46.6	63.9	55.6	<b>25.2</b>	<b>82.4</b>	<b>40.8</b>	71.9	<b>47.8</b>	
	$\pm 1.1$	$\pm 0.8$	$\pm 0.7$	$\pm 1.2$	$\pm 1.2$	$\pm 1.0$	$\pm 0.3$	$\pm 0.4$	$\pm 1.9$	$\pm 1.2$	$\pm 0.6$	$\pm 0.2$	$\pm 0.5$	$\pm 1.4$	$\pm 0.7$	
	68.6	67.9	71.4	60.7	69.9	<b>60.8</b>	<b>60.2</b>	<b>64.0</b>	62.9	<b>63.2</b>	<b>56.7</b>	65.6	<b>64.5</b>	60.9	<b>65.3</b>	
	$\pm 0.3$	$\pm 0.4$	$\pm 0.4$	$\pm 0.4$	$\pm 0.4$	$\pm 0.5$	$\pm 0.2$	$\pm 0.3$	$\pm 0.5$	$\pm 0.6$	$\pm 0.2$	$\pm 0.3$	$\pm 0.3$	$\pm 0.0$	$\pm 0.1$	
	66.2	64.2	72.6	<b>37.2</b>	61.1	<b>35.4</b>	<b>37.4</b>	<b>40.0</b>	<b>42.5</b>	<b>43.4</b>	29.4	<b>32.1</b>	44.3	47.5	51.3	
	$\pm 0.8$	$\pm 1.6$	$\pm 0.4$	$\pm 0.8$	$\pm 0.7$	$\pm 0.3$	$\pm 0.4$	$\pm 0.4$	$\pm 0.3$	$\pm 0.5$	$\pm 0.1$	$\pm 0.5$	$\pm 0.4$	$\pm 0.6$	$\pm 0.3$	

Table 6: Average classification error (%) and their corresponding standard deviations on ImageNet-C with **temporally correlated test streams**, shown per corruption. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	<i>Gaussian</i>	<i>Shot</i>	<i>Impulse</i>	<i>Defocus</i>	<i>Glass</i>	<i>Motion</i>	<i>Zoom</i>	<i>Show</i>	<i>Frost</i>	<i>Fog</i>	<i>Brightness</i>	<i>Contrast</i>	<i>Elastic</i>	<i>Pixelate</i>	<i>JPEG</i>	Avg
Source	98.4	97.7	98.4	90.6	92.5	89.8	81.8	89.5	85.0	86.4	51.1	97.2	85.3	76.9	71.7	86.1
BN Stats	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0
ONDA	98.3	<b>98.1</b>	98.4	<b>98.7</b>	<b>98.8</b>	<b>97.8</b>	<b>96.6</b>	<b>96.2</b>	<b>96.0</b>	<b>95.1</b>	<b>93.1</b>	<b>98.6</b>	<b>96.3</b>	<b>95.6</b>	<b>96.1</b>	<b>96.9</b>
PL	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0
TENT	95.1	94.7	95.0	<b>96.2</b>	<b>96.1</b>	<b>92.5</b>	<b>87.2</b>	87.4	<b>87.8</b>	<b>82.7</b>	<b>71.0</b>	<b>96.4</b>	84.9	<b>81.7</b>	<b>86.1</b>	<b>89.0</b>
CoTTA	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0
NOTE	<b>99.3</b>	<b>99.3</b>	<b>99.4</b>	<b>99.5</b>	<b>99.4</b>	<b>99.5</b>	<b>98.8</b>	<b>99.1</b>	<b>99.2</b>	<b>98.1</b>	<b>97.3</b>	<b>99.8</b>	<b>98.4</b>	<b>98.5</b>	<b>98.5</b>	<b>98.9</b>
LAME	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.1	±0.0	±0.0	±0.0	±0.0	±0.0
CoTTA	98.3	<b>98.1</b>	98.4	<b>98.7</b>	<b>98.8</b>	<b>97.8</b>	<b>96.6</b>	<b>96.2</b>	<b>96.0</b>	<b>95.1</b>	<b>93.1</b>	<b>98.6</b>	<b>96.3</b>	<b>95.6</b>	<b>96.1</b>	<b>96.9</b>
NOTE	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0
CoTTA	98.1	97.1	98.0	<b>87.9</b>	<b>90.9</b>	<b>87.1</b>	<b>78.3</b>	87.1	80.2	81.5	<b>39.8</b>	96.4	82.5	70.7	<b>64.9</b>	82.7
NOTE	±0.0	±0.0	±0.0	± <b>0.0</b>	± <b>0.0</b>	± <b>0.0</b>	± <b>0.0</b>	±0.0	±0.0	±0.0	± <b>0.0</b>	±0.0	±0.0	±0.0	± <b>0.0</b>	±0.0
CoTTA	98.2	<b>98.1</b>	98.3	<b>98.8</b>	<b>98.8</b>	<b>97.7</b>	<b>96.8</b>	<b>96.6</b>	<b>96.3</b>	<b>95.3</b>	<b>93.5</b>	<b>98.8</b>	<b>96.5</b>	<b>95.6</b>	<b>96.2</b>	<b>97.0</b>
NOTE	±0.0	± <b>0.0</b>	±0.0	±0.0	±0.0	±0.0	±0.0	±0.1	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0	±0.0
NOTE	<b>94.7</b>	<b>93.7</b>	<b>94.5</b>	<b>91.2</b>	91.0	<b>83.3</b>	79.0	<b>79.0</b>	<b>78.7</b>	<b>66.3</b>	48.0	<b>94.1</b>	<b>76.9</b>	<b>62.6</b>	<b>76.6</b>	<b>80.6</b>
NOTE	± <b>0.1</b>	± <b>0.3</b>	± <b>0.1</b>	± <b>0.1</b>	±0.2	± <b>0.1</b>	±0.2	± <b>0.4</b>	± <b>0.3</b>	± <b>0.6</b>	±0.4	± <b>0.1</b>	± <b>0.6</b>	± <b>0.7</b>	± <b>0.6</b>	±0.6

Table 7: Average classification error (%) and their corresponding standard deviations on MNIST-C with **temporally correlated test streams**, shown per corruption. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	<i>Shot</i>	<i>Impulse</i>	<i>Glass</i>	<i>Motion</i>	<i>Shear</i>	<i>Scale</i>	<i>Rotate</i>	<i>Brightness</i>	<i>Translate</i>	<i>Stripe</i>	<i>Fog</i>	<i>Spatter</i>	<i>Dotted line</i>	<i>Zigzag</i>	<i>Canny edges</i>	Avg
Source	3.7	27.3	20.4	4.6	2.2	5.1	6.5	21.1	13.8	17.4	66.6	3.8	3.7	18.2	26.4	16.1
BN Stats [29]	±0.7	±5.5	±6.4	±0.5	±0.5	±1.0	±1.0	±22.9	±1.4	±17.0	±14.7	±0.4	±0.4	±3.0	±11.4	±11.4
ONDA [27]	<b>72.0</b>	<b>75.2</b>	<b>73.7</b>	<b>72.1</b>	<b>71.2</b>	<b>71.4</b>	<b>71.2</b>	<b>71.6</b>	<b>78.5</b>	<b>72.3</b>	<b>70.8</b>	<b>71.6</b>	<b>73.8</b>	<b>74.6</b>	<b>72.3</b>	<b>72.8</b>
PL [22]	±0.6	±0.8	±1.0	±0.8	±1.1	±0.6	±0.3	±0.6	±0.2	±1.2	±1.2	±0.9	±0.7	±0.6	±0.3	±0.3
TENT [41]	53.3	59.9	59.2	54.1	51.6	53.9	54.6	50.5	65.2	57.5	54.8	54.2	55.4	61.0	56.7	56.1
CoTTA [44]	±3.0	±3.0	±3.3	±3.5	±2.2	±2.5	±2.0	±2.3	±2.1	±0.7	±2.9	±3.0	±2.8	±2.2	±2.1	±2.1
LAME [4]	73.7	76.4	75.3	74.7	72.7	73.3	73.7	73.7	78.7	74.1	75.8	72.5	75.8	76.9	74.5	74.8
NOTE	±1.0	±0.4	±0.5	±1.1	±0.9	±1.6	±0.9	±1.0	±0.3	±1.4	±2.6	±0.8	±0.6	±1.4	±0.1	±0.1
CoTTA [44]	74.7	78.1	76.6	76.1	75.8	73.7	75.2	75.4	78.9	76.7	81.4	73.9	77.3	79.2	75.8	76.6
NOTE	±1.1	±0.9	±0.6	±0.7	±1.1	±1.3	±1.1	±0.3	±0.2	±1.8	±1.7	±0.5	±0.7	±2.0	±1.0	±1.0
CoTTA [44]	<b>1.1</b>	17.0	<b>12.5</b>	<b>1.1</b>	<b>0.4</b>	<b>1.5</b>	<b>2.3</b>	17.2	<b>6.0</b>	<b>12.3</b>	68.3	<b>0.7</b>	<b>0.7</b>	13.2	22.1	11.8
NOTE	± <b>0.3</b>	±8.7	± <b>6.5</b>	± <b>0.3</b>	± <b>0.2</b>	± <b>0.6</b>	± <b>0.6</b>	±26.0	± <b>2.3</b>	± <b>17.2</b>	±15.8	± <b>0.3</b>	± <b>0.4</b>	±3.4	±12.3	±12.3
CoTTA [44]	76.9	79.4	79.1	77.6	75.4	76.2	77.6	76.0	81.6	76.8	78.0	77.6	79.3	80.6	77.6	78.0
NOTE	±0.5	±0.4	±0.5	±0.6	±0.4	±1.3	±0.2	±0.5	±0.9	±0.6	±0.4	±0.6	±0.4	±1.0	±0.5	±0.5
NOTE	<b>3.9</b>	<b>13.8</b>	14.3	3.3	1.7	3.8	6.5	<b>0.9</b>	8.0	14.4	<b>1.6</b>	3.9	4.5	<b>12.6</b>	<b>13.4</b>	<b>7.1</b>
NOTE	± <b>1.3</b>	± <b>2.4</b>	±1.5	±2.4	±0.2	±0.7	±0.3	± <b>0.0</b>	±1.2	±8.1	± <b>0.3</b>	±0.4	±1.2	± <b>2.5</b>	± <b>3.9</b>	±0.6



Table 8: Average classification error (%) and their corresponding standard deviations on CIFAR10-C with **uniformly distributed test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. NOTE\* indicates NOTE used directly with test batches (without using PBRS). Averaged over three runs.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	74.0	66.8	75.3	43.3	48.0	32.6	35.2	22.0	33.0	25.9	8.5	66.1	23.4	53.6	26.8	42.3
	± 3.3	± 3.5	± 4.2	± 2.7	± 2.7	± 1.2	± 2.6	± 0.4	± 2.5	± 0.9	± 0.3	± 1.8	± 0.7	± 0.7	± 0.7	
BN Stats [29]	33.1	31.1	39.8	12.3	34.8	13.7	12.6	18.3	19.9	14.5	<b>9.3</b>	13.0	23.3	20.8	<b>28.0</b>	21.6
	± 0.9	± 1.0	± 0.9	± 0.4	± 0.3	± 0.3	± 0.4	± 0.7	± 0.6	± 0.6	± 0.3	± 0.3	± 0.3	± 0.2	± 0.6	
ONDA [27]	33.4	31.3	40.0	12.3	34.6	13.7	12.4	18.3	19.8	14.3	<b>9.1</b>	14.0	23.3	20.9	<b>28.0</b>	21.7
	± 0.6	± 0.9	± 1.1	± 0.4	± 0.7	± 0.3	± 0.5	± 0.6	± 0.8	± 0.4	± 0.0	± 0.2	± 0.4	± 0.2	± 0.7	
PL [22]	29.4	26.3	36.8	13.7	36.5	14.0	13.5	19.7	21.2	15.6	<b>10.0</b>	14.8	<b>24.5</b>	20.1	<b>27.4</b>	21.6
	± 1.1	± 1.0	± 1.6	± 0.4	± 1.1	± 1.0	± 0.2	± 0.8	± 0.6	± 1.5	± 0.6	± 0.2	± 2.0	± 0.9	± 1.3	
TENT [41]	25.3	23.1	32.1	<b>11.7</b>	33.1	13.2	<b>11.2</b>	15.9	18.8	<b>12.9</b>	<b>8.6</b>	14.4	21.7	<b>16.5</b>	23.6	18.8
	± 0.8	± 1.1	± 1.2	± <b>0.6</b>	± 3.0	± 1.1	± <b>0.1</b>	± 0.3	± 0.7	± <b>0.8</b>	± <b>0.3</b>	± 0.6	± 0.9	± <b>0.8</b>	± 0.7	
LAME [4]	78.2	<b>70.6</b>	<b>80.5</b>	<b>46.6</b>	48.0	<b>34.2</b>	<b>37.4</b>	20.8	30.5	<b>26.9</b>	<b>9.8</b>	<b>71.9</b>	<b>24.2</b>	<b>56.4</b>	25.8	<b>44.1</b>
	± 3.6	± <b>4.0</b>	± <b>4.5</b>	± <b>1.9</b>	± 3.8	± <b>0.4</b>	± <b>1.5</b>	± 0.8	± 4.1	± <b>1.8</b>	± <b>0.2</b>	± <b>1.0</b>	± <b>0.9</b>	± <b>0.8</b>	± <b>0.9</b>	
CoTTA [44]	<b>23.1</b>	<b>21.5</b>	<b>28.0</b>	<b>11.7</b>	<b>29.2</b>	13.3	12.0	16.6	16.6	13.8	<b>8.8</b>	14.9	<b>20.6</b>	17.3	<b>19.9</b>	17.8
	± <b>0.7</b>	± <b>0.6</b>	± <b>0.3</b>	± <b>0.5</b>	± <b>0.6</b>	± 0.6	± 0.5	± 0.2	± 0.3	± 0.4	± <b>0.2</b>	± 0.5	± <b>0.7</b>	± 0.5	± <b>0.4</b>	
NOTE	33.5	30.0	38.2	12.6	34.4	<b>11.5</b>	12.9	<b>14.1</b>	15.2	14.0	<b>7.4</b>	7.8	20.7	24.7	24.2	20.1
	± 1.7	± 1.6	± 0.9	± 0.8	± 0.8	± <b>0.5</b>	± 0.6	± <b>0.2</b>	± 0.8	± 0.6	± <b>0.2</b>	± 0.2	± 0.3	± 0.7	± 0.4	
NOTE*	23.8	23.0	31.1	11.8	30.9	11.8	11.9	15.3	<b>14.0</b>	13.3	<b>8.6</b>	<b>7.5</b>	21.2	16.9	23.0	<b>17.6</b>
	± 0.7	± 0.9	± 0.3	± 0.6	± 1.3	± 0.4	± 0.7	± 1.3	± <b>0.7</b>	± 0.7	± <b>0.2</b>	± <b>0.3</b>	± 0.3	± 0.6	± 1.2	

Table 9: Average classification error (%) and their corresponding standard deviations on CIFAR100-C with **uniformly distributed test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs. NOTE\* indicates NOTE used directly with test batches (without using PBRS)

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	88.1	86.8	93.7	64.9	79.7	55.5	57.7	53.8	66.3	59.3	33.0	81.4	49.2	73.6	55.5	66.6
	± 0.2	± 0.6	± 0.6	± 0.4	± 0.9	± 0.3	± 0.2	± 0.4	± 0.8	± 0.4	± 0.3	± 0.4	± 0.4	± 1.1	± 0.3	
BN Stats [29]	60.9	59.9	65.7	33.7	57.6	36.5	35.2	46.7	46.9	42.8	<b>32.3</b>	35.6	45.8	43.6	55.5	46.6
	± 0.8	± 0.6	± 0.8	± 0.4	± 0.4	± 0.2	± 0.4	± 0.3	± 0.4	± 0.7	± <b>0.4</b>	± 0.5	± 0.3	± 0.3	± 0.2	
ONDA [27]	60.8	60.2	66.0	33.9	57.5	36.3	34.6	46.5	47.2	42.1	<b>32.1</b>	36.4	45.5	43.4	55.1	46.5
	± 0.9	± 0.5	± 0.6	± 0.4	± 0.4	± 0.4	± 0.4	± 0.3	± 0.3	± 0.6	± <b>0.5</b>	± 0.4	± 0.1	± 0.8	± 0.1	
PL [22]	52.2	50.3	59.4	33.5	54.0	35.7	33.1	42.8	44.5	39.2	30.9	35.5	45.5	39.9	50.4	43.1
	± 0.9	± 1.0	± 0.9	± 0.5	± 0.6	± 0.3	± 0.5	± 0.9	± 1.6	± 1.3	± 0.2	± 0.2	± 1.0	± 0.3	± 1.3	
TENT [41]	<b>48.7</b>	<b>47.2</b>	<b>55.6</b>	<b>31.5</b>	<b>50.9</b>	33.5	<b>31.7</b>	39.6	41.0	36.8	29.4	33.6	<b>42.3</b>	<b>36.8</b>	<b>46.4</b>	<b>40.3</b>
	± <b>0.8</b>	± <b>0.6</b>	± <b>0.9</b>	± <b>0.2</b>	± <b>0.5</b>	± 0.4	± <b>0.2</b>	± 0.3	± 0.1	± 0.7	± 0.3	± 0.4	± <b>0.6</b>	± <b>0.5</b>	± <b>0.5</b>	
LAME [4]	<b>91.0</b>	89.5	<b>95.2</b>	<b>68.1</b>	<b>82.7</b>	<b>57.1</b>	<b>60.2</b>	<b>54.7</b>	<b>68.9</b>	<b>61.8</b>	<b>33.7</b>	<b>85.2</b>	<b>50.3</b>	<b>76.7</b>	<b>56.2</b>	<b>68.8</b>
	± <b>1.0</b>	± 1.0	± <b>0.7</b>	± <b>0.9</b>	± <b>1.1</b>	± <b>0.5</b>	± <b>0.3</b>	± <b>0.3</b>	± <b>1.2</b>	± <b>0.6</b>	± <b>0.5</b>	± <b>0.4</b>	± <b>0.2</b>	± <b>1.3</b>	± <b>0.5</b>	
CoTTA [44]	52.8	51.0	56.9	35.8	53.9	37.9	36.8	45.2	44.5	44.0	32.2	41.3	46.1	39.7	46.9	44.3
	± 0.7	± 0.4	± 0.6	± 0.4	± 0.2	± 0.5	± 0.1	± 0.5	± 0.1	± 0.2	± 0.5	± 1.4	± 0.1	± 0.3	± 0.7	
NOTE	65.6	62.6	72.0	36.8	60.5	34.9	36.7	39.6	41.7	42.3	<b>28.6</b>	32.3	43.8	47.7	50.9	46.4
	± 1.0	± 0.7	± 0.2	± 0.7	± 0.7	± 0.5	± 0.2	± 0.2	± 0.6	± 0.3	± <b>0.2</b>	± 0.9	± 0.2	± 0.4	± 0.2	
NOTE*	51.8	50.0	60.7	32.6	54.4	<b>33.0</b>	33.5	<b>38.5</b>	<b>38.6</b>	<b>36.7</b>	29.7	<b>27.3</b>	43.2	37.1	47.6	41.0
	± 1.0	± 0.3	± 0.4	± 0.2	± 0.3	± <b>0.2</b>	± 0.4	± <b>0.3</b>	± <b>0.1</b>	± <b>0.3</b>	± 0.5	± <b>0.3</b>	± 0.4	± 0.2	± 0.9	

Table 10: Average classification error (%) and their corresponding standard deviations on ImageNet-C with **temporally correlated test streams**, shown per corruption. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	<i>Gaussian</i>	<i>Shot</i>	<i>Impulse</i>	<i>Defocus</i>	<i>Glass</i>	<i>Motion</i>	<i>Zoom</i>	<i>Snow</i>	<i>Frost</i>	<i>Fog</i>	<i>Brightness</i>	<i>Contrast</i>	<i>Elastic</i>	<i>Pixelate</i>	<i>JPEG</i>	Avg	
Source	98.4	97.7	98.4	90.6	92.5	89.8	81.8	89.5	85.0	86.4	51.1	97.2	85.3	76.9	71.7	86.1	
BN Stats	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	89.4
ONDA	89.2	88.2	89.0	<b>90.8</b>	90.0	81.3	69.8	72.6	73.8	62.6	44.3	92.1	64.5	60.3	70.7	76.0	
PL	89.8	86.1	88.5	<b>93.0</b>	92.5	82.2	64.6	70.2	79.7	<b>55.8</b>	43.9	97.2	<b>57.8</b>	<b>52.7</b>	<b>60.5</b>	74.4	
TENT	±1.9 ±0.9	±1.6 ±1.1	±0.6 ±0.0	±0.3 ±0.6	±0.4 ±0.2	±0.1 ±0.0	±0.1 ±0.0	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	±0.0 ±0.1	89.2
LAME	91.1	89.7	91.0	<b>93.1</b>	92.2	84.7	72.4	73.3	78.7	59.8	44.5	95.2	61.6	56.4	67.4	76.5	
CoTTA	±2.4 ±1.6	±2.5 ±3.2	±3.2 ±4.9	±3.5 ±1.1	±6.9 ±4.0	±0.5 ±4.3	±4.3 ±5.6	±4.7 ±7.1	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	89.8
NOTE	87.6	85.7	87.2	<b>83.3</b>	<b>83.2</b>	<b>73.6</b>	65.4	<b>65.0</b>	<b>68.6</b>	57.9	43.5	<b>75.9</b>	61.2	54.1	62.8	<b>70.3</b>	
NOTE*	±0.1 ±0.1	±0.2 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.1 ±0.1	±0.0 ±0.0	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.0 ±0.0	±0.1 ±0.1	±0.1 ±0.1	89.5
	89.5	87.9	88.9	84.6	83.7	74.4	66.6	66.1	71.2	58.2	44.7	78.8	61.2	54.8	64.8	71.7	
	±0.4 ±0.2	±0.3 ±0.2	±0.2 ±0.2	±0.2 ±0.2	±0.1 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.2 ±0.2	±0.0 ±0.0	±0.1 ±0.1	±0.1 ±0.1	89.5

Table 11: Average classification error (%) and their corresponding standard deviations on MNIST-C with **uniformly distributed test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs. NOTE\* indicates NOTE used directly with test batches (without using PBRs).

Method	<i>Shot</i>	<i>Impulse</i>	<i>Glass</i>	<i>Motion</i>	<i>Shear</i>	<i>Scale</i>	<i>Rotate</i>	<i>Brightness</i>	<i>Translate</i>	<i>Stripe</i>	<i>Fog</i>	<i>Spatter</i>	<i>Dotted line</i>	<i>Zigzag</i>	<i>Canny edges</i>	Avg	
Source	3.7	27.3	20.4	4.6	2.2	5.1	6.5	21.1	13.8	17.4	66.6	3.8	3.7	18.2	26.4	16.1	
BN Stats [29]	±0.7 ±5.5	±6.4 ±0.5	±0.5 ±1.0	±1.0 ±1.0	±22.9 ±1.4	±17.0 ±14.7	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	±0.4 ±0.4	2.9
ONDA [27]	2.9	7.0	9.1	3.0	2.0	3.8	6.1	1.1	12.5	6.5	2.2	3.3	2.5	11.4	6.7	5.3	
PL [22]	±0.7 ±1.6	±1.0 ±0.8	±0.3 ±0.2	±0.7 ±0.1	±0.8 ±0.8	±2.6 ±0.5	±0.3 ±0.2	±0.2 ±0.2	±0.9 ±0.4	±0.2 ±0.4	±0.2 ±0.4	±0.2 ±0.4	±0.2 ±0.4	±0.2 ±0.4	±0.2 ±0.4	±0.2 ±0.4	2.6
TENT [41]	1.6	3.5	4.8	1.7	1.5	2.3	4.9	0.8	6.8	2.7	1.0	2.2	1.7	5.3	3.9	3.0	
LAME [4]	±0.3 ±0.7	±0.8 ±0.0	±0.0 ±0.1	±0.7 ±0.1	±0.8 ±0.8	±0.6 ±0.0	±0.3 ±0.2	±0.4 ±0.9	±0.1 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	1.4
CoTTA [44]	±0.1 ±0.4	±0.5 ±0.0	±0.0 ±0.1	±0.2 ±0.1	±0.1 ±0.7	±0.1 ±0.8	±0.6 ±0.0	±0.3 ±0.2	±0.4 ±0.9	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	3.0
NOTE	3.0	<b>30.7</b>	18.9	3.4	1.9	4.2	6.3	<b>25.9</b>	<b>13.9</b>	<b>18.5</b>	<b>78.2</b>	3.3	3.2	<b>19.3</b>	<b>28.0</b>	<b>17.2</b>	
NOTE*	±0.8 ±8.3	±5.8 ±0.5	±0.3 ±0.5	±0.9 ±29.8	±1.9 ±21.2	±9.8 ±0.7	±0.3 ±3.2	±12.7 ±0.2	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	2.6
	2.5	10.7	10.9	2.0	1.5	2.4	5.5	0.9	5.5	12.1	1.2	2.8	3.0	10.9	9.1	5.4	
	±0.8 ±1.9	±2.0 ±0.3	±0.0 ±0.1	±0.3 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.2 ±0.2	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	2.5
	<b>1.3</b>	<b>2.7</b>	<b>3.8</b>	<b>1.3</b>	<b>1.1</b>	<b>1.6</b>	<b>3.5</b>	<b>0.7</b>	<b>2.8</b>	2.2	<b>0.7</b>	<b>1.7</b>	1.4	4.8	3.5	<b>2.2</b>	
	±0.2 ±0.1	±0.5 ±0.1	±0.1 ±0.1	±0.1 ±0.1	±0.0 ±0.0	±0.1 ±0.1	±0.0 ±0.0	±0.0 ±0.0	±0.0 ±0.0	±0.1 ±0.1	±0.1 ±0.1	±0.4 ±0.2	±0.2 ±0.2	±0.1 ±0.1	±0.1 ±0.1	±0.1 ±0.1	2.6

## B.2 Real distributions with domain shift

Since the adaptation is done from a single source domain to a single target domain in KITTI, no further per-domain tables are specified here.

Table 12: Average classification error (%) and their corresponding standard deviations on HARTH with **real test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Averaged over three runs.

Method	S008	S018	S019	S021	S022	S028	S029	Avg
Source	86.2 ± 1.3	44.7 ± 2.1	50.4 ± 9.5	74.8 ± 3.8	72.0 ± 2.6	53.0 ± 24.0	57.0 ± 16.7	62.6
BN Stats [29]	70.3 ± 1.4	<b>73.8 ± 1.3</b>	<b>68.1 ± 3.0</b>	64.9 ± 0.9	68.5 ± 0.3	<b>65.5 ± 0.5</b>	<b>69.4 ± 1.4</b>	<b>68.6</b>
ONDA [27]	75.3 ± 4.0	<b>60.4 ± 0.9</b>	<b>63.1 ± 4.6</b>	67.9 ± 0.4	70.0 ± 3.8	<b>73.6 ± 0.7</b>	<b>74.5 ± 4.4</b>	<b>69.3</b>
PL [22]	60.4 ± 1.3	<b>71.4 ± 1.5</b>	<b>62.9 ± 1.9</b>	61.8 ± 1.2	63.1 ± 0.4	<b>64.5 ± 0.8</b>	<b>69.4 ± 2.0</b>	<b>64.8</b>
TENT [41]	<b>59.5 ± 0.3</b>	<b>71.0 ± 1.6</b>	<b>62.2 ± 1.9</b>	<b>61.1 ± 1.1</b>	<b>61.7 ± 0.4</b>	<b>64.1 ± 0.5</b>	<b>69.3 ± 2.1</b>	<b>64.1</b>
LAME [4]	85.5 ± 1.7	43.4 ± 2.0	48.8 ± 10.9	73.2 ± 3.8	70.7 ± 2.6	51.2 ± 29.4	54.1 ± 20.6	61.0
CoTTA [44]	70.4 ± 1.4	<b>73.8 ± 1.3</b>	<b>68.2 ± 2.9</b>	64.9 ± 1.0	68.5 ± 0.2	<b>65.5 ± 0.5</b>	<b>69.4 ± 1.4</b>	<b>68.7</b>
NOTE	84.8 ± 0.7	<b>32.9 ± 1.8</b>	<b>36.3 ± 10.9</b>	69.1 ± 2.4	67.1 ± 1.2	<b>30.0 ± 13.8</b>	<b>36.6 ± 9.8</b>	<b>51.0</b>

Table 13: Average classification error (%) and their corresponding standard deviations on Extrasensory with **real test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors, while **Red** fonts show performance degradation after adaptation. Due to the length of the name of each domain, denoted here with the first three characters. Averaged over three runs.

Method	4FC	598	619	797	A5D	C48	D7D	Avg
Source	34.6 ± 2.5	40.1 ± 0.7	63.8 ± 5.7	45.3 ± 2.4	64.6 ± 3.7	39.6 ± 6.8	63.0 ± 3.9	50.2
BN Stats[29]	<b>61.7 ± 4.2</b>	<b>50.1 ± 5.1</b>	51.6 ± 1.5	<b>59.4 ± 1.1</b>	54.4 ± 1.0	<b>52.4 ± 2.8</b>	62.6 ± 2.9	<b>56.0</b>
ONDA [27]	<b>36.3 ± 3.5</b>	<b>44.0 ± 2.2</b>	<b>50.8 ± 2.4</b>	<b>56.1 ± 1.9</b>	59.7 ± 2.7	<b>43.5 ± 5.9</b>	<b>46.7 ± 4.2</b>	48.2
PL [22]	<b>62.2 ± 4.3</b>	<b>50.0 ± 5.1</b>	51.7 ± 1.8	<b>59.2 ± 1.1</b>	53.9 ± 1.1	<b>52.3 ± 2.9</b>	62.8 ± 3.0	<b>56.0</b>
TENT [41]	<b>62.1 ± 4.6</b>	<b>49.8 ± 5.0</b>	51.6 ± 1.9	<b>59.4 ± 1.2</b>	53.9 ± 1.0	<b>52.2 ± 2.9</b>	62.8 ± 3.0	<b>56.0</b>
LAME [4]	<b>33.1 ± 2.4</b>	<b>37.8 ± 0.4</b>	<b>68.0 ± 8.8</b>	<b>37.1 ± 6.7</b>	<b>73.2 ± 2.6</b>	39.0 ± 7.6	<b>66.4 ± 4.0</b>	<b>50.7</b>
CoTTA [44]	<b>61.7 ± 4.2</b>	<b>50.0 ± 4.9</b>	51.6 ± 1.5	<b>59.4 ± 1.1</b>	54.4 ± 1.0	<b>52.4 ± 2.8</b>	62.6 ± 2.9	<b>56.0</b>
NOTE	<b>41.7 ± 5.9</b>	<b>40.7 ± 0.8</b>	55.5 ± 10.8	<b>45.8 ± 4.6</b>	<b>45.8 ± 10.4</b>	<b>32.9 ± 1.1</b>	55.5 ± 10.4	<b>45.4</b>

### B.3 Ablation study

Table 14: Average classification error (%) and their corresponding standard deviations of varying ablation settings on CIFAR10-C with **temporally correlated test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	Gaussian	Shot	Impulse	Defocus	Glass	Motion	Zoom	Snow	Frost	Fog	Brightness	Contrast	Elastic	Pixelate	JPEG	Avg
Source	74.0	66.8	75.3	43.3	48.0	32.6	35.2	22.0	33.0	25.9	8.5	66.1	23.4	53.6	26.8	42.3
	± 3.3	± 3.5	± 4.2	± 2.7	± 2.7	± 1.2	± 2.6	± 0.4	± 2.5	± 0.9	± 0.3	± 1.8	± 0.7	± 0.7	± 0.7	
IABN	44.5	41.3	48.0	16.3	39.9	13.8	16.1	14.9	17.8	16.3	7.6	8.8	22.5	34.0	26.7	24.6
	± 2.7	± 2.3	± 1.9	± 1.0	± 0.1	± 0.7	± 0.7	± 0.3	± 0.6	± 0.6	± 0.2	± 0.3	± 0.3	± 1.2	± 0.6	
PBRs	45.2	38.5	46.8	24.5	38.2	19.1	20.0	16.5	19.1	16.5	<b>7.1</b>	34.4	21.5	39.8	<b>25.2</b>	27.5
	± 3.0	± 4.9	± 3.3	± 2.2	± 2.8	± 0.9	± 0.2	± 0.2	± 2.4	± 0.4	± <b>0.7</b>	± 3.0	± 0.5	± 4.7	± <b>0.4</b>	
IABN + RS	<b>33.7</b>	<b>30.0</b>	<b>37.6</b>	<b>13.6</b>	<b>34.9</b>	12.4	<b>14.5</b>	<b>13.9</b>	<b>15.0</b>	<b>14.0</b>	7.2	<b>7.4</b>	21.1	<b>26.2</b>	25.9	<b>20.5</b>
	± <b>6.4</b>	± <b>6.7</b>	± <b>2.9</b>	± <b>0.3</b>	± <b>1.9</b>	± 1.2	± <b>1.7</b>	± <b>1.1</b>	± <b>3.1</b>	± <b>1.3</b>	± 0.0	± <b>0.7</b>	± 0.9	± <b>4.4</b>	± 1.1	
IABN + PBRs	34.9	32.3	39.6	<b>13.6</b>	35.8	<b>11.8</b>	<b>14.5</b>	14.1	15.2	14.2	7.7	7.6	<b>20.8</b>	27.7	26.4	21.1
	± 1.6	± 3.1	± 2.5	± <b>0.5</b>	± 1.9	± <b>0.8</b>	± <b>0.5</b>	± 0.6	± 1.3	± 0.6	± 0.3	± 0.6	± <b>0.7</b>	± 2.6	± 0.5	

Table 15: Average classification error (%) and their corresponding standard deviations of varying ablation settings on CIFAR100-C with **temporally correlated test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	<i>Gaussian</i>	<i>Shot</i>	<i>Impulse</i>	<i>Defocus</i>	<i>Glass</i>	<i>Motion</i>	<i>Zoom</i>	<i>Snow</i>	<i>Frost</i>	<i>Fog</i>	<i>Brightness</i>	<i>Contrast</i>	<i>Elastic</i>	<i>Pixelate</i>	<i>JPEG</i>	Avg
Source	88.1	86.8	93.7	64.9	79.7	55.5	57.7	53.8	66.3	59.3	33.0	81.4	49.2	73.6	55.5	66.6
	±0.2	±0.6	±0.6	±0.4	±0.9	±0.3	±0.2	±0.4	±0.8	±0.4	±0.3	±0.4	±0.4	±1.1	±0.3	
IABN	79.3	77.2	84.2	45.0	69.6	40.9	43.1	42.5	48.6	52.5	30.4	40.5	47.6	59.8	56.2	54.5
	±0.7	±0.7	±1.0	±0.6	±0.3	±0.3	±0.6	±0.4	±0.3	±0.5	±0.1	±0.7	±0.5	±1.1	±0.4	
PBRS	68.8	66.2	73.3	46.2	64.9	41.8	41.7	44.2	48.5	44.7	<b>28.3</b>	60.1	<b>44.2</b>	51.9	<b>50.5</b>	51.7
	±0.6	±0.4	±0.9	±0.6	±1.5	±0.6	±0.3	±0.4	±0.7	±0.2	<b>±0.2</b>	±0.4	<b>±0.4</b>	±0.8	<b>±0.5</b>	
IABN + RS	66.8	65.2	73.1	38.7	63.0	36.6	38.0	41.9	43.9	44.6	29.5	33.5	46.0	49.9	52.4	48.2
	±2.1	±0.3	±1.0	±0.4	±0.9	±0.0	±0.2	±0.8	±0.4	±0.5	±0.3	±0.7	±0.5	±0.9	±0.4	
IABN + PBRS	<b>66.2</b>	<b>64.2</b>	<b>72.6</b>	<b>37.2</b>	<b>61.1</b>	<b>35.4</b>	<b>37.4</b>	<b>40.0</b>	<b>42.5</b>	<b>43.4</b>	29.4	<b>32.1</b>	44.3	<b>47.5</b>	51.3	<b>47.0</b>
	<b>±0.8</b>	<b>±1.6</b>	<b>±0.4</b>	<b>±0.8</b>	<b>±0.7</b>	<b>±0.3</b>	<b>±0.4</b>	<b>±0.4</b>	<b>±0.3</b>	<b>±0.5</b>	±0.1	<b>±0.5</b>	±0.4	<b>±0.6</b>	±0.3	

Table 16: Average classification error (%) and their corresponding standard deviations of varying ablation settings on CIFAR10-C with **uniformly distributed test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	<i>Gaussian</i>	<i>Shot</i>	<i>Impulse</i>	<i>Defocus</i>	<i>Glass</i>	<i>Motion</i>	<i>Zoom</i>	<i>Snow</i>	<i>Frost</i>	<i>Fog</i>	<i>Brightness</i>	<i>Contrast</i>	<i>Elastic</i>	<i>Pixelate</i>	<i>JPEG</i>	Avg
Source	74.0	66.8	75.3	43.3	48.0	32.6	35.2	22.0	33.0	25.9	8.5	66.1	23.4	53.6	26.8	42.3
	±3.3	±3.5	±4.2	±2.7	±2.7	±1.2	±2.6	±0.4	±2.5	±0.9	±0.3	±1.8	±0.7	±0.7	±0.7	
IABN	44.5	41.4	48.1	16.3	39.9	13.9	16.2	14.9	17.9	16.4	7.6	8.8	22.5	34.1	26.7	24.6
	±2.7	±2.3	±1.9	±1.0	±0.1	±0.7	±0.7	±0.3	±0.6	±0.5	±0.2	±0.3	±0.4	±1.2	±0.6	
PBRS	43.4	37.9	46.2	21.8	36.8	18.1	17.6	16.1	19.3	15.2	<b>7.1</b>	32.5	<b>20.0</b>	30.7	<b>23.8</b>	25.8
	±0.8	±0.6	±1.5	±2.0	±1.0	±0.3	±0.8	±0.1	±0.5	±0.3	<b>±0.4</b>	±1.5	<b>±0.2</b>	±0.7	<b>±0.1</b>	
IABN + RS	33.8	31.1	40.4	13.3	35.6	11.8	13.2	14.6	<b>14.9</b>	14.7	7.7	8.1	22.3	<b>24.6</b>	25.1	20.7
	±1.6	±0.9	±1.3	±0.7	±0.2	±0.6	±0.3	±0.3	<b>±0.6</b>	±0.4	±0.2	±0.4	±0.5	<b>±1.9</b>	±1.2	
IABN + PBRS	<b>33.5</b>	<b>30.0</b>	<b>38.2</b>	<b>12.6</b>	<b>34.4</b>	<b>11.5</b>	<b>12.9</b>	<b>14.1</b>	15.2	<b>14.0</b>	7.4	<b>7.8</b>	20.7	24.7	24.2	<b>20.1</b>
	<b>±1.7</b>	<b>±1.6</b>	<b>±0.9</b>	<b>±0.8</b>	<b>±0.8</b>	<b>±0.5</b>	<b>±0.6</b>	<b>±0.2</b>	±0.8	<b>±0.6</b>	±0.2	<b>±0.2</b>	±0.3	±0.7	±0.4	

Table 17: Average classification error (%) and their corresponding standard deviations of varying ablation settings on CIFAR100-C with **uniformly distributed test streams**, shown per domain. **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	<i>Gaussian</i>	<i>Shot</i>	<i>Impulse</i>	<i>Defocus</i>	<i>Glass</i>	<i>Motion</i>	<i>Zoom</i>	<i>Snow</i>	<i>Frost</i>	<i>Fog</i>	<i>Brightness</i>	<i>Contrast</i>	<i>Elastic</i>	<i>Pixelate</i>	<i>JPEG</i>	Avg
Source	88.1	86.8	93.7	64.9	79.7	55.5	57.7	53.8	66.3	59.3	33.0	81.4	49.2	73.6	55.5	66.6
	±0.2	±0.6	±0.6	±0.4	±0.9	±0.3	±0.2	±0.4	±0.8	±0.4	±0.3	±0.4	±0.4	±1.1	±0.3	
IABN	79.3	77.2	84.3	45.0	69.6	40.9	43.1	42.5	48.6	52.5	30.5	40.5	47.6	59.8	56.2	54.5
	±0.6	±0.6	±1.0	±0.5	±0.2	±0.3	±0.6	±0.4	±0.3	±0.5	±0.1	±0.7	±0.5	±1.1	±0.4	
PBRS	68.6	66.0	72.9	45.3	64.1	40.9	41.6	43.7	47.9	44.2	<b>28.3</b>	59.9	44.2	51.1	<b>50.4</b>	51.3
	±1.0	±0.3	±0.3	±0.3	±0.8	±0.5	±0.5	±0.2	±0.2	±0.3	<b>±0.3</b>	±0.7	±0.5	±1.6	<b>±0.6</b>	
IABN + RS	67.1	65.6	74.0	39.0	61.4	36.5	38.7	41.4	44.0	45.0	30.0	34.0	46.0	48.8	52.5	48.3
	±1.2	±0.3	±0.4	±0.3	±1.3	±0.1	±0.8	±0.2	±0.4	±0.2	±0.2	±0.2	±1.4	±1.3	±0.5	
IABN + PBRS	<b>65.6</b>	<b>62.6</b>	<b>72.0</b>	<b>36.8</b>	<b>60.5</b>	<b>34.9</b>	<b>36.7</b>	<b>39.6</b>	<b>41.7</b>	<b>42.3</b>	28.6	<b>32.3</b>	<b>43.8</b>	<b>47.7</b>	50.9	<b>46.4</b>
	<b>±1.0</b>	<b>±0.7</b>	<b>±0.2</b>	<b>±0.7</b>	<b>±0.7</b>	<b>±0.5</b>	<b>±0.2</b>	<b>±0.2</b>	<b>±0.6</b>	<b>±0.3</b>	±0.2	<b>±0.9</b>	<b>±0.2</b>	<b>±0.4</b>	±0.2	

## C Replacing BN with IABN during test time

Table 18: Average classification error (%) and corresponding standard deviations of varying ablation settings on CIFAR10-C/100-C under temporally correlated (non-i.i.d.) and uniformly distributed (i.i.d.) test data stream. IABN\* refers to replacing BN with IABN during test time (no pre-training with IABN layers). **Bold** fonts indicate the lowest classification errors. Averaged over three runs.

Method	Temporally correlated test stream			Uniformly distributed test stream		
	CIFAR10-C	CIFAR100-C	Avg	CIFAR10-C	CIFAR100-C	Avg
Source	42.3 ± 1.1	66.6 ± 0.1	54.4	42.3 ± 1.1	66.6 ± 0.1	54.4
<b>IABN*</b>	27.1 ± 0.4	60.8 ± 0.1	44.0	27.1 ± 0.4	60.8 ± 0.2	44.0
IABN	24.6 ± 0.6	54.5 ± 0.1	39.5	24.6 ± 0.6	54.5 ± 0.1	39.5
<b>IABN*+PBRS</b>	24.9 ± 0.2	55.9 ± 0.2	40.4	23.2 ± 0.4	55.3 ± 0.1	39.3
<b>IABN+PBRS</b>	<b>21.1 ± 0.6</b>	<b>47.0 ± 0.1</b>	<b>34.0</b>	<b>20.1 ± 0.5</b>	<b>46.4 ± 0.0</b>	<b>33.2</b>

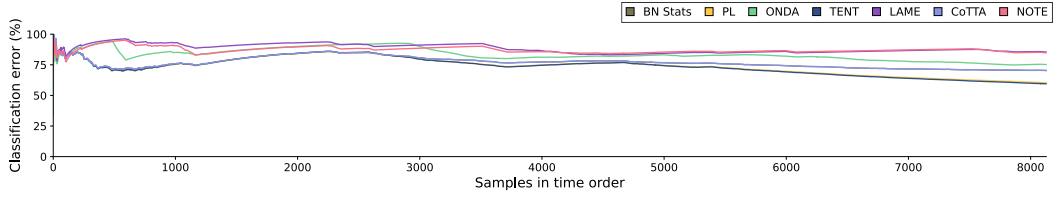
For pre-trained models with BN layers such as ResNet [12], NOTE needs to re-train the model by replacing BN layers with IABN layers in order to utilize the effectiveness of IABN. This requires the additional computational cost of re-training, which might make it inconvenient to utilize off-the-shelf models. We further investigate whether simply switching BN to IABN without re-training still leads to performance gain.

Table 18 shows the result of this experiment, where IABN\* refers to replacing BN with IABN during test time. We note that IABN\* still shows a significant reduction of errors under CIFAR10-C and CIFAR100-C datasets compared with BN (Source). We interpret this as the normalization correction in IABN is somewhat valid without re-training the model. We notice that IABN\* outperforms the baselines in CIFAR10-C with 27.1% error, while the second best (LAME) shows 36.2% error 19. In addition, IABN\* also shows improvement combined with PBRS. This implies that IABN can be used without re-training the model, which aligns with the fully test-time adaptation paradigm introduced in a recent study [41].

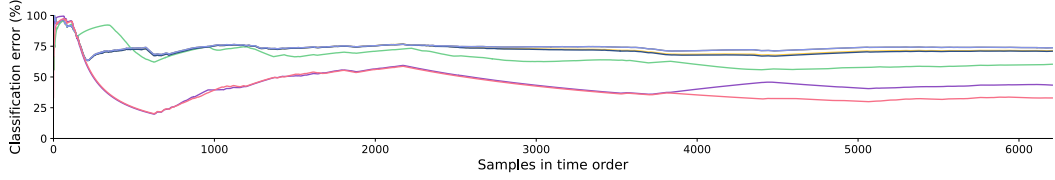
## D License of assets

**Datasets** KITTI dataset (CC-BY-NC-SA 3.0), KITTI-rain dataset (CC-BY-NC-SA 3.0), CIFAR10, 100 (MIT License), ImageNet-C (Apache 2.0), MNIST-C (CC-BY-NC-SA 4.0), HARTH dataset (MIT License), and the Extrasensory dataset (CC-BY-NC-SA 4.0)

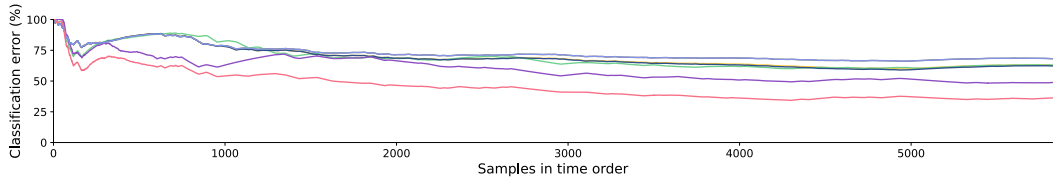
**Codes** Code for rain augmentation on the KITTI dataset (Apache 2.0), torch-vision for ResNet18 and ResNet50 (Apache 2.0), code for depth estimation used in rain augmentation on the KITTI dataset (UCLB ACP-A License), code for generating Dirichlet distributions (Apache 2.0), the official repository of CoTTA (MIT License), the official repository of TENT (MIT License), and the official repository of LAME (CC BY-NC-SA 4.0).



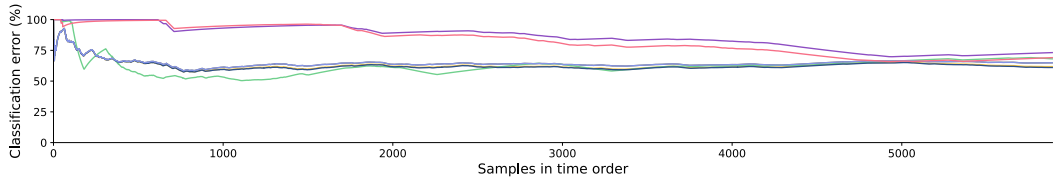
(a) S008.



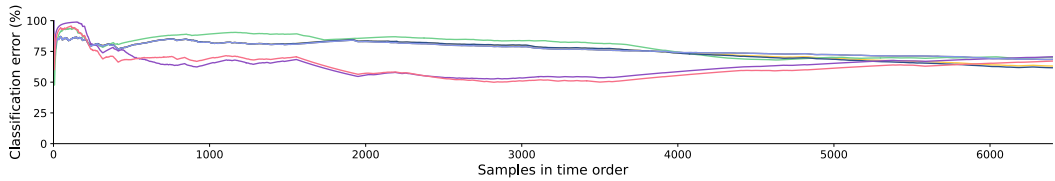
(b) S018.



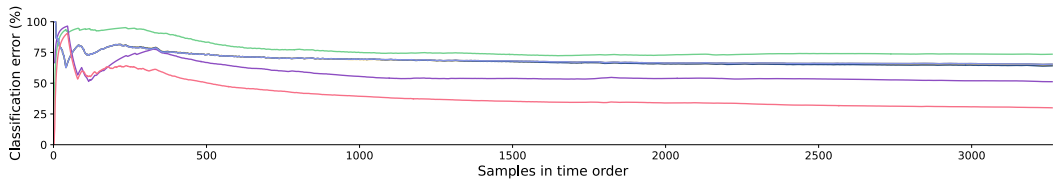
(c) S019.



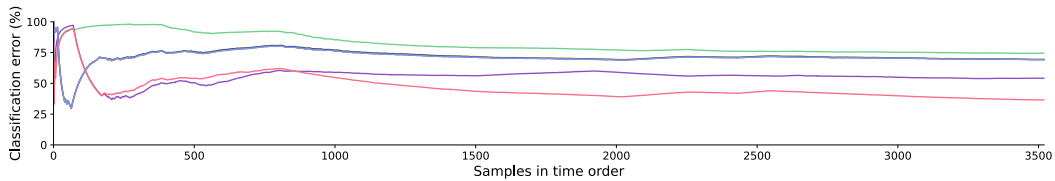
(d) S021.



(e) S022.



(f) S028.



(g) S029.

Figure 7: Illustration of the real-time cumulative classification error change of different methods on the HARTH dataset. The x-axis denotes the samples in order, whereas the y-axis denotes the error rate in percentage. Note that some lines are not clearly visible due to overlap.



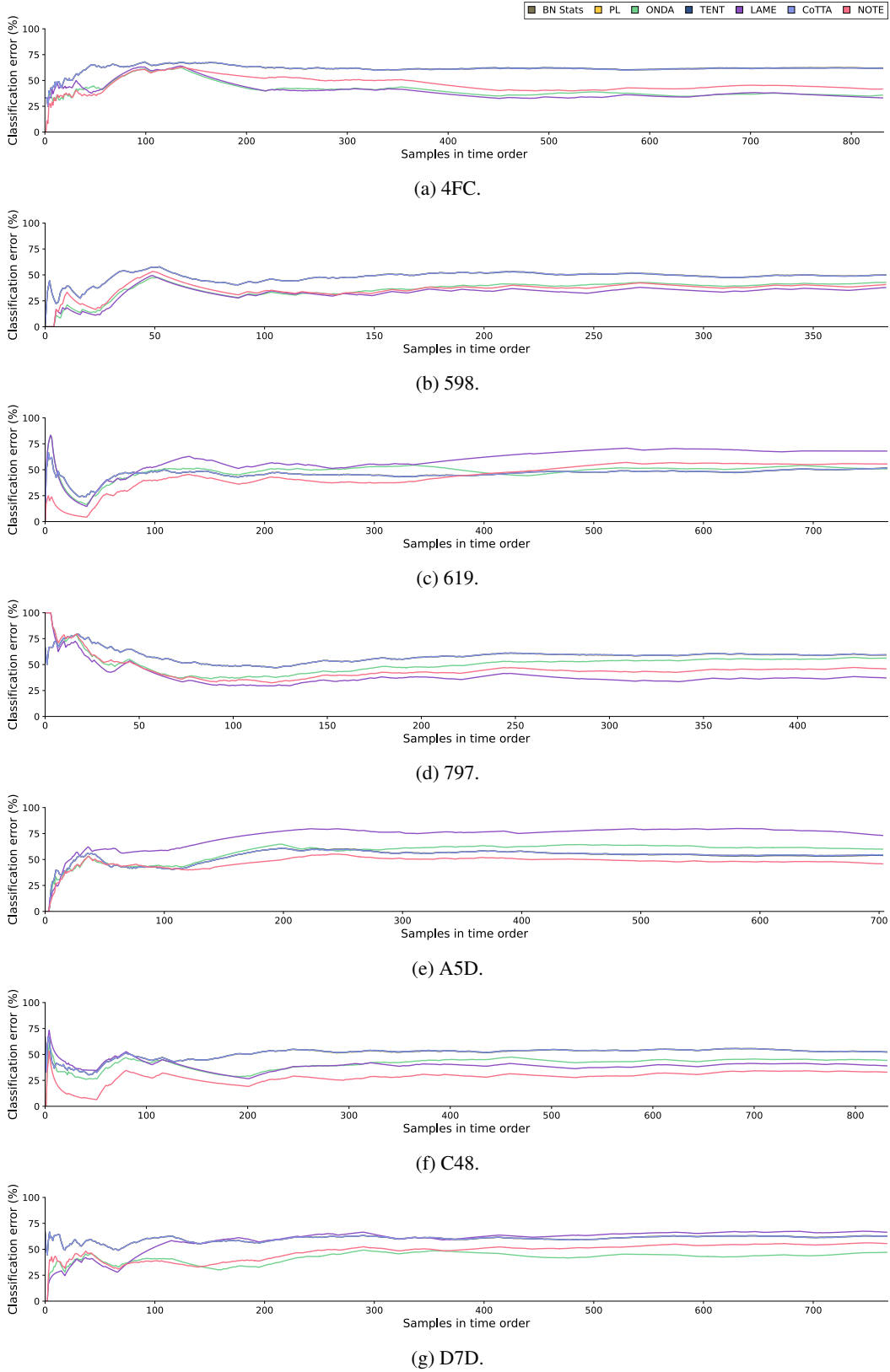


Figure 8: Illustration of the real-time cumulative classification error change of different methods on the Extrasensory dataset. The x-axis denotes the samples in order, whereas the y-axis denotes the error rate in percentage. Note that some lines are not clearly visible due to overlap.