

# COME: TEST-TIME ADAPTION BY CONSERVATIVELY MINIMIZING ENTROPY

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Machine learning models must continuously self-adjust themselves for novel data distribution in the open world. As the predominant principle, entropy minimization (EM) has been proven to be a simple yet effective cornerstone in existing test-time adaption (TTA) methods. While unfortunately its fatal limitation (i.e., overconfidence) tends to result in model collapse. For this issue, we propose to **Conservatively Minimize the Entropy (COME)**, which is a simple drop-in replacement of traditional EM to elegantly address the limitation. In essence, COME explicitly models the uncertainty by characterizing a Dirichlet prior distribution over model predictions during TTA. By doing so, COME naturally regularizes the model to favor conservative confidence on unreliable samples. Theoretically, we provide a preliminary analysis to reveal the ability of COME in enhancing the optimization stability by introducing a data-adaptive lower bound on the entropy. Empirically, our method achieves state-of-the-art performance on commonly used benchmarks, showing significant improvements in terms of classification accuracy and uncertainty estimation under various settings including standard, life-long and open-world TTA, i.e., up to 34.5% improvement on accuracy and 15.1% on false positive rate. Our code is available at: <https://anonymous.4open.science/r/anonymous-9F46>.

## 1 INTRODUCTION

Endowing machine learning models with self-adjust ability is essential for their deployment in the open world, such as autonomous vehicle control and embodied AI systems. To this end, test-time adaption (TTA) emerges as a promising strategy to enhance the performance in the open world which often encounters unexpected noise or corruption (e.g., data from rainy or snowy weather). Unsupervised losses play a crucial role in model adaptation, which can improve the accuracy of a model on novel distributional test data without the need for additional labeled training data. The representative strategy entropy minimization (EM) adapts classifiers by iteratively increasing the probabilities assigned to the most likely classes, and is an integral part in the state-of-the-art TTA methods (Press et al., 2024; Wang et al., 2021; Zhang et al., 2022; Niu et al., 2022; Wang et al., 2022b; Iwasawa & Matsuo, 2021; Niu et al., 2023; Yang et al., 2024). The initial intuition behind using entropy minimization, given by (Wang et al., 2021) is based on the observation that models tend to be more accurate on samples for which they make predictions with higher confidence. The natural extension of this observation is to encourage models to bolster the confidence on test samples.

However, this intuition may not always be true since there always exists irreducible uncertainty which arises from the natural complexity of the data or abnormal outliers. Naturally, one might expect a machine learning model to adapt itself to test data and favor higher confidence on right prediction, but of course not absolute certainty for the erroneous. This contradiction challenges the suitability of EM in TTA tasks, which greedily pursues low-entropy on all test samples. A notable example in recent research concerns that EM can be highly unstable and frequently lead to model collapse when the models encounter unreliable samples in the wild (Niu et al., 2023). In this work, we hypothesize that due to the nature of EM, previous TTA methods tend to be highly overconfident ignoring the reliability of various test samples, which further results in the unsatisfactory performance.

For the above issues, we propose a simple yet effective model-agnostic learning principle, termed **Conservatively Minimizing Entropy (COME)** to stabilize TTA. We first consider the model output as *opinion* which explicitly models the uncertainty of each sample from a Theory of Evidence

perspective. Then, we encourage the model to favor definitive opinions for TTA and meanwhile take the uncertainty information into consideration. This offers two-fold advantages compared to EM learning principle. First, our COME leverages subjective logic (Jsang, 2018), which is an off-the-shelf uncertainty tool in Bayesian toolbox to effectively perceive the uncertainty raised upon varying test samples without altering the original model architecture or training strategy. Second, when encountering unreliable outliers, the model is regularized to favor conservative confidence and be able to explicitly express "*I do not know*", i.e., reject to classify them to any known classes, which meets our expectation on model trustworthiness. Theoretically, our COME takes inspiration from Bayesian framework, and can be proved to correspond with a data-adaptive upper bound on the model confidence, which is a desirable property for TTA where the reliability of test samples are often varying from time. The contributions of this work are summarized as follows:

- As a principled alternative beyond entropy minimization, we propose a simple yet effective driven strategy for test-time adaption called Conservatively Minimizing Entropy (COME) which improves previous methods by exploring and exploiting the uncertainty.
- We provide theoretical analysis with insight in contrast to EM, the model confidence of our COME is provably upper bounded in a data-adaptive manner, which enables TTA methods to focus on reliable samples and conservatively handle abnormal test samples.
- We perform extensive experiments under various settings, including standard, open-world and lifelong TTA, where the proposed COME achieves excellent performance in terms of both classification accuracy and uncertainty quantification.

## 2 RELATED WORK

**Test-time adaption** aims to bridge the gaps between source and target domains during test-time without accessing the training-time source data. **Entropy minimization** performs an important role in test-time adaption, which has been integrated as a part of numerous TTA methods (Press et al., 2024; Wang et al., 2021; Niu et al., 2022; Wang et al., 2022b; Iwasawa & Matsuo, 2021; Yang et al., 2024; Chen et al., 2022). However, it has been observed that the performance of EM can be highly sub-optimal and unstable when encounter unreliable environments. To this end, previous works incorporate many strategies including i) Samples selection, which selectively filter out the high-entropy unreliable samples by manually setting an entropy threshold (Iwasawa & Matsuo, 2021; Niu et al., 2023). ii) Constrained optimization, which heuristically enforces that the updated parameters do not diverge too much during adaption. iii) Model recovery, which lively monitor the state of the adapting model and frequently reset it when detecting performance collapse (Niu et al., 2023; Wang et al., 2022b). Although these strategies have shown promising performance, the underlying issues of the EM are still largely unexplored. In contrast, this work aims to handle the inherent issues of EM learning objective and validate the necessity and effectiveness in TTA settings.

**Uncertainty quantification** is one key aspect of the model reliability. With accurate uncertainty estimation ability, further processing can be taken to improve the performance of machine learning systems (e.g., human assistance) when the predictive uncertainty is high. This is especially useful in high-stake scenarios such as medical diagnosis (Wang et al., 2023). To obtain the uncertainty, Bayesian neural networks (BNNs) (Denker & LeCun, 1990; Mackay, 1992) have been proposed to replace the deterministic weight parameters of model with distribution. Unlike BNNs, ensemble-based methods obtain the uncertainty by training multiple models and ensembling them (Rahaman et al., 2021; Abe et al., 2022). In this paper, we focus on estimating and exploiting uncertainty under the theory of subjective logic (SL, (Jsang, 2018)). Unlike BNNs or ensemble, SL explicitly models the uncertainty in a single forward pass without modifying the training strategy or model architecture, which meets our expectation of computational effectiveness for TTA tasks.

## 3 MOTIVATION

We consider the fully test-time adaption setting in  $K$ -classification task where  $\mathcal{X}$  is the input space and  $\mathcal{Y} = \{1, 2, \dots, K\}$  denotes the target space. Given a classifier  $f : \mathcal{X} \rightarrow \mathbb{R}^K$  parameterized by  $\theta$  which has been pretrained on training distribution  $P^{\text{train}}$ , our goal is to boost  $f$  by updating its parameters  $\theta$  online on each batch of test data drawn from test distribution  $P^{\text{test}}$ . Note that in fully

TTA setting, the training data  $P^{\text{train}}$  is unavailable and one can only tune  $\theta$  on unlabeled test data. This is derived from realistic concerns of privacy, bandwidth or profit. Compared to other closely related setting, i.e., Source-Free Domain Adaptation (SFDA), TTA focuses on online adjusting during the testing while SFDA generally perform offline. That is, TTA method aims on (unsupervised) adaptation online and the inference latency matters. Entropy minimization (EM) algorithm iteratively optimizes the model to minimize the predictive entropy on test sample  $x$

$$H(p(y|x)) = - \sum_{k=1}^K p(y = k|x) \log p(y = k|x), \tag{1}$$

where  $p(y|x)$  is the class distribution calculated by normalizing the output logits  $f(x)$  with softmax function, i.e.,  $p(y = k|x) = \frac{\exp f_k(x)}{\sum \exp f(x)}$ .  $H$  is the Shannon’s entropy.

**Other learning objectives.** Besides EM, there also exists several TTA methods which explore other unsupervised learning objectives. Notable examples include 1) Pseudo label (PL):  $\mathcal{L}_{\text{PL}} = -\mathbb{E} \log p(y = \hat{y}|x)$  which encourages the adapted model to fit the pseudo label  $\hat{y}$  predicted by the pretrained model, 2) Module adjustment (T3A) which adjusts the parameters in the last fully connected layer, and can be viewed as an implicit way to minimize entropy (Iwasawa & Matsuo, 2021), 3) Energy minimization:  $\mathcal{L}_{\text{TEA}} = -\mathbb{E} \log \sum_{k=1}^K \exp f(x)$  which aims to minimize the free energy during adaption, and takes inspiration from energy model (Yige et al., 2024), 4) Contrastive learning objective:  $\mathcal{L}_{\text{infoNCE}} = -\log \frac{\exp \text{query} \cdot \text{key}^+}{\sum \exp \text{query} \cdot \text{key}}$  which strives to minimize the cosine distance between the query and positive samples (key<sup>+</sup>) while maximizing the cosine distances between query and negative samples (Chen et al., 2022), 5) The recent advanced FOA (Niu et al., 2024) which uses evolution strategy to minimize the test-training statistic discrepancy and model prediction entropy.

**The overconfident issue of EM.** We begin by testing EM in standard TTA setting, and put forward the following observations to detail its unsatisfying performance.

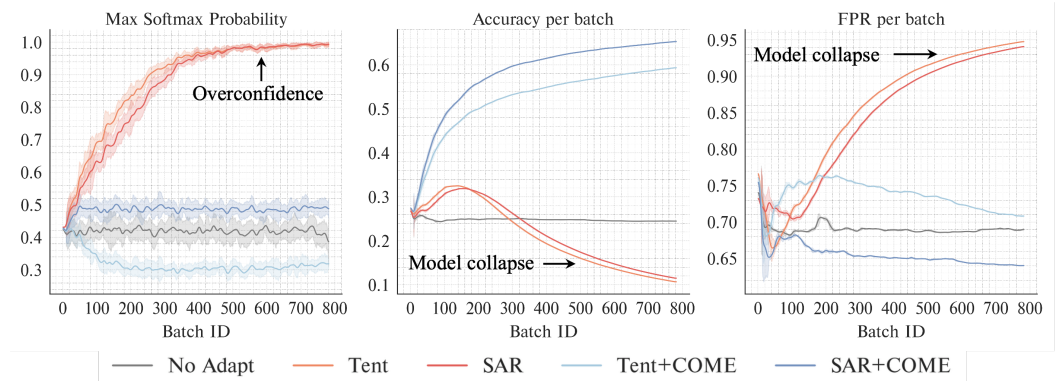


Figure 1: Empirical observations of Entropy Minimization when equipped to two representative TTA methods, i.e., the seminal Tent (Wang et al., 2021) and recent advanced SAR (Niu et al., 2023). Along the TTA process, the uncertainty of models tuned with EM quickly drops, and the false positive rate decreases temporarily for a very short time horizon before quickly increasing. Along the same adaption trajectory, the model accuracy also improves for a short time compared to the initial model and then quickly decreases, after which the model collapses to a trivial solution. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate. Besides, the model confidence of our COME is much more conservative, which implies fewer risks of overconfidence and a more accurate uncertainty awareness.

As shown in Figure 1, EM tends to overconfident prediction and assign extremely high probability to one certain class. The absolutely high confidence on all test samples is obviously a rather undesired characteristic. We test on ImageNet-C under snow corruption of severity level 5 as a typical showcase, and refer interested readers to Appendix C.11 for more similar results.

## 4 METHODOLOGY

We propose to conservatively minimize the entropy under uncertainty modeling, a simple alternative to EM algorithm. The key idea of COME is to quantify and then regularize the uncertainty during TTA without altering the model architecture or training strategy, which avoids the overconfident nature of EM at minimal cost. We first introduce uncertainty quantification by the subjective logic and then present how to regularize the uncertainty during TTA.

### 4.1 MODELING UNCERTAINTY BY THE SUBJECTIVE LOGIC

Given a well trained classifier  $f: \mathcal{X} \rightarrow \mathcal{Y}$ , the most simple way to quantify the uncertainty of each sample is using the softmax probability as confidence in prediction. A few pioneer works propose to filter out the test samples with high-entropy predicted softmax probability for stable TTA (Niu et al., 2023; 2022). However, it has been shown that softmax probability often leads to overconfident predictions, even when the predictions are wrong or the inputs are abnormal outliers (Moon et al., 2020; Van Amersfoort et al., 2020). Thus this simple strategy may not be satisfied enough and highlights the necessity of better uncertainty modeling. To this end, we propose to obtain the uncertainty through subjective logic, which defines a framework for obtaining the probabilities (belief masses) of different classes and the overall uncertainty (uncertainty mass) based on the *evidence*<sup>1</sup> collected from data. Specifically, in  $K$  classification task, SL formalizes the belief assignments over a frame of discernment as a Dirichlet distribution. In contrast to softmax function that directly normalizes the logits  $f(x)$  to model the class distribution  $p(y|x)$ , SL considers the model output as evidence (denoted as  $e$ ) to model a Dirichlet distribution which represents the density of all possible probability assignment  $\boldsymbol{\mu} = [p(y = 1|x), p(y = 2|x), \dots, p(y = K|x)]$ . That is, the predicted categoricals  $\boldsymbol{\mu}$  is also a random variable itself, which yields a Dirichlet distribution as follow

$$p(\boldsymbol{\mu}|x) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \mu_k^{\alpha_k - 1}, \quad \boldsymbol{\alpha} = \mathbf{e} + 1, \quad (2)$$

where  $\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha})$  is the Dirichlet distribution characterized by parameters  $\boldsymbol{\alpha}$ . The summation of all  $\alpha_k \in \boldsymbol{\alpha}$  is so called the strength  $S$  of the Dirichlet distribution, i.e.,  $S = \sum_k \alpha_k = \sum_k e_k + 1$ . Then SL tries to assign a belief mass  $b_k$  to each class label  $k$  and an overall uncertainty mass  $u$  to the whole frame based on the collected evidence as follow

$$b_k = \frac{e_k}{S} = \frac{\alpha_k - 1}{S} \quad \text{and} \quad u = \frac{K}{S}, \quad \text{subject to} \quad u + \sum_{k=1}^K b_k = 1, \quad (3)$$

where  $S$  is the Dirichlet strength which denotes the total evidence we collected and  $K$  is the total classes number. Eq. 3 actually describes the phenomenon where the more evidence observed for the  $k$ -th category, the greater the belief mass assigned to the  $k$ -th class. Correspondingly, the less total evidence  $S$  observed, the greater the total uncertainty  $u$ . Such assignment is so called the subjective opinion

$$\mathcal{M}(x) = [b_1, b_2, \dots, b_k, u], \quad (4)$$

which not only describes the belief of assigning  $x$  to each class  $k$  but also explicitly models the uncertainty due to lack-of-evidence. To ensure that  $\alpha \geq 1$  which meets the requirements of valid parameters of a Dirichlet, we first deploy ReLU function on the logits and then use exponential output function for obtaining the parameters of Dirichlet distribution from the model output  $f(x)$  (i.e., logits), where  $\boldsymbol{\alpha} = \exp(\text{ReLU}(f(x)))$  and  $\mathbf{e} = \boldsymbol{\alpha} - 1$ . We defer the discussion of other alternatives to realize subjective opinion to the Appendix.

**Benefits of modeling uncertainty based on subjective opinion for TTA.** It has been widely recognized that using the softmax output as the confidence often leads to overconfidence phenomenon (Guo et al., 2017; Hendrycks & Gimpel, 2017; Liu et al., 2020). Other advanced uncertainty measurement methods such as MC-dropout (Gal & Ghahramani, 2016), Bayesian neural networks (Huang et al., 2022), deep ensemble (Lakshminarayanan et al., 2017; Rahaman et al., 2021), calibration (Guo et al., 2017; Han et al., 2024) usually require additional computations during inference or a separated validation set, which are not applicable in unsupervised TTA task. By contrast, the introduced SL

<sup>1</sup>In Bayesian context, evidence refers to the metrics collected from the input to support the classification.

216 directly deduce an additional uncertainty mass through one single forward pass, which is model-  
 217 agnostic, light-weight and can be seamlessly integrated in standard pretrained classifier. We defer  
 218 the empirical comparisons between the proposed method and other Bayesian methods or learning  
 219 objectives originate from unsupervised domain adaption to Table 7 in the Appendix.

#### 221 4.2 MODEL ADAPTION BY SHARPENING THE OPINION

222 Vanilla EM minimizes the softmax entropy of the predicted class distribution  $p(y|x)$ , which inevitably  
 223 results in assigning rather high probability to one certain class. In contrast, we propose the learning  
 224 principle to minimize the entropy of opinion

$$225 \underset{\theta}{\text{minimize}} H(\mathcal{M}(x)) = - \sum_{k=1}^K b_k \log b_k - u \log u. \quad (5)$$

226 Compared to entropy minimization on softmax probability which ultimately assigns all the probability  
 227 to one certain class, the above learning principle offers the model with an additional option, i.e.,  
 228 express high overall uncertainty and reject to adapt when the observed total evidence is insufficient.  
 229 In other words, by assigning all belief masses (probability) to uncertainty  $u$ , the model can now  
 230 express “I do not know” as its predicted opinion. Alternatively, an additional hyperparameter  $\lambda$   
 231 can be introduced to weight the last term as  $-\sum_{k=1}^K b_k \log b_k - \lambda u \log u$ . The magnitude reflects  
 232 our confidence in whether the test sample should be adapted by the model or not. We validate the  
 233 effectiveness of this generalization to Table 9 in the Appendix.

#### 238 4.3 REGULARIZING UNCERTAINTY IN AN UNSUPERVISED MANNER

239 While subjective logic offers an opportunity to modeling uncertainty and reject to classify unreliable  
 240 samples, naively minimizing the entropy of opinion for TTA may still be problematic. As shown in  
 241 previous works, model pretrained with softmax output function frequently suffers from overconfi-  
 242 dence issue (Guo et al., 2017; Nguyen et al., 2015; Hendrycks et al., 2019). Therefore, the belief  
 243 mass assigned to the one certain class  $k$  by the pretrained model is usually much larger than the  
 244 uncertainty mass  $u$ . This results in the model tendency of increasing the belief mass during the  
 245 entropy minimization process, while neglecting the uncertainty function  $u$ . Motivated by the above  
 246 analysis, our next goal is to devise an effective regularization strategy for the uncertainty mass. In  
 247 supervised learning tasks, previous works leverage labeled training data to constrain the uncertainty  
 248 mass (Sensoy et al., 2018; Malinin & Gales, 2018). However, these strategies is not applicable due  
 249 to the unsupervised nature of TTA task where the training data is unavailable. This motivates us  
 250 to explore the uncertainty information lies in the pretrained model itself for regularization without  
 251 additional supervision. As one of the simplest yet effective design choices, we propose to constrain  
 252 the uncertainty mass predicted by the adapted model not to diverge too far from the pretrained model.  
 253 This results in the following constrained optimization objective

$$254 \underset{\theta}{\text{minimize}} H(\mathcal{M}(x)) \text{ subject to } |u_{\theta}(x) - u_{\theta_0}(x)| \leq \delta, \quad (6)$$

255 where  $\theta, \theta_0$  denote the adapted and pretrained model respectively, and  $u$  is the uncertainty estimated  
 256 by Eq. 3 and  $\delta$  is a threshold to prevent overly extreme model uncertainty. The magnitude of  $\delta$   
 257 represents our tolerance for uncertainty divergence.

258 **Rethink why we constrain on  $|u(x) - u_0(x)| \leq \delta$ .** The uncertainty estimated by pretrained  
 259 model may not be ideal. However, since in fully TTA task, we can only access unlabeled test data  
 260 coming online and the inference efficiency matters. Thus traditional methods devised for handling  
 261 overconfidence like calibration, ensembling, and other Bayesian methods are not applicable. The  
 262 only practically available choice is to explore the uncertainty information contained within the model  
 263 itself. As shown in previous works (Sensoy et al., 2018; Malinin & Gales, 2018), while the softmax  
 264 probability of pretrained model tend to be overconfident, subjective logic is much more reliable, which  
 265 can support the proposed regularization. The most straightforward way to realize the aforementioned  
 266 constraint is to storage the model before and after adaption and explicitly compare the uncertainty  
 267 mass. However, this may violate the efficiency requirements of TTA. Thus considering the difficulty  
 268 of constrained optimization in modern neural networks, our next target is to find a way to convert  
 269 Eq. 6 into an unconstrained form. To this end, we introduce the following Lemma.

**Lemma 1.** For any  $x \in \mathcal{X}$ , we have

$$\frac{K}{\|f(x)\|_p + \log K} \leq u \leq \frac{K^{1+1/p}}{\|f(x)\|_p}, \quad (7)$$

where  $f(x)$  is the model output logits,  $K$  is the total class number and  $\|\cdot\|_p$  denotes the  $p$ -norm.

Lemma 1 shows that the uncertainty mass of subjective opinion is bounded by the norm of the total evidence collected from the model output. Thus instead of directly constraining on  $u(x)$ , we can alternatively constrain on the  $p$ -norm of model output logits, which is more flexible. Taking inspiration from previous work in supervised learning literature (Wei et al., 2022), this can be achieved by factorize  $f(x)$  into  $f(x)/\|f(x)\|_p \cdot \|f(x)\|_p$  and then enforcing the gradient on the second term to be equal to zero during optimization. Specifically, the final minimizing objective of COME is

$$\text{minimize } H(\mathcal{M}(x)) \text{ where } f(x) = \frac{f(x)}{\|f(x)\|_p} \cdot \|f(x)\|_p^{\text{no-grad}} \cdot \tau, \quad (8)$$

and  $\|f(x)\|_p^{\text{no-grad}}$  is the  $p$ -norm of  $f(x)$  with zero gradient. This can be achieved by applying the detach operation which is a common used function in modern deep learning toolbox like PyTorch and TensorFlow. By doing so, minimizing the entropy of opinion would not influence  $\|f(x)\|_p$ .  $\tau$  is a hyper-parameter which controls the magnitude of recovered logits. Our COME can be implemented by modifying only a few lines of code in the original EM algorithm (shown as Algorithm 1).

**How to set  $\tau$  and  $p$  in practice?** Noted that the tightness for the upper and lower bounds in Lemma 1 is determined by the choice of  $p$ . The ratio between the upper and lower bound is minimized by  $p = \infty$ . The strictness of such constraint should be selected per need by the user via trial and error: if users are extremely cautious about unreliable TTA,  $p$  should be tuned up; otherwise, if a better performance is required. In our experiments, we choose  $p = 2$  and  $\tau = 1$  in accordance with Occam’s Razor. Experiments on different hypeparameters can be found in Table 10 in the Appendix.

---

**Algorithm 1:** Pseudo code of COME in a PyTorch-like style.

---

```

298 # x: the output logits, model: the test model
299 def entropy_of_opinion(x):
300     belief = exp(x) - 1 / sum(exp(x)) # belief mass
301     uncertainty = K / sum(exp(x)) # uncertainty mass
302     opinion = cat([belief, uncertainty]) # subjective opinion
303     return -sum(opinion * log(opinion)) # entropy of opinion
304
305 for data in test_loader: # load a minibatch data
306     x = model(data) # forward
307     x = x / norm(x, p=2) * norm(x, p=2).detach() # constraint in Eq.9
308     loss = entropy_of_opinion(x) # calculate loss
309     # ... [backwards and update the parameters]
```

---

**Stability of COME.** We provide preliminary theoretical understanding of the superiority of COME. As we mentioned before, one notable limitation of EM is that it enforces low entropy for all test samples while ignores the instinct complexity of wild test data. Thus at the end of TTA progress, EM ultimately produces model that yields overconfident prediction. Our COME resolves this issue and introduces an upper bound for each test sample  $x$  according to its trustworthiness. This property is formalized as follows

**Theorem 1** (Model confidence upper bound). For any  $x \in \mathcal{X}$ , if  $|u(x) - u_0(x)| \leq \delta$  holds, then we have

$$\max_k p(y = k|x) \leq \frac{1}{1 + (K - 1) \exp(-\frac{K}{u_0 - \delta})}, \quad (9)$$

where  $\max_k p(y = k|x)$  is the model confidence (class probability assigned to the most likely class) and  $K$  is the total class number.  $u_0$  is the shorthand of  $u_{\theta_0}(x)$ .

From Theorem 1, we find that the model confidence in COME has a sample-wise upper bound according to  $u_0(x)$ . In particular, it implies that the model confidence upper bound of the most

likely class decreases according to  $u_0(x)$ . For this reason, one can suspect that if the test model is uncertain about some sample  $x$  (with a rather large  $u_0$ ), it will be difficult to further increase the model confidence on such  $x$ , which is a desirable property for TTA in the wild.

## 5 EXPERIMENTS

We conduct experiments on multiple datasets with distributional shift to answer the following questions. Q1. In the standard TTA setting, does the proposed method outperform other algorithms? Q2. How does COME perform in more realistic TTA settings, such as open-world TTA or lifelong TTA? Q3. Uncertainty quantification is both the motivation behind COME and the reason for its effectiveness, does our method achieves more reliable uncertainty estimation during TTA? Q4. Ablation study - what is the key factor of performance improvement in our method?

### 5.1 SETUP

**Datasets.** Following the common practice (Niu et al., 2022; 2023), we conduct experiments on standard covariate-shifted distribution dataset ImageNet-C (a large-scale benchmark with 15 types of diverse corruption). Besides, we also consider open-world test-time adaption setting, where the test data distribution  $P^{\text{test}}$  is a mixture of both normal covariate-shifted data  $P^{\text{Cov}}$  and abnormal outliers  $P^{\text{Outlier}}$  of which the true labels do not belong to any known classes in  $P^{\text{train}}$ . [Following previous work in open-set OOD generalization literature \(Lee et al., 2023; Bai et al., 2023; Baek et al., 2024\)](#),  $P^{\text{Outlier}}$  is a suit of diverse datasets introduced by (Yang et al., 2022), including iNaturalist, Open-Image, NINCO and SSB-Hard. **Compared methods.** We compare our COME with a board line of test-time adaption methods, including both EM-based and non-EM based TTA methods. ◦ EM-based methods choose entropy minimization as their learning objective, including Tent (Wang et al., 2021), EATA (Niu et al., 2022), CoTTA (Wang et al., 2022b) and recent advanced SAR (Niu et al., 2023). ◦ Non-EM methods employ other learning objectives including Pseudo Label (PL), module adjustment (Iwasawa & Matsuo, 2021) (T3A) and energy minimization (Yige et al., 2024) (TEA). Following (Niu et al., 2023), we use the ViT-base architecture as our backbone and defer the results on ResNet50 to Table 14 in the Appendix. The test batch size is 64. When equipped to previous EM-based TTA baselines, we only replace the learning objective with our COME and keep all the other configures consistent to the official implementation. **Tasks and Metrics.** For classification performance comparison, we report the accuracy (Acc) on covariate-shifted data. Besides, for open-world TTA, we report the false positive rate (FPR). The mis-classified samples and outliers are considered as positive samples which should be of higher uncertainty compared to correct classification that is considered as negative. For all experiments, we run multiple times and report the average. [Full results with standard deviation are deferred to the Appendix.](#)

### 5.2 EXPERIMENTAL RESULTS

**Performance comparison in standard TTA settings (Q1).** As shown in Table 1, our COME establishes strong overall performance in terms of both classification accuracy. We highlight a few essential observations. Compared to EM learning principle, our COME consistently outperforms it when equipped to the same baseline methods, including Tent (Wang et al., 2021), EATA (Niu et al., 2022), CoTTA (Wang et al., 2022b) and SAR (Niu et al., 2023). As an example of our method’s improved performance, when equipped to the recent SAR, our method yields an accuracy of 64.2%, which outperforms the original implementation based on EM of 10.1%. Besides, we also compare to Non-EM TTA methods, including TEA, T3A and PL. These methods do not rely on EM learning objective, yet are less effective than EM in terms of classification performance. [For a more comprehensive evaluation, we conduct additional experiments on ImageNet-R and ImageNet-S, the results can be found in Table 11 and Table 12 in the Appendix.](#)

**Performance comparison in open-world and lifelong TTA settings (Q2).** In Table 2 and 3, we present the results under open-world and lifelong TTA settings respectively. In open-world TTA, the test data distribution is a mixture of both normal covariate-shifted data and abnormal outliers. The mixture ratio of  $P^{\text{Cov}}$  and  $P^{\text{Outlier}}$  is 0.5 following previous work (Bai et al., 2023), i.e.,  $P^{\text{test}} = 0.5P^{\text{Cov}} + 0.5P^{\text{Outlier}}$ . Such outliers arise from unknown classes that are not present in training data, which should not be of high uncertainty for model trustworthiness. According to

Table 1: Classification accuracy comparison on ImageNet-C (level 5). Substantial ( $\geq 0.5$ ) **improvement** and **degradation** compared to the baseline are highlighted in blue or brown respectively. The detailed results with standard deviation are defer to Table 5 in Appendix C.1.

Methods	COME	Noise			Blur				Weather				Digital			
		Gauss.	Shot	Impul.	Defoc	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elast.	Pixel	JPEG
No Adapt	$\times$	35.1	32.2	35.9	31.4	25.3	39.4	31.6	24.5	30.1	54.7	64.5	49.0	34.2	53.2	56.5
PL	$\times$	49.9	48.8	50.9	48.5	41.2	52.7	42.7	24.6	42.5	63.6	73.1	65.3	44.1	63.9	62.6
T3A	$\times$	34.6	31.5	35.5	32.7	27.5	40.7	33.5	25.6	30.8	56.5	64.9	50.8	38.0	54.3	58.4
TEA	$\times$	44.5	39.3	45.8	37.6	35.4	46.4	31.4	8.9	46.4	60.1	72.5	59.5	45.7	62.3	58.9
LAME	$\times$	34.8	31.9	35.5	30.9	24.4	38.9	30.7	23.4	29.5	53.3	64.2	41.0	32.7	52.8	56.0
FOA	$\times$	46.6	43.9	48.3	47.1	40.7	49.3	43.7	53.8	52.8	64.2	76.2	63.8	48.5	62.1	64.0
Tent	$\times$	52.5	52.1	53.4	52.8	47.4	56.7	47.4	10.5	26.4	67.2	74.3	67.3	50.4	66.5	64.6
	$\checkmark$	53.9	53.9	55.2	55.8	51.8	59.8	52.6	58.4	61.3	71.3	78.2	68.8	57.9	70.5	68.2
	Improve	$\Delta 1.4$	$\Delta 1.9$	$\Delta 1.8$	$\Delta 2.9$	$\Delta 4.4$	$\Delta 3.2$	$\Delta 5.2$	$\Delta 47.9$	$\Delta 34.9$	$\Delta 4.1$	$\Delta 3.9$	$\Delta 1.5$	$\Delta 7.6$	$\Delta 4.0$	$\Delta 3.6$
EATA	$\times$	56.0	56.1	57.1	54.5	54.8	59.6	58.7	61.8	60.1	71.4	75.3	68.5	62.7	69.0	66.5
	$\checkmark$	56.1	56.5	57.2	58.0	57.9	62.6	59.3	65.6	63.5	72.6	78.0	69.5	66.6	72.5	70.5
	Improve	$\Delta 0.1$	$\Delta 0.3$	$\Delta 0.1$	$\Delta 3.4$	$\Delta 3.1$	$\Delta 3.1$	$\Delta 0.6$	$\Delta 3.9$	$\Delta 3.4$	$\Delta 1.3$	$\Delta 2.7$	$\Delta 1.0$	$\Delta 3.9$	$\Delta 3.5$	$\Delta 4.0$
SAR	$\times$	51.9	51.7	52.8	51.5	48.9	55.5	49.5	22.2	46.9	66.2	72.9	65.8	50.9	64.0	62.8
	$\checkmark$	56.4	56.6	57.4	58.3	56.9	62.9	58.3	65.3	64.5	72.7	78.5	69.6	64.0	71.9	69.7
	Improve	$\Delta 4.5$	$\Delta 4.9$	$\Delta 4.6$	$\Delta 6.8$	$\Delta 8.1$	$\Delta 7.4$	$\Delta 8.8$	$\Delta 43.1$	$\Delta 17.5$	$\Delta 6.5$	$\Delta 5.5$	$\Delta 3.8$	$\Delta 13.2$	$\Delta 7.9$	$\Delta 6.9$
CoTTA	$\times$	40.3	37.6	41.7	34.3	28.3	44.0	35.6	38.0	43.0	58.8	70.3	58.4	39.8	58.1	59.9
	$\checkmark$	43.5	40.9	45.5	36.9	29.7	48.1	37.8	40.7	42.0	62.3	73.6	58.9	42.8	63.5	63.8
	Improve	$\Delta 3.1$	$\Delta 3.4$	$\Delta 3.8$	$\Delta 2.6$	$\Delta 1.4$	$\Delta 4.0$	$\Delta 2.2$	$\Delta 2.7$	$\nabla 1.0$	$\Delta 3.5$	$\Delta 3.3$	$\Delta 0.5$	$\Delta 3.0$	$\Delta 5.4$	$\Delta 3.9$
MEMO	$\times$	39.7	36.5	39.8	32.4	25.8	40.3	34.7	27.5	32.8	53.5	66.2	56.0	35.7	55.9	58.2
	$\checkmark$	40.6	37.5	40.6	33.4	26.7	41.2	35.4	28.7	33.7	54.7	67.1	55.9	36.6	57.2	59.3
	Improve	$\Delta 0.8$	$\Delta 1.0$	$\Delta 0.8$	$\Delta 1.0$	$\Delta 0.9$	$\Delta 1.0$	$\Delta 0.7$	$\Delta 1.2$	$\Delta 0.9$	$\Delta 1.2$	$\Delta 0.8$	$\nabla 0.1$	$\Delta 0.9$	$\Delta 1.3$	$\Delta 1.1$

Table 2: Classification and uncertainty estimation comparisons under **open-world** TTA settings, where  $P^{\text{test}}$  is a mixture of both covariate-shifted samples (Gaussian noise of severity level 3) and a suit of diverse abnormal outliers. Additional results with standard deviation are in Appendix C.2.

Method	COME	None		NINCO		iNaturalist		SSB-Hard		Texture		Places	
		Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$
No Adapt	$\times$	64.4	63.7	64.5	69.9	64.4	69.5	64.4	72.5	64.8	65.6	64.3	56.8
PL	$\times$	69.1	62.8	65.6	71.6	68.8	69.5	68.4	75.9	66.1	66.0	66.6	59.7
T3A	$\times$	64.4	71.2	64.3	70.0	64.2	75.0	63.7	80.7	64.4	69.0	63.8	69.5
TEA	$\times$	63.9	63.8	60.5	72.5	62.3	74.6	63.3	79.4	61.0	67.8	61.9	64.5
LAME	$\times$	64.1	64.4	64.1	72.3	64.1	72.4	64.2	74.0	64.7	68.8	64.0	61.3
FOA	$\times$	67.8	61.2	66.4	70.5	67.4	66.1	67.1	75.6	65.8	61.4	66.6	54.1
Tent	$\times$	70.8	63.2	66.2	71.9	69.9	70.7	69.8	77.4	66.4	66.5	68.3	59.9
	$\checkmark$	72.6	64.7	68.9	64.3	72.5	63.7	72.7	70.7	68.4	60.4	70.7	45.9
	Improve	$\Delta 1.7$	$\Delta 1.6$	$\Delta 2.7$	$\nabla 7.6$	$\Delta 2.6$	$\nabla 7.0$	$\Delta 2.9$	$\nabla 6.7$	$\Delta 2.1$	$\nabla 6.1$	$\Delta 2.4$	$\nabla 14.0$
EATA	$\times$	70.3	63.7	66.4	68.6	70.3	71.5	70.0	77.4	67.3	67.1	68.8	61.9
	$\checkmark$	73.4	62.7	70.1	60.5	73.2	63.3	73.0	70.5	70.5	55.8	72.3	45.6
	Improve	$\Delta 3.1$	$\nabla 1.0$	$\Delta 3.7$	$\nabla 8.2$	$\Delta 2.9$	$\nabla 8.2$	$\Delta 3.0$	$\nabla 6.9$	$\Delta 3.2$	$\nabla 11.3$	$\Delta 3.6$	$\nabla 16.3$
SAR	$\times$	69.7	62.3	64.9	71.4	66.9	70.9	67.7	78.1	64.4	64.5	66.1	58.6
	$\checkmark$	73.1	62.9	69.8	66.3	73.2	65.2	73.5	71.8	69.5	59.4	72.3	49.8
	Improve	$\Delta 3.5$	$\Delta 0.6$	$\Delta 4.9$	$\nabla 5.1$	$\Delta 6.3$	$\nabla 5.6$	$\Delta 5.9$	$\nabla 6.3$	$\Delta 5.1$	$\nabla 5.1$	$\Delta 6.3$	$\nabla 8.8$
CoTTA	$\times$	67.6	63.4	65.3	69.7	70.4	69.5	70.3	76.0	65.8	66.2	66.6	59.2
	$\checkmark$	70.5	62.5	66.2	68.8	72.4	73.5	72.2	78.7	66.5	64.7	68.9	55.0
	Improve	$\Delta 2.9$	$\nabla 0.9$	$\Delta 0.9$	$\nabla 0.9$	$\Delta 2.0$	$\Delta 4.0$	$\Delta 2.0$	$\Delta 2.7$	$\Delta 0.7$	$\nabla 1.5$	$\Delta 2.3$	$\nabla 4.2$
MEMO	$\times$	64.8	69.8	64.8	77.5	64.7	71.9	64.8	77.3	65.0	79.3	64.6	71.5
	$\checkmark$	65.2	67.8	65.9	76.4	65.2	70.8	65.3	75.0	65.4	76.7	65.3	67.6
	Improve	$\Delta 0.5$	$\nabla 1.9$	$\Delta 1.1$	$\nabla 1.1$	$\Delta 0.5$	$\nabla 1.1$	$\Delta 0.5$	$\nabla 2.3$	$\Delta 0.4$	$\nabla 2.6$	$\Delta 0.7$	$\nabla 3.9$

the experimental results, it is observed that our COME can consistently improve the performance of existing TTA methods.

**Reliability of uncertainty estimation (Q3).** We visualize the distribution of model confidence, i.e., the maximum predicted class probability<sup>2</sup> in open-world TTA setting, where the covariate-shifted

<sup>2</sup>For our COME, the class probability is calculated according to Eq. 5.



samples is ImageNet-C (Gaussian noise level 3), and outliers are NINCO. As shown in Figure 2, the model confidence of our COME can effectively perceive incorrect predictions, which establishes an distinguishable margin. In contrast to the model confidence of EM which is almost identical and nearly 100% for all test samples, the model confidence of our method can provide more meaningful information with which to differentiate correct and wrong predictions.

Table 3: Classification and uncertainty estimation comparisons under **lifelong** TTA settings. The model is online adapted and the parameters will never be reset, yet the test input distribution might exhibit a continual shift over time.

Methods	COME	Noise			Blur				Weather				Digital			
		Gauss.	Shot	Impul.	Defoc	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elast.	Pixel	JPEG
No Adapt	✗	35.1	32.2	35.9	31.4	25.3	39.4	31.6	24.5	30.1	54.7	64.5	49.0	34.2	53.2	56.5
PL	✗	49.9	53.4	56.7	46.7	46.1	56.6	51.4	52.6	60.2	68.1	77.7	64.5	51.0	69.0	68.8
FOA	✗	46.5	47.0	51.1	44.8	45.5	52.3	48.1	49.7	57.9	68.0	76.4	63.6	51.7	62.2	65.3
Tent	✗	52.4	56.2	58.7	50.9	51.1	57.7	52.7	54.7	60.5	68.4	77.3	64.6	53.8	69.3	68.6
	✓	54.7	59.0	60.3	51.2	53.7	60.6	57.4	64.1	65.8	71.0	78.6	66.9	62.7	71.8	70.5
	Improve	Δ2.3	Δ2.8	Δ1.6	Δ0.4	Δ2.6	Δ2.9	Δ4.7	Δ9.5	Δ5.3	Δ2.6	Δ1.3	Δ2.3	Δ9.0	Δ2.5	Δ1.9
EATA	✗	55.9	59.5	60.9	56.2	59.2	63.0	61.6	65.7	67.5	72.5	78.6	66.8	67.1	72.2	71.7
	✓	57.9	60.6	61.5	57.0	59.8	63.7	62.3	67.2	68.3	73.7	78.8	69.7	68.6	73.1	71.8
	Improve	Δ2.0	Δ1.0	Δ0.6	Δ0.8	Δ0.6	Δ0.8	Δ0.7	Δ1.5	Δ0.8	Δ1.2	Δ0.3	Δ2.9	Δ1.4	Δ0.9	Δ0.1
SAR	✗	52.0	54.8	56.2	50.2	52.3	56.1	52.8	50.8	26.5	0.1	3.1	0.1	0.1	0.1	0.3
	✓	56.0	60.0	61.1	56.4	58.1	62.9	60.4	66.1	67.4	72.2	78.7	68.0	66.0	72.6	70.9
	Improve	Δ4.0	Δ5.3	Δ4.9	Δ6.2	Δ5.8	Δ6.8	Δ7.6	Δ15.3	Δ41.0	Δ72.1	Δ75.6	Δ67.9	Δ65.9	Δ72.5	Δ70.6
CoTTA	✗	40.3	49.2	57.1	39.8	50.4	55.6	48.3	53.1	61.3	63.9	73.3	62.0	56.5	67.5	66.7
	✓	49.7	61.7	64.2	45.7	57.0	59.0	51.1	58.2	63.1	66.0	73.4	62.9	58.0	68.9	68.2
	Improve	Δ9.4	Δ12.4	Δ7.1	Δ5.9	Δ6.6	Δ3.4	Δ2.8	Δ5.1	Δ1.8	Δ2.1	Δ0.1	Δ1.0	Δ1.5	Δ1.3	Δ1.5

**Ablation study.** Finally, we conduct the ablation study on different components in our COME, i.e., with and without the uncertainty constraint in Eq. 6. The experimental results are shown in Table 4, where SL indicates minimizing entropy of the subjective logic opinion and UC means the uncertainty constraint described by Eq. 6. Compared with non-constrained optimization, naively minimizing the entropy of subjective opinion can only slightly improve uncertainty estimation performance. Combining with the uncertainty constraint, the average and worst-case accuracy can be both substantially improved, which indicates the optimal design of our COME.

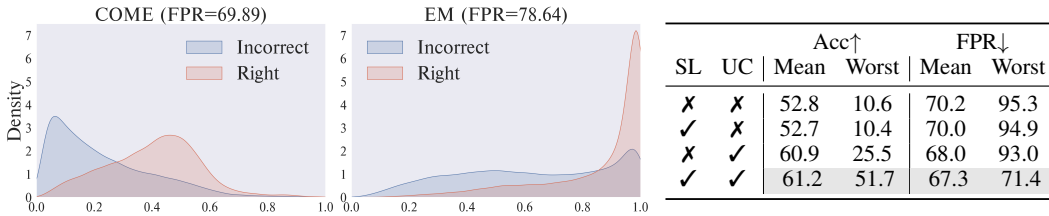


Figure 2: Distribution of model confidence.

SL	UC	Acc↑		FPR↓	
		Mean	Worst	Mean	Worst
✗	✗	52.8	10.6	70.2	95.3
✓	✗	52.7	10.4	70.0	94.9
✗	✓	60.9	25.5	68.0	93.0
✓	✓	61.2	51.7	67.3	71.4

Table 4: Ablation study.

## 6 CONCLUSION

In this paper, we propose a novel learning principle called COME to improve existing TTA methods. Our COME explicitly models the uncertainty raising upon unreliable test samples using the theory of evidence, and then regularizes the model to in favor of conservative prediction confidence during inference time. Our method takes inspiration from Bayesian framework, and consistently outperforms previous EM-based TTA methods on commonly-used benchmarks. The simplicity of the uncertainty regularization used in our implementation is both an advantage and limitation. On the one hand, constraint the uncertainty mass close to the pretrain model is easy-to-deployed and meets the efficiency requirement of TTA. On the other hand, this regularization may be less effective when the pretrain model is also overconfident. We identify this as a limitation of our work. Exploration more effective regularization techniques for better trade-off between the practical requirements of TTA and accurate uncertainty estimation can be a promising future direction.

## REFERENCES

- 486  
487  
488 Taiga Abe, Estefany Kelly Buchanan, Geoff Pleiss, Richard Zemel, and John P Cunningham. Deep  
489 ensembles work, but are they necessary? *Advances in Neural Information Processing Systems*, 35:  
490 33646–33660, 2022.
- 491 Eunsu Baek, Keondo Park, Jiyeon Kim, and Hyung-Sin Kim. Unexplored faces of robustness and  
492 out-of-distribution: Covariate shifts in environment and sensor domains. In *Proceedings of the*  
493 *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22294–22303, 2024.
- 494 Haoyue Bai, Gregory Canal, Xuefeng Du, Jeongyeol Kwon, Robert D Nowak, and Yixuan Li. Feed  
495 two birds with one scone: Exploiting wild data for both out-of-distribution generalization and  
496 detection. In *International Conference on Machine Learning*, 2023.
- 497 Julian Bitterwolf, Maximilian Mueller, and Matthias Hein. In or out? fixing imagenet out-of-  
498 distribution detection evaluation. In *International Conference on Machine Learning*, 2023.
- 500 Malik Boudiaf, Romain Mueller, Ismail Ben Ayed, and Luca Bertinetto. Parameter-free online  
501 test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
502 *Recognition (CVPR)*, pp. 8344–8353, June 2022.
- 503 Dian Chen, Dequan Wang, Trevor Darrell, and Sayna Ebrahimi. Contrastive test-time adaptation.  
504 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
505 295–305, 2022.
- 506 Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describ-  
507 ing textures in the wild. In *Proceedings of the IEEE conference on computer vision and pattern*  
508 *recognition*, 2014.
- 509 John Denker and Yann LeCun. Transforming neural-net output levels to probability distributions.  
510 *Advances in neural information processing systems*, 3, 1990.
- 511 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas  
512 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,  
513 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.  
514 *ICLR*, 2021.
- 515 Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model  
516 uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059.  
517 PMLR, 2016.
- 518 Yarin Gal et al. Uncertainty in deep learning. 2016.
- 519 Zhengqing Gao, Xu-Yao Zhang, and Cheng-Lin Liu. Unified entropy optimization for open-set  
520 test-time adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
521 *Recognition*, pp. 23975–23984, 2024.
- 522 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
523 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- 524 Zongbo Han, Changqing Zhang, Huazhu Fu, and Joey Tianyi Zhou. Trusted multi-view classification  
525 with dynamic evidential fusion. *IEEE transactions on pattern analysis and machine intelligence*,  
526 45(2):2551–2566, 2022.
- 527 Zongbo Han, Yifeng Yang, Changqing Zhang, Linjun Zhang, Joey Tianyi Zhou, Qinghua Hu, and  
528 Huaxiu Yao. Selective learning: Towards robust calibration with dynamic regularization. *arXiv*  
529 *preprint arXiv:2402.08384*, 2024.
- 530 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
531 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,  
532 2016.
- 533 Dan Hendrycks and Thomas G Dietterich. Benchmarking neural network robustness to common  
534 corruptions and surface variations. *International Conference on Learning Representations*, 2019.

- 540 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution  
541 examples in neural networks. *International Conference on Learning Representations*, 2017.
- 542
- 543 Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier  
544 exposure. *International Conference on Learning Representations*, 2019.
- 545
- 546 Dan Hendrycks, Norman Mu, Ekin D. Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshmi-  
547 narayanan. AugMix: A simple data processing method to improve robustness and uncertainty.  
548 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- 549 Hengguan Huang, Xiangming Gu, Hao Wang, Chang Xiao, Hongfu Liu, and Ye Wang. Extrapolative  
550 continuous-time bayesian neural network for fast training-free test-time adaptation. *Advances in*  
551 *Neural Information Processing Systems*, 35:36000–36013, 2022.
- 552
- 553 Yusuke Iwasawa and Yutaka Matsuo. Test-time classifier adjustment module for model-agnostic  
554 domain generalization. *Advances in Neural Information Processing Systems*, 34:2427–2440, 2021.
- 555
- 556 Audun Jsang. *Subjective Logic: A formalism for reasoning under uncertainty*. Springer Publishing  
557 Company, Incorporated, 2018.
- 558
- 559 Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive  
560 uncertainty estimation using deep ensembles. *Advances in neural information processing systems*,  
30, 2017.
- 561
- 562 Jungsoo Lee, Debasmit Das, Jaegul Choo, and Sungha Choi. Towards open-set test-time adaptation  
563 utilizing the wisdom of crowds in entropy minimization. In *Proceedings of the IEEE/CVF*  
*International Conference on Computer Vision*, pp. 16380–16389, 2023.
- 564
- 565 Jian Liang, Lijun Sheng, Zhengbo Wang, Ran He, and Tieniu Tan. Realistic unsupervised clip fine-  
566 tuning with universal entropy optimization. In *Forty-first International Conference on Machine*  
*Learning*.
- 567
- 568 Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source  
569 hypothesis transfer for unsupervised domain adaptation. In *International conference on machine*  
570 *learning*, pp. 6028–6039. PMLR, 2020.
- 571
- 572 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection.  
573 *Advances in neural information processing systems*, 2020.
- 574
- 575 David John Cameron Mackay. *Bayesian methods for adaptive models*. California Institute of  
576 Technology, 1992.
- 577
- 578 Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. *Advances in*  
*neural information processing systems*, 31, 2018.
- 579
- 580 Jooyoung Moon, Jihyo Kim, Younghak Shin, and Sangheum Hwang. Confidence-aware learning for  
581 deep neural networks. In *international conference on machine learning*, pp. 7034–7044. PMLR,  
2020.
- 582
- 583 Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence  
584 predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision*  
*and pattern recognition*, 2015.
- 585
- 586 Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofu Chen, Shijian Zheng, Peilin Zhao, and Mingkui  
587 Tan. Efficient test-time model adaptation without forgetting. In *International conference on*  
*machine learning*, pp. 16888–16905. PMLR, 2022.
- 588
- 589 Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Zhiquan Wen, Yaofu Chen, Peilin Zhao, and Mingkui  
590 Tan. Towards stable test-time adaptation in dynamic wild world. *arXiv preprint arXiv:2302.12400*,  
591 2023.
- 592
- 593 Shuaicheng Niu, Chunyan Miao, Guohao Chen, Pengcheng Wu, and Peilin Zhao. Test-time model  
adaptation with only forward passes. *International conference on machine learning*, 2024.

- 594 Ori Press, Ravid Shwartz-Ziv, Yann LeCun, and Matthias Bethge. The entropy enigma: Success and  
595 failure of entropy minimization. *arXiv preprint arXiv:2405.05012*, 2024.  
596
- 597 Rahul Rahaman et al. Uncertainty quantification and deep ensembles. *Advances in neural information*  
598 *processing systems*, 34:20063–20075, 2021.
- 599 Murat Sensoy, Lance Kaplan, and Melih Kandemir. Evidential deep learning to quantify classification  
600 uncertainty. *Advances in neural information processing systems*, 31, 2018.  
601
- 602 Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a  
603 single deep deterministic neural network. In *International conference on machine learning*, pp.  
604 9690–9700. PMLR, 2020.
- 605 Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam,  
606 Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In  
607 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018.  
608
- 609 Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Open-set recognition: A good  
610 closed-set classifier is all you need? *International Conference on Learning Representations*, 2022.
- 611 Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-  
612 time adaptation by entropy minimization. *International Conference on Learning Representations*,  
613 2021.
- 614 Haoqi Wang, Zhizhong Li, Litong Feng, and Wayne Zhang. Vim: Out-of-distribution with virtual-  
615 logit matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*  
616 *recognition*, pp. 4921–4930, 2022a.
- 617 Meng Wang, Tian Lin, Lianyu Wang, Aidi Lin, Ke Zou, Xinxing Xu, Yi Zhou, Yuanyuan Peng,  
618 Qingquan Meng, Yiming Qian, et al. Uncertainty-inspired open set learning for retinal anomaly  
619 identification. *Nature Communications*, 14(1):6757, 2023.
- 620 Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation.  
621 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
622 7201–7211, 2022b.
- 623 Hongxin Wei, Renchunzi Xie, Hao Cheng, Lei Feng, Bo An, and Yixuan Li. Mitigating neural  
624 network overconfidence with logit normalization. In *International conference on machine learning*,  
625 pp. 23631–23644. PMLR, 2022.
- 626 Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi  
627 Wang, Guangyao Chen, Bo Li, Yiyong Sun, et al. Openood: Benchmarking generalized out-of-  
628 distribution detection. *Advances in Neural Information Processing Systems*, 35:32598–32611,  
629 2022.
- 630 Mouxing Yang, Yunfan Li, Changqing Zhang, Peng Hu, and Xi Peng. Test-time adaptation against  
631 multi-modal reliability bias. In *The Twelfth International Conference on Learning Representations*,  
632 2024.
- 633 Yuan Yige, Xu Bingbing, Hou Liang, Sun Fei, Shen Huawei, and Cheng Xueqi. Tea: Test-time  
634 energy adaptation. In *IEEE/CVF Computer Vision and Pattern Recognition Conference*, 2024.  
635
- 636 Jingyang Zhang, Jingkang Yang, Pengyun Wang, Haoqi Wang, Yueqian Lin, Haoran Zhang, Yiyong  
637 Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Yixuan Li, Ziwei Liu, Yiran Chen, and Hai  
638 Li. Openood v1.5: Enhanced benchmark for out-of-distribution detection. *arXiv preprint*  
639 *arXiv:2306.09301*, 2023.
- 640 Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and  
641 augmentation. *Advances in neural information processing systems*, 35:38629–38642, 2022.  
642
- 643 Aurick Zhou and Sergey Levine. Bayesian adaptation for covariate shift. *Advances in neural*  
644 *information processing systems*, 34:914–927, 2021.  
645  
646  
647

648	APPENDICES	
649		
650	<b>A Proofs</b>	<b>13</b>
651		
652		
653	<b>B Experimental details</b>	<b>15</b>
654	B.1 Datasets . . . . .	15
655	B.2 Implementation details . . . . .	15
656		
657		
658	<b>C Additional results</b>	<b>16</b>
659		
660	C.1 Full results of accuracy with standard deviation (supplementary to Table 1) . . . . .	16
661	C.2 Full results of open-world TTA (supplementary to Table 2) . . . . .	17
662	C.3 Comparison with source-free domain adaption . . . . .	17
663	C.4 Comparison with other other Bayesian methods . . . . .	18
664	C.5 Influence of different hyperparameters . . . . .	18
665	C.6 Comparison on ImageNet-R and ImageNet-S . . . . .	18
666	C.7 In-Distribution performance . . . . .	19
667	C.8 Comparison on ResNet-50 . . . . .	19
668	C.9 Time-consuming comparison . . . . .	19
669	C.10 Mixed distributional shifts performance. . . . .	19
670	C.11 More visualization results (supplementary to Figure 1) . . . . .	19
671		
672		
673		
674		
675		
676		
677	<b>D Discussion</b>	<b>25</b>
678	D.1 Alternative Design Choice . . . . .	25
679	D.2 Limitations and Future work . . . . .	28
680		

## A PROOFS

To proof Lemma 1 and Theorem 1, we need the following lemma firstly.

**Lemma 2.** *Let  $p, q$  be two real numbers. Assuming that  $p \leq q$ , then the  $p$ -norm (also called  $\ell^p$ -norm) and  $q$ -norm of vector  $x = (x_1, \dots, x_n)$  satisfied*

$$\|x\|_p \leq n^{(1/q-1/p)} \|x\|_q. \quad (10)$$

where  $n$  is the length of the vector.

*Proof.* Recall Hölder’s inequality

$$\sum_{i=1}^n |a_i| |b_i| \leq \left( \sum_{i=1}^n |a_i|^r \right)^{1/r} \left( \sum_{i=1}^n |b_i|^{\frac{r}{r-1}} \right)^{1-\frac{1}{r}}. \quad (11)$$

Apply this inequality to the case that  $|a_i| = |x_i|^p$ ,  $|b_i| = 1$  and  $r = q/p \geq 1$ , we can derive to

$$\sum_{i=1}^n |a_i| |b_i| \leq \left( \sum_{i=1}^n ((x_i)^p)^{\frac{q}{p}} \right)^{\frac{p}{q}} \left( \sum_{i=1}^n 1^{\frac{q}{q-p}} \right)^{1-\frac{p}{q}} = \left( \sum_{i=1}^n |x_i|^q \right)^{\frac{p}{q}} n^{1-\frac{p}{q}}. \quad (12)$$

702 Then we have

$$\begin{aligned}
703 \quad \|x\|_p &= \left(\sum_{i=1}^n |x_i|_p\right)^{1/p} \\
704 \quad &\leq \left(\sum_{i=1}^n |x_i|^q\right)^{\frac{1}{q}} n^{1-\frac{1}{q}} \\
705 \quad &= \left(\sum_{i=1}^n |x_i|^q\right)^{\frac{1}{q}} n^{\frac{1}{p}-\frac{1}{q}} \\
706 \quad &= n^{1/p-1/q} \|x\|_q
\end{aligned} \tag{13}$$

714  $\square$

715 Now we proceed to proof our main results.

717 *Proof.* Proof of Lemma 1. Let  $f_{\max} = \max_i f_i(x)$ , then we have

$$718 \quad \exp(f_{\max}) \leq \sum_{i=1}^K \exp(f_i) \leq K \exp(f_{\max}), \tag{14}$$

723 Applying the logarithm to the inequality, then

$$724 \quad f_{\max} \leq \text{LSE}(f) \leq f_{\max} + \log K, \tag{15}$$

726 where LSE is the shorthand of LogSumExp function, i.e.,  $\text{LSE}(x) := \log \sum_{i=1}^K \exp x_i$ .

728 Since we assume that all the elements in logits  $x$  are all positive, then  $f_{\max} = \|f\|_{\infty}$ . Thus combining  
729 with Lemma 1 we can derive that

$$730 \quad K^{-1/p} \|f\|_p \leq \text{LSE}(x) \leq \|f\|_p + \log K, \tag{16}$$

733 Noted that  $u = K/\text{LSE}(f)$ , then we have

$$734 \quad \frac{K}{\|f\|_p + \log K} \leq u \leq \frac{K^{1+1/p}}{\|f\|_p}. \tag{17}$$

737  $\square$

738 *Proof.* Proof of Theorem 1. Assuming that the uncertainty mass  $u$  is constrained as

$$740 \quad u_0 - \delta \leq u \leq u_0 + \delta, \tag{18}$$

742 then the LSE function of model output is also bounded by

$$743 \quad \frac{K}{u_0 + \delta} \leq \text{LSE}(f(x)) \leq \frac{K}{u_0 - \delta}. \tag{19}$$

746 Noted that

$$747 \quad \max f(x) \leq \text{LSE}(f(x)), \tag{20}$$

748 and thus

$$749 \quad \max f(x) \leq \frac{K}{u_0 - \delta}. \tag{21}$$

751 According to Eq. 5, the model confidence is calculated by

$$752 \quad \max_k \mu_k(x) = \frac{\alpha_k}{S} = \frac{\exp f_{\max}}{\sum_{i=1}^K \exp f_i}, \tag{22}$$

755 where  $\alpha_i = \exp f_i(x)$ .

Assuming the  $f_j$  is the largest element in  $f(x)$ , then

$$\begin{aligned} \max_k \mu_k(x) &= \frac{1}{1 + \sum_{i=1, i \neq j}^K \exp(f_i - f_{\max})} \\ &\leq \frac{1}{1 + (K - 1) \exp(f_{\min} - f_{\max})} \\ &\leq \frac{1}{1 + (K - 1) \exp(-\frac{K}{u_0 - \delta})} \end{aligned} \quad (23)$$

□

## B EXPERIMENTAL DETAILS

### B.1 DATASETS

**Covariate-shifted OOD generalization datasets.** We conduct experiments on ImageNet-C (Hendrycks & Dietterich, 2019), which consists of 15 types of algorithmically generated corruptions from noise, blur, weather, and digital categories. Each type of corruption has 5 levels of severity, resulting in 75 distinct corruptions. Besides, we conduct additional experiments on ImageNet-R and ImageNet-S.

**Abnormal outliers for open-world TTA experiments.** We follow the settings of (Zhang et al., 2023), where OpenImage-O (Wang et al., 2022a), SSB-hard (Vaze et al., 2022), Textures (Cimpoi et al., 2014), iNaturalist (Van Horn et al., 2018) and NINCO (Bitterwolf et al., 2023) are selected as outliers for ImageNet. ◦ OpenImage-O contains 17632 manually filtered images and is  $7.8 \times$  larger than the ImageNet dataset. ◦ SSB-hard is selected from ImageNet-21K. It consists of 49K images and 980 categories. ◦ Textures (Describable Textures Dataset, DTD) consists of 5,640 images depicting natural textures. ◦ iNaturalist consists of 859000 images from over 5000 different species of plants and animals. ◦ NINCO consists with a total of 5879 samples of 64 classes which are non-overlapped with ImageNet-C. These datasets are also used in previous open-set OOD generalization works (Gao et al., 2024; Bai et al., 2023; Liang et al.; Lee et al., 2023).

### B.2 IMPLEMENTATION DETAILS

**Pretrained models.** The pretrained ViT model is ViT-Base (Dosovitskiy et al., 2021). The model is trained on the source ImageNet-1K training set and the model weights<sup>3</sup> are directly obtained from timm repository. Specifically, the pretrained ResNet model in C.8 is ResNet-50-BN (He et al., 2016). The model is trained on the source ImageNet-1K training set and the model weights<sup>4</sup> are directly obtained from torchvision library.

**TEA<sup>5</sup> (Yige et al., 2024)** We follow all hyperparameters that are set in TEA unless it does not provide. Specifically, we use SGD as the update rule, with a momentum of 0.9, batch size of 64 and learning rate of 0.001/0.00025 for ViT/ResNet models. The trainable parameters are all affine parameters of layer/batch normalization layers for ViT/ResNet models.

**T3A<sup>6</sup> (Iwasawa & Matsuo, 2021)** We follow all hyperparameters that are set in T3A unless it does not provide. Specifically, the batch size is set to 64. The number of supports to restore M is set to 20 for all experiments.

**LAME<sup>7</sup> (Boudiaf et al., 2022)** We follow all hyperparameters that are set in LAME unless it does not provide. For fair comparison, we maintain a consistent batch size of 64 for LAME, aligning it with the same batch size used by other methods in our evaluation. We use the kNN affinity matrix with the value of  $k = 5$ .

<sup>3</sup><https://huggingface.co/google/vit-base-patch16-224>

<sup>4</sup><https://download.pytorch.org/models/resnet50-19c8e357.pth>

<sup>5</sup><https://github.com/yuanyige/tea>

<sup>6</sup><https://github.com/matsuolab/T3A>

<sup>7</sup><https://github.com/fiveai/LAME>

810 **FOA**<sup>8</sup> (Niu et al., 2024) We follow all hyperparameters that are set in FOA unless it does not provide.  
 811 Specifically, the batch size is set to 64. The number of supports to restore M is set to 20 for all  
 812 experiments.

813 **Tent**<sup>9</sup> (Wang et al., 2021) We follow all hyperparameters that are set in Tent unless it does not  
 814 provide. Specifically, we use SGD as the update rule, with a momentum of 0.9, batch size of 64  
 815 and learning rate of 0.001/0.00025 for ViT/ResNet models. The trainable parameters are all affine  
 816 parameters of layer/batch normalization layers for ViT/ResNet models.

817 **EATA**<sup>10</sup> (Niu et al., 2022) We follow all hyperparameters that are set in EATA. Specifically, the  
 818 entropy constant  $E_0$  (for reliable sample identification) is set to  $0.4 \times \ln 1000$ , where 1000 is the  
 819 number of task classes. The  $\epsilon$  for redundant sample identification is set to 0.05. The trade-off  
 820 parameter  $\beta$  for entropy loss and regularization loss is set to 2,000. The number of pre-collected  
 821 in-distribution test samples for Fisher importance calculation is 2,000. We use SGD as the update  
 822 rule, with a momentum of 0.9, batch size of 64 and learning rate of 0.001/0.00025 for ViT/ResNet  
 823 models. The trainable parameters are all affine parameters of layer/batch normalization layers for  
 824 ViT/ResNet models.

825 **SAR**<sup>11</sup> (Niu et al., 2023) We follow all hyperparameters that are set in SAR. Specifically, the entropy  
 826 threshold  $E_0$  is set to  $0.4 \times \ln 1000$ , where 1000 is the number of task classes. We use SGD as  
 827 the update rule, with a momentum of 0.9, batch size of 64 and learning rate of 0.001/0.00025  
 828 for ViT/ResNet models. For model recovery, we follow all strategy that are set in SAR(except  
 829 for the experiments of life-long).The trainable parameters are all affine parameters of layer/batch  
 830 normalization layers for ViT/ResNet models.

831 **CoTTA**<sup>12</sup> (Wang et al., 2022b) We follow all hyperparameters that are set in CoTTA unless it does  
 832 not provide. Specifically, we use SGD as the update rule, with a momentum of 0.9, batch size of 64  
 833 and learning rate of 0.001/0.01 for ViT/ResNet models. The augmentation threshold  $p_{th}$  is set to  
 834 0.1. For images below threshold, we conduct 32 augmentations including color jitter, random affine,  
 835 Gaussian blur, random horizontal flip, and Gaussian noise. The restoration probability of is set to  
 836 0.01 and the EMA factor  $\alpha$  for teacher update is set to 0.999. The trainable parameters are all affine  
 837 parameters of layer/batch normalization layers for ViT/ResNet models.

838 **MEMO**<sup>13</sup> (Zhang et al., 2022) We follow all hyperparameters that are set in MEMO. Specifically,  
 839 we use the AugMix<sup>14</sup> (Hendrycks et al., 2020) as a set of data augmentations and the augmentation  
 840 size is set to 32. We use SGD as the optimizer,with learning rate 0.00025 and no weight decay. The  
 841 trainable parameters are the entire model.

842 **Source of standard deviation.** For all the experiments, we run multiple times and report the average  
 843 performance and standard deviation. The source of the standard deviation consists 1) the order in  
 844 which the test mini-batches coming online and 2) the randomness of the stochastic optimization  
 845 methods, e.g., SGD, Adam. Since in TTA setting, the model is initialized from the publicly available  
 846 pretrained model weights (i.e., via-base-patch16-224 from timm and resnet50 from PyTorch), there is  
 847 no randomness introduced by model initialization.

## 849 C ADDITIONAL RESULTS

### 851 C.1 FULL RESULTS OF ACCURACY WITH STANDARD DEVIATION (SUPPLEMENTARY TO TABLE 852 1)

853 We provide the full results with standard deviation as supplementary to Table 1 in Table 5. The results  
 854 demonstrate that our COME method consistently achieves better performance than its counterparts.

855 <sup>8</sup><https://github.com/mr-eggplant/FOA>

856 <sup>9</sup><https://github.com/DequanWang/tent>

857 <sup>10</sup><https://github.com/mr-eggplant/EATA>

858 <sup>11</sup><https://github.com/mr-eggplant/SAR>

859 <sup>12</sup><https://github.com/qinenergy/cotta>

860 <sup>13</sup><https://github.com/zhangmarvin/memo>

861 <sup>14</sup><https://github.com/google-research/augmix>



Table 5: Classification accuracy comparison with standard deviation on ImageNet-C (level 5). Substantial ( $\geq 0.5$ ) improvement and degradation compared to the baseline are highlighted in blue or brown respectively.

Methods	COME	Noise			Blur				Weather				Digital			
		Gauss.	Shot	Impul.	Defoc	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elast.	Pixel	JPEG
No Adapt	X	35.1 $\pm$ 0.0	32.2 $\pm$ 0.0	35.9 $\pm$ 0.0	31.4 $\pm$ 0.0	25.3 $\pm$ 0.0	39.4 $\pm$ 0.0	31.6 $\pm$ 0.0	24.5 $\pm$ 0.0	30.1 $\pm$ 0.0	54.7 $\pm$ 0.0	64.5 $\pm$ 0.0	49.0 $\pm$ 0.0	34.2 $\pm$ 0.0	53.2 $\pm$ 0.0	56.5 $\pm$ 0.0
PL	X	49.9 $\pm$ 0.2	48.8 $\pm$ 0.2	50.9 $\pm$ 0.1	48.5 $\pm$ 0.4	41.2 $\pm$ 0.3	52.7 $\pm$ 0.2	42.7 $\pm$ 0.5	24.6 $\pm$ 1.9	42.5 $\pm$ 1.9	63.6 $\pm$ 0.5	73.1 $\pm$ 0.1	65.3 $\pm$ 0.2	44.1 $\pm$ 0.6	63.9 $\pm$ 0.1	62.6 $\pm$ 0.1
T3A	X	34.6 $\pm$ 0.0	31.5 $\pm$ 0.0	35.5 $\pm$ 0.1	32.7 $\pm$ 0.0	27.5 $\pm$ 0.0	40.7 $\pm$ 0.1	33.5 $\pm$ 0.1	25.6 $\pm$ 0.0	30.8 $\pm$ 0.2	56.5 $\pm$ 0.1	64.9 $\pm$ 0.0	50.8 $\pm$ 0.1	38.0 $\pm$ 0.0	54.3 $\pm$ 0.0	58.4 $\pm$ 0.0
TEA	X	44.5 $\pm$ 0.1	39.3 $\pm$ 0.1	45.8 $\pm$ 0.2	37.6 $\pm$ 0.4	35.4 $\pm$ 0.7	46.4 $\pm$ 0.2	31.4 $\pm$ 0.7	8.9 $\pm$ 0.3	46.4 $\pm$ 0.3	60.1 $\pm$ 0.2	72.5 $\pm$ 0.1	59.5 $\pm$ 0.8	45.7 $\pm$ 0.3	62.3 $\pm$ 0.1	58.9 $\pm$ 0.2
LAME	X	34.8 $\pm$ 0.0	31.9 $\pm$ 0.0	35.5 $\pm$ 0.0	30.9 $\pm$ 0.0	24.4 $\pm$ 0.0	38.9 $\pm$ 0.1	30.7 $\pm$ 0.0	23.4 $\pm$ 0.0	29.5 $\pm$ 0.0	53.3 $\pm$ 0.0	64.2 $\pm$ 0.0	41.0 $\pm$ 0.1	32.7 $\pm$ 0.0	52.8 $\pm$ 0.0	56.0 $\pm$ 0.0
FOA	X	46.6 $\pm$ 0.7	43.9 $\pm$ 0.9	48.3 $\pm$ 0.1	47.1 $\pm$ 0.4	40.7 $\pm$ 0.6	49.3 $\pm$ 0.3	43.7 $\pm$ 1.0	53.8 $\pm$ 0.5	52.8 $\pm$ 0.3	64.2 $\pm$ 1.2	76.2 $\pm$ 0.4	63.8 $\pm$ 0.4	48.5 $\pm$ 0.6	62.1 $\pm$ 0.7	64.0 $\pm$ 0.0
Tent	X	52.5 $\pm$ 0.1	52.1 $\pm$ 0.1	53.4 $\pm$ 0.1	52.8 $\pm$ 0.1	47.4 $\pm$ 0.5	56.7 $\pm$ 0.0	47.4 $\pm$ 0.1	10.5 $\pm$ 1.2	26.4 $\pm$ 2.1	67.2 $\pm$ 0.1	74.3 $\pm$ 0.1	67.3 $\pm$ 0.0	50.4 $\pm$ 0.3	66.5 $\pm$ 0.1	64.6 $\pm$ 0.0
	✓	53.9 $\pm$ 0.0	53.9 $\pm$ 0.0	55.2 $\pm$ 0.1	55.8 $\pm$ 0.1	51.8 $\pm$ 0.1	59.8 $\pm$ 0.0	52.6 $\pm$ 0.0	58.4 $\pm$ 0.5	61.3 $\pm$ 0.1	71.3 $\pm$ 0.1	78.2 $\pm$ 0.0	68.8 $\pm$ 0.1	57.9 $\pm$ 0.4	70.5 $\pm$ 0.1	68.2 $\pm$ 0.1
	Improve	$\Delta$ 1.4	$\Delta$ 1.9	$\Delta$ 1.8	$\Delta$ 2.9	$\Delta$ 4.4	$\Delta$ 3.2	$\Delta$ 5.2	$\Delta$ 47.9	$\Delta$ 34.9	$\Delta$ 4.1	$\Delta$ 3.9	$\Delta$ 1.5	$\Delta$ 7.6	$\Delta$ 4.0	$\Delta$ 3.6
EATA	X	56.0 $\pm$ 0.2	56.1 $\pm$ 0.3	57.1 $\pm$ 0.1	54.5 $\pm$ 1.9	54.8 $\pm$ 1.7	59.6 $\pm$ 1.6	58.7 $\pm$ 0.1	61.8 $\pm$ 0.3	60.1 $\pm$ 0.1	71.4 $\pm$ 0.2	75.3 $\pm$ 0.0	68.5 $\pm$ 0.1	62.7 $\pm$ 0.2	69.0 $\pm$ 0.3	66.5 $\pm$ 0.2
	✓	56.1 $\pm$ 0.1	56.5 $\pm$ 0.1	57.2 $\pm$ 0.0	58.0 $\pm$ 0.1	57.9 $\pm$ 0.2	62.6 $\pm$ 0.1	59.3 $\pm$ 0.2	65.6 $\pm$ 0.2	63.5 $\pm$ 0.3	72.6 $\pm$ 0.1	78.0 $\pm$ 0.1	69.5 $\pm$ 0.1	66.6 $\pm$ 0.3	72.5 $\pm$ 0.2	70.5 $\pm$ 0.1
	Improve	$\Delta$ 0.1	$\Delta$ 0.3	$\Delta$ 0.1	$\Delta$ 3.4	$\Delta$ 3.1	$\Delta$ 3.1	$\Delta$ 0.6	$\Delta$ 3.9	$\Delta$ 3.4	$\Delta$ 1.3	$\Delta$ 2.7	$\Delta$ 1.0	$\Delta$ 3.9	$\Delta$ 3.5	$\Delta$ 4.0
SAR	X	51.9 $\pm$ 0.1	51.7 $\pm$ 0.1	52.8 $\pm$ 0.1	51.5 $\pm$ 0.5	48.9 $\pm$ 0.2	55.5 $\pm$ 0.1	49.5 $\pm$ 0.2	22.2 $\pm$ 0.8	46.9 $\pm$ 2.2	66.2 $\pm$ 0.4	72.9 $\pm$ 0.1	65.8 $\pm$ 0.1	50.9 $\pm$ 0.5	64.0 $\pm$ 0.0	62.8 $\pm$ 0.0
	✓	56.4 $\pm$ 0.1	56.6 $\pm$ 0.1	57.4 $\pm$ 0.1	58.3 $\pm$ 0.2	56.9 $\pm$ 0.1	62.9 $\pm$ 0.1	58.3 $\pm$ 0.2	65.3 $\pm$ 0.1	64.5 $\pm$ 0.1	72.7 $\pm$ 0.1	78.5 $\pm$ 0.0	69.6 $\pm$ 0.0	64.0 $\pm$ 0.2	71.9 $\pm$ 0.1	69.7 $\pm$ 0.0
	Improve	$\Delta$ 4.5	$\Delta$ 4.9	$\Delta$ 4.6	$\Delta$ 6.8	$\Delta$ 8.1	$\Delta$ 7.4	$\Delta$ 8.8	$\Delta$ 43.1	$\Delta$ 17.5	$\Delta$ 6.5	$\Delta$ 5.5	$\Delta$ 3.8	$\Delta$ 13.2	$\Delta$ 7.9	$\Delta$ 6.9
COTTA	X	40.3 $\pm$ 0.2	37.6 $\pm$ 0.1	41.7 $\pm$ 0.1	34.3 $\pm$ 0.4	28.3 $\pm$ 0.7	44.0 $\pm$ 0.1	35.6 $\pm$ 0.2	38.0 $\pm$ 0.1	43.0 $\pm$ 0.2	58.8 $\pm$ 0.3	70.3 $\pm$ 0.3	58.4 $\pm$ 0.5	39.8 $\pm$ 0.2	58.1 $\pm$ 0.2	59.9 $\pm$ 0.1
	✓	43.1 $\pm$ 0.0	40.9 $\pm$ 0.3	45.5 $\pm$ 0.4	36.9 $\pm$ 0.2	29.7 $\pm$ 0.2	48.1 $\pm$ 1.2	37.8 $\pm$ 0.3	40.7 $\pm$ 1.4	42.0 $\pm$ 0.4	62.3 $\pm$ 0.6	73.6 $\pm$ 0.2	58.9 $\pm$ 2.4	42.8 $\pm$ 0.3	63.5 $\pm$ 0.2	63.8 $\pm$ 0.1
	Improve	$\Delta$ 3.1	$\Delta$ 3.4	$\Delta$ 3.8	$\Delta$ 2.6	$\Delta$ 1.4	$\Delta$ 4.0	$\Delta$ 2.2	$\Delta$ 2.7	$\nabla$ 1.0	$\Delta$ 3.5	$\Delta$ 3.3	$\Delta$ 0.5	$\Delta$ 3.0	$\Delta$ 5.4	$\Delta$ 3.9

## C.2 FULL RESULTS OF OPEN-WORLD TTA (SUPPLEMENTARY TO TABLE 2)

We provide the full results with standard deviation as supplementary to Table 2 in Table 6. The results demonstrate that our COME method consistently achieves better performance than its counterparts.

Table 6: Classification and uncertainty estimation comparisons with standard deviation under **open-world** TTA settings, where  $P^{\text{test}}$  is a mixture of both covariate-shifted samples (Gaussian noise of severity level 3).

Method	COME	None		NINCO		iNaturalist		SSB-Hard		Texture		Places	
		Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$
No Adapt	X	64.4 $\pm$ 0.0	63.7 $\pm$ 0.0	64.6 $\pm$ 0.3	70.3 $\pm$ 1.0	64.4 $\pm$ 0.0	69.5 $\pm$ 0.0	64.4 $\pm$ 0.0	72.5 $\pm$ 0.1	64.7 $\pm$ 0.2	65.0 $\pm$ 0.4	64.5 $\pm$ 0.3	57.4 $\pm$ 0.6
PL	X	69.1 $\pm$ 0.1	63.0 $\pm$ 0.3	65.6 $\pm$ 0.0	70.9 $\pm$ 1.7	68.8 $\pm$ 0.1	69.3 $\pm$ 0.2	68.3 $\pm$ 0.0	76.2 $\pm$ 0.3	66.1 $\pm$ 0.1	65.7 $\pm$ 0.3	67.1 $\pm$ 0.3	59.1 $\pm$ 0.6
T3A	X	64.4 $\pm$ 0.0	71.0 $\pm$ 0.3	64.3 $\pm$ 0.2	71.4 $\pm$ 1.6	64.1 $\pm$ 0.1	74.1 $\pm$ 0.6	63.7 $\pm$ 0.0	80.1 $\pm$ 0.4	64.4 $\pm$ 0.2	67.4 $\pm$ 1.4	64.0 $\pm$ 0.3	69.6 $\pm$ 0.1
TEA	X	64.1 $\pm$ 0.3	63.3 $\pm$ 0.5	60.2 $\pm$ 0.6	72.9 $\pm$ 0.3	62.4 $\pm$ 0.1	74.5 $\pm$ 0.2	63.5 $\pm$ 0.1	79.3 $\pm$ 0.2	60.6 $\pm$ 0.4	69.0 $\pm$ 0.9	62.0 $\pm$ 0.2	65.8 $\pm$ 1.0
LAME	X	64.1 $\pm$ 0.0	64.1 $\pm$ 0.3	64.2 $\pm$ 0.3	72.3 $\pm$ 0.1	64.1 $\pm$ 0.0	72.2 $\pm$ 0.2	64.1 $\pm$ 0.0	73.9 $\pm$ 0.1	64.6 $\pm$ 0.2	68.5 $\pm$ 0.3	64.3 $\pm$ 0.3	61.4 $\pm$ 0.3
FOA	X	67.8 $\pm$ 0.0	62.0 $\pm$ 0.6	65.9 $\pm$ 0.3	69.2 $\pm$ 0.9	67.6 $\pm$ 0.2	66.8 $\pm$ 0.5	67.5 $\pm$ 0.2	76.0 $\pm$ 0.3	65.8 $\pm$ 0.0	60.6 $\pm$ 0.5	66.7 $\pm$ 0.1	52.8 $\pm$ 0.9
Tent	X	70.9 $\pm$ 0.1	63.6 $\pm$ 0.5	66.4 $\pm$ 0.3	71.5 $\pm$ 0.3	70.0 $\pm$ 0.1	70.8 $\pm$ 0.4	69.8 $\pm$ 0.0	77.4 $\pm$ 0.1	66.7 $\pm$ 0.4	67.3 $\pm$ 0.7	68.0 $\pm$ 0.4	60.8 $\pm$ 0.8
	✓	72.6 $\pm$ 0.0	64.9 $\pm$ 0.1	69.0 $\pm$ 0.1	63.9 $\pm$ 0.3	72.6 $\pm$ 0.1	63.6 $\pm$ 0.1	72.7 $\pm$ 0.0	70.6 $\pm$ 0.0	68.5 $\pm$ 0.0	60.5 $\pm$ 0.0	70.7 $\pm$ 0.0	46.3 $\pm$ 0.2
	Improve	$\Delta$ 1.7	$\Delta$ 1.3	$\Delta$ 2.6	$\nabla$ 7.6	$\Delta$ 2.6	$\nabla$ 7.2	$\Delta$ 2.8	$\nabla$ 6.8	$\Delta$ 1.8	$\nabla$ 6.8	$\Delta$ 2.6	$\nabla$ 14.5
EATA	X	70.2 $\pm$ 0.1	63.5 $\pm$ 0.1	66.2 $\pm$ 0.1	68.8 $\pm$ 0.2	70.2 $\pm$ 0.1	71.5 $\pm$ 0.1	69.9 $\pm$ 0.1	77.5 $\pm$ 0.1	67.6 $\pm$ 0.5	66.5 $\pm$ 0.5	68.6 $\pm$ 0.2	62.0 $\pm$ 0.2
	✓	73.3 $\pm$ 0.0	63.1 $\pm$ 0.3	70.1 $\pm$ 0.0	60.5 $\pm$ 0.0	73.2 $\pm$ 0.0	63.2 $\pm$ 0.1	73.0 $\pm$ 0.0	70.6 $\pm$ 0.1	70.3 $\pm$ 0.2	56.8 $\pm$ 0.7	72.3 $\pm$ 0.0	45.8 $\pm$ 0.2
	Improve	$\Delta$ 3.1	$\nabla$ 0.4	$\Delta$ 3.9	$\nabla$ 8.2	$\Delta$ 3.0	$\nabla$ 8.3	$\Delta$ 3.1	$\nabla$ 6.9	$\Delta$ 2.7	$\nabla$ 9.7	$\Delta$ 3.7	$\nabla$ 16.2
SAR	X	69.6 $\pm$ 0.1	62.3 $\pm$ 0.1	65.5 $\pm$ 0.5	70.3 $\pm$ 0.9	67.1 $\pm$ 0.1	70.7 $\pm$ 0.2	67.6 $\pm$ 0.1	77.9 $\pm$ 0.2	65.5 $\pm$ 0.8	65.8 $\pm$ 1.0	66.2 $\pm$ 0.1	59.8 $\pm$ 1.0
	✓	73.1 $\pm$ 0.0	62.3 $\pm$ 0.4	70.0 $\pm$ 0.2	65.9 $\pm$ 0.2	73.3 $\pm$ 0.1	65.3 $\pm$ 0.0	73.9 $\pm$ 0.0	71.5 $\pm$ 0.3	69.4 $\pm$ 0.0	59.3 $\pm$ 0.0	72.0 $\pm$ 0.2	48.5 $\pm$ 0.9
	Improve	$\Delta$ 3.5	$\Delta$ 0.1	$\Delta$ 4.5	$\nabla$ 4.3	$\Delta$ 6.2	$\nabla$ 5.4	$\Delta$ 6.0	$\nabla$ 6.5	$\Delta$ 3.9	$\nabla$ 6.5	$\Delta$ 5.8	$\nabla$ 11.3
COTTA	X	67.7 $\pm$ 0.1	63.4 $\pm$ 0.1	65.3 $\pm$ 0.4	70.4 $\pm$ 0.8	70.4 $\pm$ 0.2	69.5 $\pm$ 0.1	70.1 $\pm$ 0.1	75.8 $\pm$ 0.4	65.7 $\pm$ 0.2	65.0 $\pm$ 0.9	67.0 $\pm$ 0.3	58.8 $\pm$ 0.3
	✓	70.4 $\pm$ 0.1	62.8 $\pm$ 0.2	66.5 $\pm$ 0.3	68.0 $\pm$ 0.6	72.4 $\pm$ 0.0	72.9 $\pm$ 0.5	72.1 $\pm$ 0.1	78.5 $\pm$ 0.2	66.4 $\pm$ 0.2	63.8 $\pm$ 0.7	68.9 $\pm$ 0.4	54.5 $\pm$ 0.5
	Improve	$\Delta$ 2.8	$\nabla$ 0.6	$\Delta$ 1.2	$\nabla$ 2.4	$\Delta$ 2.0	$\Delta$ 3.5	$\Delta$ 2.0	$\Delta$ 2.7	$\Delta$ 0.7	$\nabla$ 1.2	$\Delta$ 1.9	$\nabla$ 4.4

## C.3 COMPARISON WITH SOURCE-FREE DOMAIN ADAPTION

There is a strong connection between Source-Free Domain Adaptation (SFDA) and Test-Time Adaptation (TTA). The primary difference is that TTA focuses on **online** adjusting during the testing. On the other hand, SFDA approaches generally perform **offline**. That is, the inference is deferred until the optimization is done. In contrast, our TTA method can achieve adaption and inference at the same time.

To further validate the applicability of our method, we report the classification accuracy on ImageNet-C Gaussian noise level 5 under **source-free domain adaption settings**. The results are in Table 7. The baselines we considered include pseudo label (PL), mutual information maximization (IM), and entropy minimization (EM) following (Liang et al., 2020). “-” means the model accuracy collapses to random guess level.

Table 7: Classification accuracy comparison on ImageNet-C Gaussian noise (level 5) under source-free domain adaption settings. “-” means the classification accuracy collapses to random guess level. The accuracy of the original pretrained model is 35.1.

EPOCH	PL	EM	IM	TENT	EATA	SAR
COME	✗	✗	✗	✓	✓	✓
1	34.2 $\pm$ 3.2	-	60.7 $\pm$ 0.0	66.5 $\pm$ 0.0	66.4 $\pm$ 0.3	68.5 $\pm$ 0.1
2	-	-	63.8 $\pm$ 0.1	68.4 $\pm$ 0.0	67.2 $\pm$ 0.2	69.7 $\pm$ 0.1
3	-	-	65.1 $\pm$ 0.1	69.2 $\pm$ 0.0	67.8 $\pm$ 0.2	70.2 $\pm$ 0.1

#### C.4 COMPARISON WITH OTHER OTHER BAYESIAN METHODS

There exists a few Bayesian inspired TTA methods closely related to our method. (Zhou & Levine, 2021) explores Bayesian model ensembling Zhou & Levine (2021) for TTA, which introduces noticeable inference latency. Since (Zhou & Levine, 2021) has not made their source code publicly available. For a fair comparison, we implement the proposed COME using the same backbone (ResNet50-v2) and dataset (ImageNet-C) as in (Zhou & Levine, 2021) and report the average accuracy of all corruptions and levels. The results of (Zhou & Levine, 2021) are directly copied from the original paper. The results are in Table 8.

Table 8: Classification accuracy comparison with other bayesian inspired TTA methods on ImageNet-C under TTA settings averaged over all corruption types and levels. The result of BACS is copied from the original paper.

	No Adapt	BN	BACS	TENT	TENT	EATA
COME	✗	✗	✗	✗	✓	✓
	47.3 $\pm$ 0.0	47.3 $\pm$ 0.0	56.1	48.9 $\pm$ 0.0	51.1 $\pm$ 0.0	58.1 $\pm$ 0.0

#### C.5 INFLUENCE OF DIFFERENT HYPERPARAMETERS

We conduct additional experiments to investigate the influence of different hyperparameters. The results are in Table 9 and Table 10. Our COME generally outperforms EM with moderate hyperparameters. For  $\tau$ , as we mentioned before, it actually controls the magnitude of uncertainty mass. If we have some prior knowledge that most test samples should be rejected to adapt during TTA, we should choose a relatively small  $\tau$  and making the model confidence more conservative in this circumstance. As for  $\delta$  (or  $p$ ) which represents the tolerance of uncertainty divergence, it should be selected per need by the user via trial and error: if users are extremely cautious about unreliable TTA,  $\delta$  should be tuned down and  $p$  should be tuned up; otherwise, if a better performance is required. Besides, we introduce an additional hyperparameter  $\lambda$  to the learning objective. The experiments results in Table 9 shows that an additional performance improvement can be observed when  $\lambda \in [1, 100]$ .

Table 9: Additional results with different  $\lambda$ . We report the accuracy when our COME is equipped to Tent on ImageNet-C Gaussian noise level 5 with different  $\lambda$ . The accuracy of the original Tent using entropy minimization is 52.6.

$\lambda$	0.10	1.00	10.0	30.0	50.0	80.0	100	150
Acc.	53.7 $\pm$ 0.1	53.9 $\pm$ 0.0	54.6 $\pm$ 0.1	55.3 $\pm$ 0.1	55.8 $\pm$ 0.1	56.2 $\pm$ 0.1	56.1 $\pm$ 0.1	10.6 $\pm$ 0.3

#### C.6 COMPARISON ON IMAGENET-R AND IMAGENET-S

We introduce two more additional datasets, i.e., ImageNet-R and ImageNet-S, to further evaluate the proposed method on commonly used OOD generalization benchmarks. The results are in Table 11 and Table 12.

Table 10: Additional results with different hyperparameters. We report the accuracy on ImageNet-C Gaussian noise level 5 with different hyperparameters. We implement our COME with Tent. The accuracy of the original Tent using entropy minimization is 52.6.

	$p = 1$	$p = 2$	$p = 3$	$p = \infty$
$\tau = 0.5$	37.8 $\pm$ 0.0	38.5 $\pm$ 0.0	39.1 $\pm$ 0.0	41.1 $\pm$ 0.0
$\tau = 1.0$	53.5 $\pm$ 0.0	53.8 $\pm$ 0.1	53.2 $\pm$ 0.0	47.2 $\pm$ 0.1
$\tau = 1.2$	54.7 $\pm$ 0.0	54.7 $\pm$ 0.1	54.0 $\pm$ 0.1	48.9 $\pm$ 0.0
$\tau = 1.5$	54.2 $\pm$ 0.2	54.2 $\pm$ 0.1	53.2 $\pm$ 0.1	49.1 $\pm$ 0.1

Table 11: Additional results on ImageNet-R

	Tent	EATA	SAR	COTTA
EM	37.73 $\pm$ 0.03	36.11 $\pm$ 0.12	51.90 $\pm$ 0.35	36.04 $\pm$ 0.06
COME	39.05 $\pm$ 0.05	38.22 $\pm$ 0.23	56.40 $\pm$ 0.18	37.39 $\pm$ 0.08

### C.7 IN-DISTRIBUTION PERFORMANCE

We compare the in-distribution performance of proposed COME to EM-based methods. As shown in Table 13, our method consistently outperforms entropy minimization.

### C.8 COMPARISON ON RESNET-50

As shown in previous work (Niu et al., 2023), the TTA performance can be influenced by different model architectures, especially the type of normalization layers, i.e., batch normalization, group normalization, layer normalization, and instance normalization. To further evaluate the proposed method, we conduct additional experiments on ResNet-50 with batch normalization layers under open-world TTA settings. The experimental results in Table 14 show that our COME still achieves superior performance compared to entropy minimization learning principle when equipped to Tent.

### C.9 TIME-CONSUMING COMPARISON

We compare the time-cost of proposed COME to EM-based methods and Non-EM based methods in Table 15. We run all the experiments on one single NVIDIA 4090 GPU. Our COME does not introduce noticeably extra cost of computation.

### C.10 MIXED DISTRIBUTIONAL SHIFTS PERFORMANCE.

We evaluate the proposed COME in two additional settings introduced by (Niu et al., 2023). These scenarios include 1) online imbalanced label distribution shifts, where the test data are sorted by class, and 2) mixed domain shifts, where the test data stream includes several randomly mixed domains with different distribution shifts. As shown in Table 16 and Table 17, our COME consistently outperforms entropy minimization with an exception of a slightly suboptimal uncertainty estimation performance compared to CoTTA.

### C.11 MORE VISUALIZATION RESULTS (SUPPLEMENTARY TO FIGURE 1)

We provide more visualization results on two representative TTA methods, i.e., the seminal Tent (Wang et al., 2021) and recent SOTA SAR (Niu et al., 2023). We observe that our COME enjoys a more stable TTA progress with less risk of model collapse and overconfidence across various types of corruption. We test on ImageNet-C level 5.

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

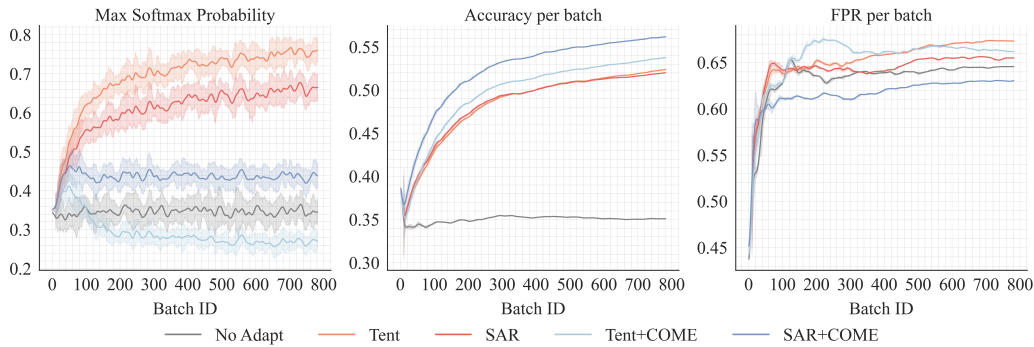


Figure 3: Comparison on two representative TTA methods on ImageNet-C under **Gaussian Noise** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

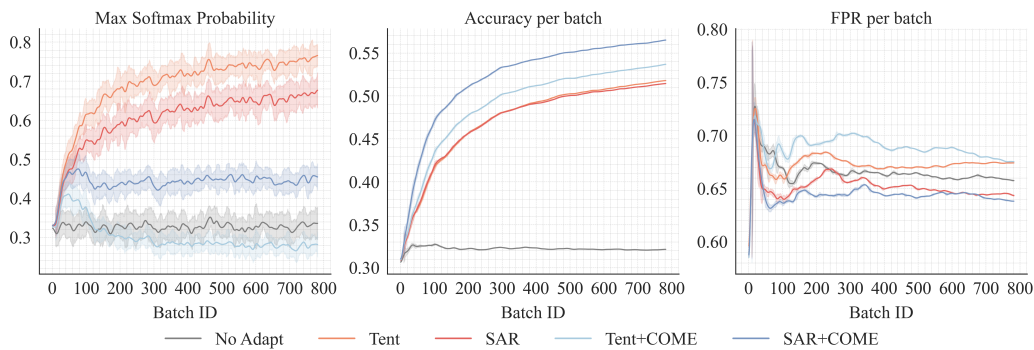


Figure 4: Comparison on two representative TTA methods on ImageNet-C under **Shot Noise** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

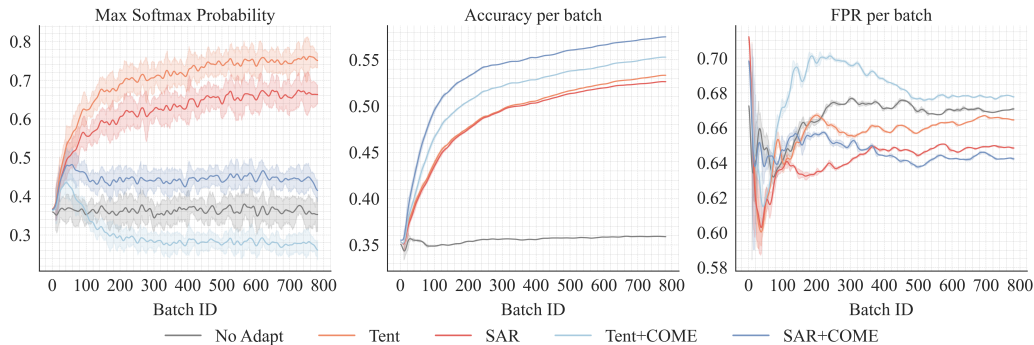


Figure 5: Comparison on two representative TTA methods on ImageNet-C under **Impulse Noise** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

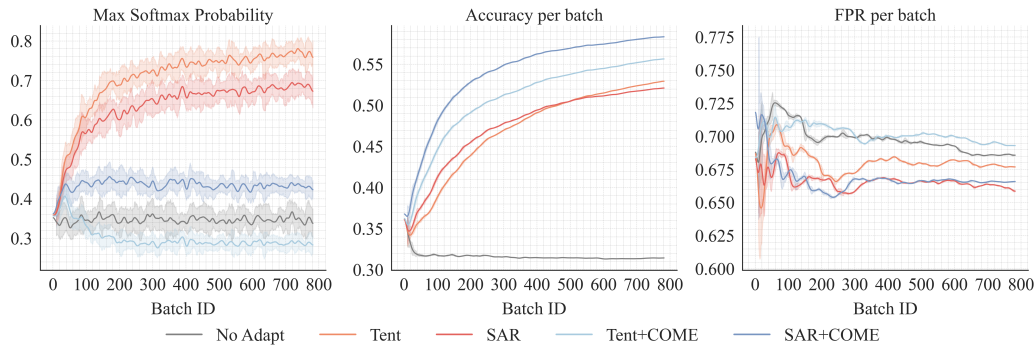


Figure 6: Comparison on two representative TTA methods on ImageNet-C under **Defocus Blur** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

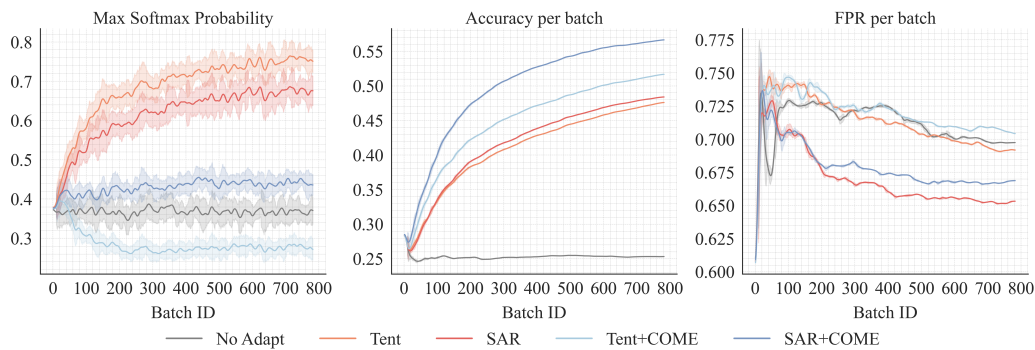


Figure 7: Comparison on two representative TTA methods on ImageNet-C under **Glass Blur** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

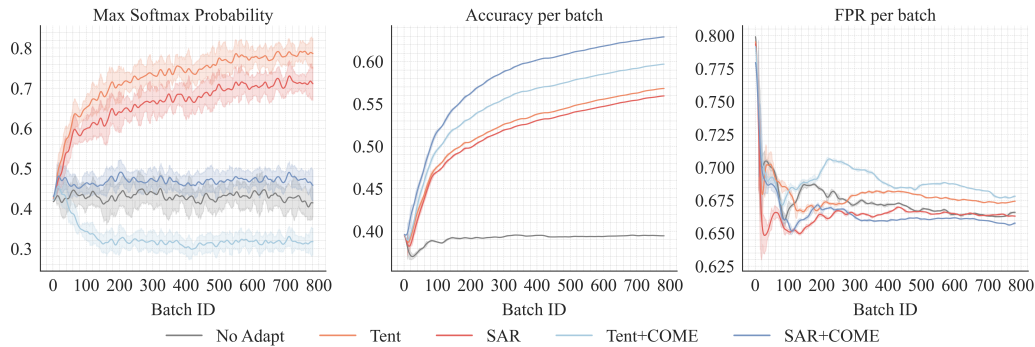


Figure 8: Comparison on two representative TTA methods on ImageNet-C under **Motion Blur** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

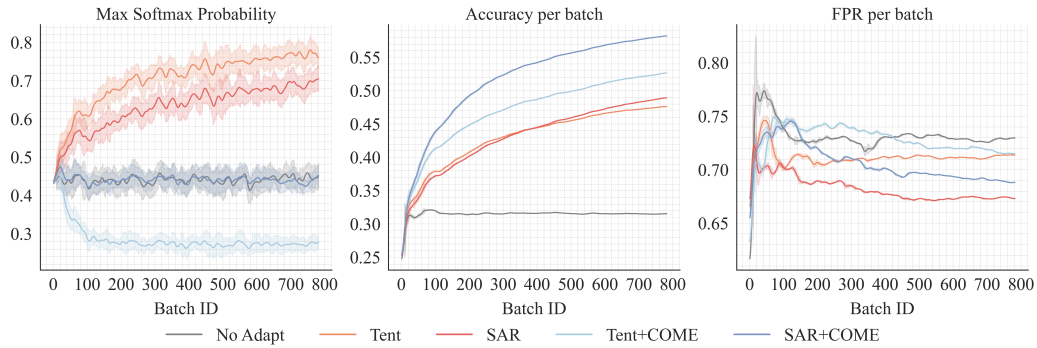
1130

1131

1132

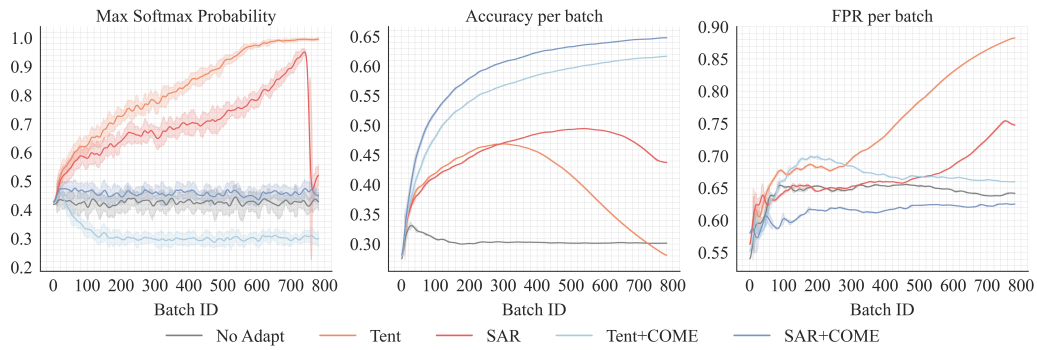
1133

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147



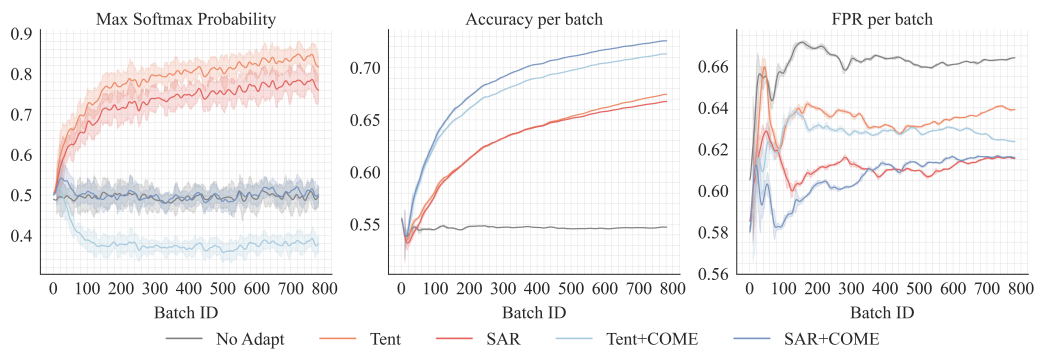
1148 Figure 9: Comparison on two representative TTA methods on ImageNet-C under **Zoom Blur**  
1149 corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with  
1150 consistently improved classification accuracy and false positive rate.

1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164



1165 Figure 10: Comparison on two representative TTA methods on ImageNet-C under **Frost**  
1166 corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently  
1167 improved classification accuracy and false positive rate. Although the SAR method can recover  
1168 the model when it collapses to a trivial solution, its performance remains poor. Our COME method  
1169 addresses the issue of overconfidence that leads to model collapse.

1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184



1185 Figure 11: Comparison on two representative TTA methods on ImageNet-C under **Fog** corruption of  
1186 severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently  
1187 improved classification accuracy and false positive rate.

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

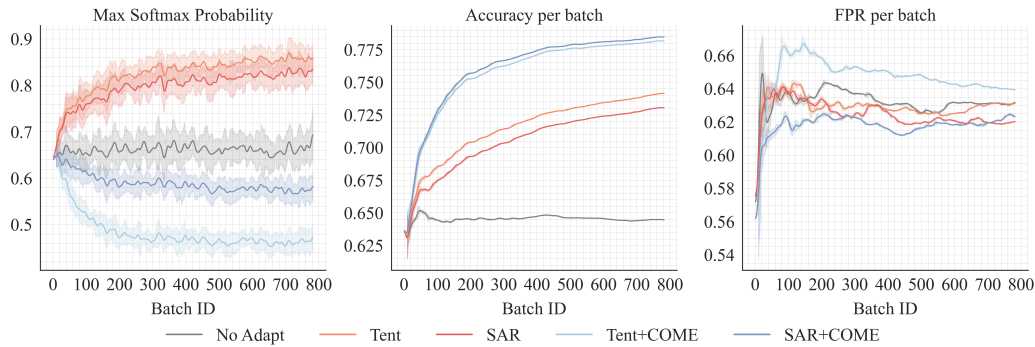


Figure 12: Comparison on two representative TTA methods on ImageNet-C under **Brightness** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

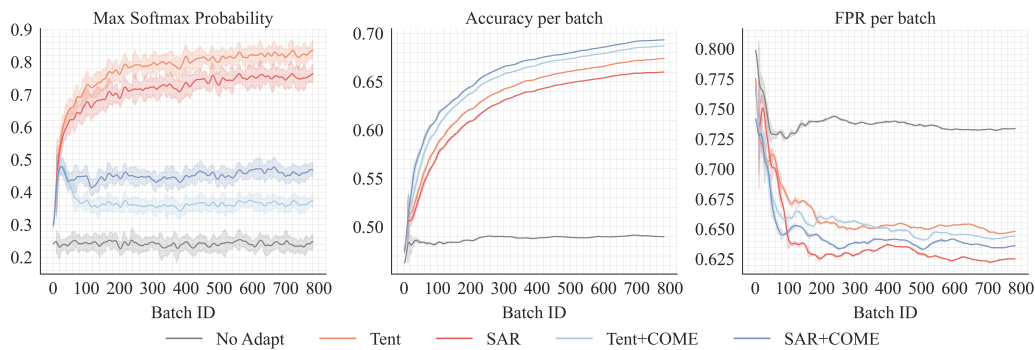


Figure 13: Comparison on two representative TTA methods on ImageNet-C under **Contrast** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

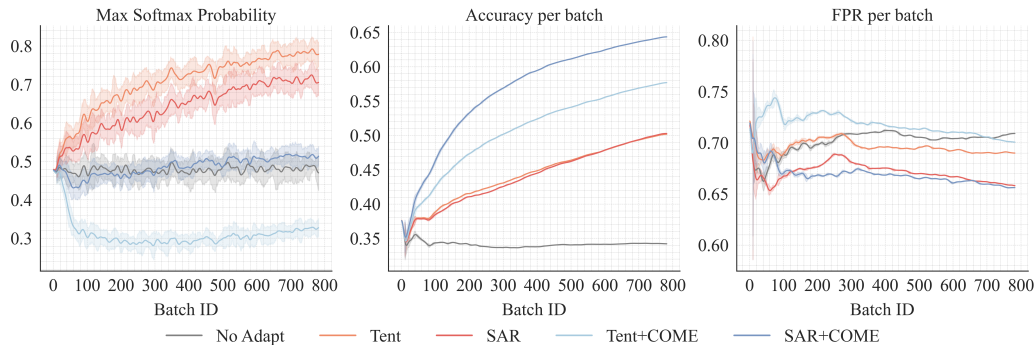


Figure 14: Comparison on two representative TTA methods on ImageNet-C under **Elastic Transform** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

1241

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

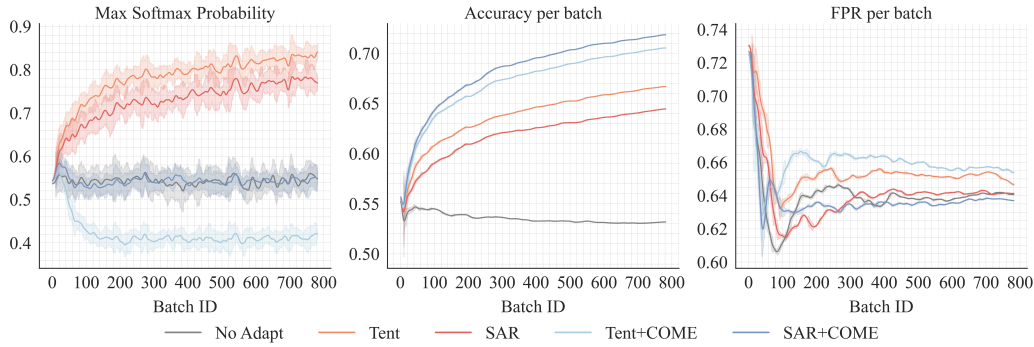


Figure 15: Comparison on two representative TTA methods on ImageNet-C under **Pixelate** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.

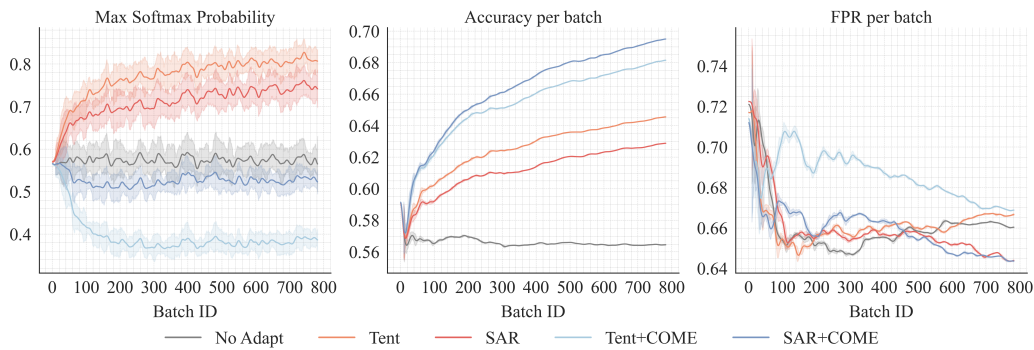


Figure 16: Comparison on two representative TTA methods on ImageNet-C under **Jpeg Compression** corruption of severity level 5. By contrast to EM, our COME establishes a stable TTA process with consistently improved classification accuracy and false positive rate.



Table 12: Additional results on ImageNet-S.

	Tent	EATA	SAR	COTTA
EM	31.63 ± 0.07	36.11 ± 0.19	33.92 ± 0.69	30.84 ± 0.35
COME	39.22 ± 0.04	39.22 ± 0.32	43.52 ± 0.52	33.82 ± 0.73

Table 13: Comparison w.r.t. **in-distribution performance**, *i.e.*, on clean/original ImageNet validation set, with ViT as the base model. Substantial ( $\geq 0.5$ ) **improvement** and **degradation** compared to the baseline are highlighted in blue or red respectively.

	TENT	SAR	EATA	CoTTA	MEMO	Avg.
COME	Acc↑	Acc↑	Acc↑	Acc↑	Acc↑	Acc↑
$\times$	81.4	80.7	81.3	82.1	80.3	81.2
$\checkmark$	83.1	83.1	83.1	82.8	80.6	82.6
Improve	$\Delta 1.7$	$\Delta 2.5$	$\Delta 1.8$	$\Delta 0.7$	-	$\Delta 1.4$

## D DISCUSSION

### D.1 ALTERNATIVE DESIGN CHOICE

**Choices of transformation function to obtain the opinion.** By definitions, the parameters of a Dirichlet distribution  $\alpha$  must be greater than 1 and the evidence  $e$  should be non-negative. This can be achieved by applying ReLU activation function or exponential function to the output logits as suggested in previous works (Han et al., 2022; Malinin & Gales, 2018). That is, we can get the evidence by

$$e = \text{ReLU}(f(x)) \quad (24)$$

or

$$e = \exp f(x) - 1. \quad (25)$$

In this paper, we choose the exponential function. Since we assume the pretrain model is trained with standard cross-entropy loss, using exponential function to get the evidence can keep the training strategy unchanged. Besides, based on our early empirical findings, using exponential function can achieve better classification performance compared to ReLU.

We refer interested readers to (Malinin & Gales, 2018) and Gal’s PhD Thesis (Gal et al., 2016) for more detailed implementation instructions and math deviations.

**Choices of uncertainty constraint.** In Lemma 1, we prove that by constraining on the model output logits, we can control the uncertainty mass  $u$  not to diverge too far from the pretrained model. Previous work (Wei et al., 2022) proposes to mitigate the overconfidence issue by normalizing the logits during pretrain progress in supervised learning tasks. Following their implementation, we propose to optimize on the direction vector of  $f(x)$ , *i.e.*,  $f(x)/\|f(x)\|_p$ , and thus we can expect that the optimization progress is not related to the magnitude of  $f(x)$ , *i.e.*, its norm. Different from Wei et al. (2022), we recover the magnitude by multiplying the direction vector with its norm (detached), rather than a constant to avoid an additional hyperparameter. However, the uncertainty estimated by pretrain model may not be ideal. However, please kindly remind that in fully TTA task, we can only access unlabeled test data coming online and the inference efficiency matters. Thus traditional methods devised for handling overconfidence like calibration (Guo et al., 2017), ensembling (Zhou & Levine, 2021), BNNs (Huang et al., 2022) and other Bayesian methods like dropout (Gal & Ghahramani, 2016) are not applicable. The only practically available choice is to explore the uncertainty information contained within the model itself. As shown in previous works, while the softmax probability of pretrained model tend to be overconfident, subjective logic is much more reliable (Sensoy et al., 2018; Malinin & Gales, 2018), which can support the proposed regularization. Exploring more effective and efficient regularization is an interesting future research direction.

**Choices of  $p$ -norm.** The tightness for the upper and lower bounds in Lemma 1 is determined by the choice of  $p$ . By considering the simple model where  $f(x)$  outputs the same logits for all classes, the ratio between the upper and lower bound is minimized by  $p = \infty$ . A larger  $p$  can lead to a more strict constraint on  $|u - u_0| \leq \delta$ . We conduct additional experiments on varying  $p$ . When using

Table 14: Classification and uncertainty estimation comparisons under **open-world** TTA settings with **ResNet-50-BN**, where  $P^{\text{test}} = 0.5P^{\text{Cov}} + 0.5P^{\text{Sem}}$  (Gaussian noise of severity level 3) and a suit of diverse abnormal outliers as same with Table 2. Substantial ( $\geq 0.5$ ) **improvement** and **degradation** compared to the baseline are highlighted in blue or red respectively.

		None		NINCO		iNaturalist		SSB-Hard		Texture		Places		Avg.	
Method	COME	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$
No Adapt	$\times$	3.0	81.6	3.0	91.2	3.0	89.4	3.0	90.2	3.0	88.3	2.8	90.7	3.0	88.6
PL	$\times$	26.9	71.2	16.1	84.3	12.9	88.4	15.9	87.8	18.1	86.1	16.9	82.9	17.8	83.5
TEA	$\times$	28.5	73.5	17.6	84.0	9.6	84.5	11.8	88.3	19.9	84.2	16.8	82.9	17.4	82.9
Tent	$\times$	52.5	67.5	43.7	79.6	52.1	78.3	51.9	82.1	44.4	78.6	48.7	73.2	48.9	76.6
	$\checkmark$	55.0	67.6	46.3	75.5	54.3	75.1	54.4	81.9	45.8	75.6	50.8	64.0	51.1	73.3
	Improve	$\Delta 2.6$	-	$\Delta 2.6$	$\nabla 4.2$	$\Delta 2.2$	$\nabla 3.2$	$\Delta 2.5$	-	$\Delta 1.4$	$\nabla 3.1$	$\Delta 2.2$	$\nabla 9.2$	$\Delta 2.2$	$\nabla 3.3$
EATA	$\times$	55.9	68.2	47.8	80.8	53.1	78.4	52.2	82.0	48.7	75.3	49.3	74.5	51.2	76.5
	$\checkmark$	58.0	66.2	52.9	74.8	57.6	73.2	57.4	81.3	52.5	70.7	55.4	62.9	55.6	71.5
	Improve	$\Delta 2.0$	$\nabla 2.0$	$\Delta 5.1$	$\nabla 6.0$	$\Delta 4.5$	$\nabla 5.1$	$\Delta 5.2$	$\nabla 0.7$	$\Delta 3.9$	$\nabla 4.6$	$\Delta 6.0$	$\nabla 11.6$	$\Delta 4.5$	$\nabla 5.0$
SAR	$\times$	51.8	64.6	42.4	78.3	47.6	81.3	48.1	84.4	42.7	79.1	46.0	76.7	46.4	77.4
	$\checkmark$	56.3	64.0	46.7	77.9	55.3	77.1	55.1	81.6	46.4	77.5	52.5	68.1	52.0	74.4
	Improve	$\Delta 4.5$	$\nabla 0.6$	$\Delta 4.3$	$\nabla 0.3$	$\Delta 7.7$	$\nabla 4.2$	$\Delta 6.9$	$\nabla 2.8$	$\Delta 3.7$	$\nabla 1.6$	$\Delta 6.5$	$\nabla 8.6$	$\Delta 5.6$	$\nabla 3.0$
COTTA	$\times$	22.6	70.7	14.4	87.4	21.1	78.6	19.7	84.1	15.5	87.3	15.8	82.0	18.2	81.7
	$\checkmark$	24.5	69.4	14.7	86.1	21.4	81.6	19.4	86.1	16.0	85.7	16.4	82.6	18.7	81.9
	Improve	$\Delta 1.8$	$\nabla 1.3$	-	$\nabla 1.3$	-	$\Delta 2.9$	-	$\Delta 2.0$	$\Delta 0.5$	$\nabla 1.6$	$\Delta 0.6$	$\Delta 0.7$	$\Delta 0.5$	-
MEMO	$\times$	8.0	83.6	7.5	89.0	7.9	87.9	7.9	89.8	7.8	88.6	7.7	88.4	7.8	87.9
	$\checkmark$	9.1	77.9	8.7	90.2	9.0	87.3	9.1	89.2	9.1	87.4	8.7	88.8	9.0	86.8
	Improve	$\Delta 1.1$	$\nabla 5.7$	$\Delta 1.2$	$\Delta 1.2$	$\Delta 1.2$	$\nabla 0.7$	$\Delta 1.2$	$\nabla 0.6$	$\Delta 1.3$	$\nabla 1.2$	$\Delta 1.1$	-	$\Delta 1.2$	$\nabla 1.1$

Table 15: Comparisons w.r.t. computation complexity. Accuracy (%) and FPR (%) are average results on ImageNet-C (level 5) with ViT-Base. The Wall-Clock Time (seconds) and Memory Usage (MB) are measured for processing 50,000 images of ImageNet-C on a single RTX 4090 GPU.

Method	COME	Acc $\uparrow$	FPR $\downarrow$	Memory	Run Time
No Adapt	$\times$	39.8	67.5	853	59
LAME	$\times$	38.7	69.7	853	62
T3A	$\times$	41.0	67.7	984	179
PL	$\times$	51.3	69.1	6393	128
FOA	$\times$	53.7	63.6	869	1687
TEA	$\times$	46.9	68.3	17266	2865
Tent	$\times$	52.8	70.1	6393	129
	$\checkmark$	61.2	66.5	6393	130
	Improve	$\Delta 8.4$	$\nabla 3.6$	-	-
EATA	$\times$	62.1	65.1	6394	135
	$\checkmark$	64.5	63.8	6394	134
	Improve	$\Delta 2.4$	$\nabla 1.3$	-	-
SAR	$\times$	54.2	66.7	6393	253
	$\checkmark$	64.2	63.8	6393	254
	Improve	$\Delta 10.1$	$\nabla 2.9$	-	-
COTTA	$\times$	46.1	67.9	19612	738
	$\checkmark$	49.1	67.5	19611	739
	Improve	$\Delta 3.0$	$\nabla 0.3$	-	-
MEMO	$\times$	42.3	72.1	5392	20576
	$\checkmark$	43.2	70.8	5392	20530
	Improve	$\Delta 0.9$	$\nabla 1.3$	-	-

Table 16: Comparison w.r.t. imbalanced label shifts performance. Results obtained on ViT and ImageNet-C (level 5) under **imbalanced label shifts** TTA setting, where the imbalance ratio is  $\infty$ . Substantial ( $\geq 0.5$ ) **improvement** and **degradation** compared to the baseline are highlighted in blue or red respectively.

Methods	COME	Noise			Blur				Weather				Digital				Avg.
		Gauss.	Shot	Impul.	Defoc	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elast.	Pixel	JPEG	Acc $\uparrow$
No Adapt	$\times$	35.1	32.2	35.9	31.4	25.3	39.4	31.6	24.5	30.1	54.7	64.5	49.0	34.2	53.2	56.5	39.8
PL	$\times$	49.7	48.6	50.9	49.8	41.5	53.0	41.9	26.6	49.0	64.3	73.6	65.6	45.2	63.9	63.0	52.4
T3A	$\times$	33.4	30.3	34.2	31.3	26.8	38.7	32.1	25.1	29.3	54.5	62.8	48.8	37.4	51.9	56.2	39.5
TEA	$\times$	44.9	40.3	46.3	39.8	35.2	46.0	12.1	14.3	46.9	60.3	72.7	60.2	48.6	62.7	58.8	45.9
LAME	$\times$	47.0	43.3	48.2	39.8	31.8	50.3	39.4	30.5	37.1	66.0	75.4	63.5	42.0	65.1	68.1	49.8
FOA	$\times$	41.5	39.2	43.6	42.5	33.7	45.5	41.0	44.9	44.5	60.1	67.7	58.8	45.7	57.3	62.7	48.6
Tent	$\times$	52.4	51.9	53.3	53.8	48.1	57.0	46.2	10.3	53.5	67.9	74.2	67.1	52.3	66.5	64.9	54.6
	$\checkmark$	55.0	55.0	56.2	57.1	54.6	61.6	49.3	62.9	64.0	72.3	78.1	69.3	62.7	71.3	69.0	62.6
	Improve	$\Delta 2.5$	$\Delta 3.2$	$\Delta 2.9$	$\Delta 3.4$	$\Delta 6.5$	$\Delta 4.6$	$\Delta 3.1$	$\Delta 52.6$	$\Delta 10.6$	$\Delta 4.4$	$\Delta 4.0$	$\Delta 2.2$	$\Delta 10.4$	$\Delta 4.9$	$\Delta 4.1$	$\Delta 7.9$
SAR	$\times$	51.8	51.7	52.7	51.9	48.2	55.6	47.8	20.3	52.9	66.8	73.2	66.0	52.2	64.1	62.8	54.5
	$\checkmark$	56.0	56.0	57.2	58.0	56.3	62.3	54.1	64.0	64.3	72.4	78.3	69.6	64.0	71.5	69.1	63.5
	Improve	$\Delta 4.2$	$\Delta 4.4$	$\Delta 4.5$	$\Delta 6.2$	$\Delta 8.1$	$\Delta 6.6$	$\Delta 6.3$	$\Delta 43.8$	$\Delta 11.4$	$\Delta 5.7$	$\Delta 5.1$	$\Delta 3.6$	$\Delta 11.8$	$\Delta 7.4$	$\Delta 6.2$	$\Delta 9.0$
EATA	$\times$	52.0	53.6	53.9	49.3	49.5	54.4	55.6	58.1	56.9	69.6	74.9	63.6	61.1	68.0	64.2	59.0
	$\checkmark$	54.9	56.4	54.7	56.5	56.3	62.1	59.0	67.0	65.4	73.4	78.4	68.0	68.0	73.0	70.4	64.2
	Improve	$\Delta 2.8$	$\Delta 2.8$	$\Delta 0.8$	$\Delta 7.3$	$\Delta 6.8$	$\Delta 7.7$	$\Delta 3.4$	$\Delta 8.9$	$\Delta 8.5$	$\Delta 3.9$	$\Delta 3.5$	$\Delta 4.4$	$\Delta 6.8$	$\Delta 5.0$	$\Delta 6.1$	$\Delta 5.3$
COTTA	$\times$	42.9	40.0	44.6	36.0	29.7	44.8	37.2	42.3	46.4	60.7	72.9	65.0	45.4	61.6	62.9	48.8
	$\checkmark$	51.6	49.0	52.9	41.7	37.0	51.6	43.8	46.7	53.2	65.9	74.4	65.6	52.8	66.7	65.9	54.6
	Improve	$\Delta 8.6$	$\Delta 9.0$	$\Delta 8.3$	$\Delta 5.7$	$\Delta 7.4$	$\Delta 6.8$	$\Delta 6.6$	$\Delta 4.4$	$\Delta 6.9$	$\Delta 5.2$	$\Delta 1.5$	$\Delta 0.5$	$\Delta 7.4$	$\Delta 5.1$	$\Delta 3.0$	$\Delta 5.8$
MEMO	$\times$	39.7	36.5	39.8	32.4	25.8	40.3	34.7	27.5	32.8	53.5	66.2	56.0	35.7	55.9	58.2	42.3
	$\checkmark$	40.6	37.5	40.6	33.4	26.7	41.2	35.4	28.7	33.7	54.7	67.1	55.9	36.6	57.2	59.2	43.2
	Improve	$\Delta 0.8$	$\Delta 1.0$	$\Delta 0.8$	$\Delta 1.0$	$\Delta 0.9$	$\Delta 1.0$	$\Delta 0.7$	$\Delta 1.2$	$\Delta 0.9$	$\Delta 1.2$	$\Delta 0.8$	-	$\Delta 0.9$	$\Delta 1.3$	$\Delta 1.1$	$\Delta 0.9$

Table 17: Comparison w.r.t. mixed shifts performance. Results obtained on ViT and ImageNet-C (level 5) under **mixed shifts** TTA setting, the performance is evaluated on a single data stream consisting of 15 mixed corruptions. Substantial ( $\geq 0.5$ ) **improvement** and **degradation** compared to the baseline are highlighted in blue or red respectively.

COME	TENT		SAR		EATA		CoTTA		Avg.	
	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$	Acc $\uparrow$	FPR $\downarrow$
$\times$	58.0	72.3	53.6	68.2	58.8	71.3	62.0	69.7	58.1	70.4
$\checkmark$	61.2	67.9	62.3	66.9	61.8	67.0	65.1	70.7	62.6	68.1
Improve	$\Delta 3.2$	$\nabla 4.4$	$\Delta 8.6$	$\nabla 1.3$	$\Delta 3.0$	$\nabla 4.3$	$\Delta 3.1$	$\Delta 0.9$	$\Delta 4.5$	$\nabla 2.3$

1458 infinity norm, a suboptimal classification accuracy is observed. We suppose this is because an overly  
1459 strict constraint can be harmful to TTA. Since on reliable test samples, we still expect to reduce the  
1460 uncertainty (in a conservative manner).  
1461

## 1462 D.2 LIMITATIONS AND FUTURE WORK 1463

1464 Many state-of-the-art TTA methods are equipped with entropy minimization learning principle;, but  
1465 the potential pitfalls lie in this optimization objective is not well understood. In this paper, we provide  
1466 empirical analysis towards understanding the failure mode. These findings motivate us to further  
1467 explore the connection between uncertainty learning and reliable TTA progress, which further implies  
1468 a principle to design novel learning principle as an alternative to entropy minimization. Finally,  
1469 we perform extensive experiments on multiple benchmarks to support our findings. In the work,  
1470 a simple yet effective regularization on the uncertainty mass is devised, and other regularization  
1471 techniques could be explored. Another interesting direction is further explore the relationship between  
1472 overconfidence issue and model collapse theoretically.  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511