

ETAGE: ENHANCED TEST TIME ADAPTATION WITH INTEGRATED ENTROPY AND GRADIENT NORMS FOR ROBUST MODEL PERFORMANCE

Afshar Shamsi[†], Rejisa Becirovic^{††}, Ahmadreza Argha^{††}, Ehsan Abbasnejad[‡]
Hamid Alinejad-Rokny^{††}, Arash Mohammadi[†]

[†] Concordia Institute of Information Systems Engineering, Concordia University, Montreal, Canada

[‡] Australian Institute for Machine Learning (AIML), University of Adelaide, Australia

^{††} School of Biomedical Engineering, UNSW Sydney, Sydney, Australia

ABSTRACT

Test time adaptation (TTA) equips deep learning models to handle unseen test data that deviates from the training distribution, even when source data is inaccessible. While traditional TTA methods often rely on entropy as a confidence metric, its effectiveness can be limited, particularly in biased scenarios. Extending existing approaches like the Pseudo Label Probability Difference (PLPD), we introduce ETAGE, a refined TTA method that integrates entropy minimization with gradient norms and PLPD, to enhance sample selection and adaptation. Our method prioritizes samples that are less likely to cause instability by combining high entropy with high gradient norms out of adaptation, thus avoiding the overfitting to noise often observed in previous methods. Extensive experiments on CIFAR-10-C and CIFAR-100-C datasets demonstrate that our approach outperforms existing TTA techniques, particularly in challenging and biased scenarios, leading to more robust and consistent model performance across diverse test scenarios. The codebase for ETAGE is available on <https://github.com/afsharshamsi/ETAGE>.

Index Terms— Test time adaptation, distribution shift, entropy minimization

1. INTRODUCTION

Deep learning models [1–3] have demonstrated significant success across various tasks, particularly when both training and testing data share the same distribution. In real-world applications, however, such models often face data that deviates from training distribution [4, 5], a challenge known as domain shift or dataset shift. This discrepancy between the training (source) data and the testing (target) data can severely undermine model performance [6], thereby limiting their practical efficiency. Tackling this issue requires innovative strategies that enable models to adapt to new data distributions without requiring additional training data or supervision.

Literature Review: Unsupervised Domain Adaptation (UDA) [7, 8] has been one approach to address the aforementioned issue. Generally speaking, in the UDA approach, the knowledge from labeled source data is transferred to unlabeled target data using both datasets during the training. By analyzing the distribution of the target set, UDA allows the model to learn domain-invariant features [9] that generalize well across different distributions. While UDA has shown promise in various scenarios, it still relies on the availability of source data during the adaptation process. Test Time Adaptation (TTA) [10] bridges this gap by enabling models to adapt solely at

test time, without access to the source data. This makes TTA particularly appealing in situations where retraining is impractical due to privacy concerns, computational constraints, or the urgent need for adaptation.

Test-Time Entropy Minimization (TENT) [11] is among the first strategies proposed for TTA, which focused on reducing the entropy of model predictions during test time. TENT enhances model robustness to corrupted datasets and domain adaptation scenarios by directly minimizing prediction entropy, i.e., updating model parameters through entropy minimization without altering the training process. Entropy-Aware Test-time Adaptation (EATA) [12] is another noteworthy TTA approach that extended TENT by stabilizing model weights to address the “forgetting” phenomenon, i.e., degradation of adapted models’ performance on in-distribution test samples. EATA starts with an initial pass over the target data to estimate the distribution, applying Fisher-based [13] weighting to identify crucial parameters for adaptation. This ensures that critical model weights remain stable and minimizes the impact of noisy or irrelevant samples. Sharpness-Aware Minimization (SAR) [14] further improves stability by optimizing only the most reliable features. By guiding reliable samples towards a flat minimum, it enhances the model’s stability and resilience, even under severe distribution shifts. SAR is based on the insight that models trained in flatter regions of the loss landscape tend to be more robust to perturbations [15], thereby, ensuring consistent performance even under challenging conditions. Despite the advancements made, existing approaches to TTA focus primarily on entropy as the key measure of confidence. This focus, however, has limitations, particularly in scenarios where spurious correlation shift makes entropy an unreliable confidence metric [16]. Recently, DeYO [17] proposed the concept of Pseudo Label Probability Difference (PLPD) to better identify harmful samples that may compromise model’s performance. DeYO enhances performance by intentionally degrading the shape of objects in images, ensuring the model’s judgments are based on generalizable features rather than misleading patterns.

Contributions: Building upon SAR’s focus on mitigating noisy samples through sharpness aware minimization and DeYO’s use of shape information in TTA, we propose the enhanced test time adaptation with integrated entropy and gradient norms (ETAGE) method. ETAGE introduces a refined test-time adaptation strategy that integrates entropy minimization with gradient norms and PLPD directly addresses the limitations of SAR and DeYO. By considering gradient norms, we capture the model’s sensitivity to noisy samples more effectively, enhancing the stability and efficacy of the adaptation process. Additionally, we provide a mathematical analysis

This work was partially supported by the Natural Sciences and Engineering Research Council (NSERC) of Canada through the NSERC Discovery Grant RGPIN-2023-05654

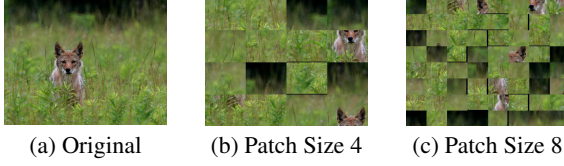


Fig. 1. Patch augmentation for estimating PLPD.

demonstrating why PLPD alone may miss noisy gradient samples, and propose a filtering method to eliminate these samples by combining high entropy with high gradient norms out of adaptation. The proposed ETAGE method not only avoids overfitting to noise but also ensures more robust performance across different test sets, as supported by empirical evidence/results from our experiments. In short, the paper makes the following key contributions:

- Introduction of a refined test-time adaptation approach, ETAGE, that couples entropy minimization with gradient norms and PLPD. ETAGE filters out noisy gradient samples by combining high entropy and high gradient norms, avoiding overfitting to noise.
- Theoretically illustrating shortcomings of PLPD in identifying noisy gradient samples, and mathematically demonstrating how to address this limitation.
- To the best of our knowledge, this is the first implementation of Contrastive Language-Image Pre-training (CLIP) foundation model [18] for TTA.

For performance evaluations, extensive experiments were conducted on CIFAR-10-C and CIFAR-100-C. It is observed that ETAGE achieves superior generalization and more consistent performance across different test sets compared to its state-of-the-art counterparts. The rest of the paper is organized as follows: First, Section 2 provides background information required for presentation of the proposed TTA approach. The ETAGE method is then introduced in Section 3. Section 4, first, presents the datasets utilized in this study, and then provides in details experimental results. Finally Section 5 concludes the paper.

2. PRELIMINARIES

In this section, we provide the required background used for development of the proposed ETAGE method.

2.1. Distribution Shift

Domain shift refers to differences between the distributions of training and testing data, which can be formally described as follows

$$\mathbb{S} = \{(\mathbf{x}, \mathbf{y}) | \forall (\mathbf{x}, \mathbf{y}) \sim p_{\mathbb{S}}\}, \quad f_{\mathbb{S}} : f_{\mathbb{S}}(\mathbf{x}) = \mathbf{y}$$

$$\mathbb{T} = \{(\mathbf{x}, \mathbf{y}) | \forall (\mathbf{x}, \mathbf{y}) \sim p_{\mathbb{T}}\}, \quad f_{\mathbb{T}} : f_{\mathbb{T}}(\mathbf{x}) = \mathbf{y}$$

where (\mathbf{x}, \mathbf{y}) represents the source data, which follows a probability distribution denoted by $p_{\mathbb{S}}$. Function $f_{\mathbb{S}}(\cdot)$ is a mapping that assigns \mathbf{x} to \mathbf{y} . The target data, \mathbb{T} , is defined in a similar fashion. Consider a Neural Network (NN) denoted by $f(\mathbf{x}, \boldsymbol{\theta}) : \mathbf{x} \rightarrow \mathbf{y}$, where $f(\cdot)$ represents a fixed architecture with parameters $\boldsymbol{\theta}$. Generally speaking, the objective is to minimize the loss function given by

$$\varepsilon_{\mathbb{T}}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\mathbb{T}}} [\ell(f(\mathbf{x}, \boldsymbol{\theta}), \mathbf{y})].$$

Such a minimization task becomes particularly challenging in a TTA setup, where access to the source data is restricted.

2.2. Entropy Minimization

To tackle the above mentioned issue, where the model is required to adapt to new, unseen data during inference without access to the source data, one effective strategy is entropy minimization. The key

idea behind this approach is to adjust the model’s parameters in real-time, focusing on minimizing the entropy of the output predictions. Entropy, in this context, measures the uncertainty of the model’s predictions. By minimizing entropy, the model becomes more confident in its predictions, effectively adapting to the new distribution presented by the test data. To achieve this, typically, the focus is on adapting only specific layers within the model, particularly normalization layers, which are believed to serve as proxies for the source data by retaining certain statistics of the source domain. This approach was first introduced with Batch Normalization (BN) layers in TENT, and has since been extended to other normalization techniques such as Group Normalization (GN) and Layer Normalization (LN) [14]. By adjusting these normalization layers during test time, the model can adapt to the new data distribution without requiring access to the source data, therefore, maintaining its performance across varying domains.

2.3. Pseudo Label Probability Difference (PLPD)

PLPD focuses on estimating the probability associated with an input for the same class before and after applying noise or augmentation, and is computed as follows

$$\text{PLPD}_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}') = (\mathbf{P}_{\boldsymbol{\theta}}(\mathbf{x}) - \mathbf{P}_{\boldsymbol{\theta}}(\mathbf{x}'))_{\hat{y}'}, \quad (1)$$

where \mathbf{x} and \mathbf{x}' are the sample input before and after, respectively, applying the augmentation/noise. The PLPD measures the sensitivity of the model’s predictions to slight changes in the input. For example, the patch shuffle technique, as shown in Figure 1, rearranges small regions (patches) in the input to alter spatial dependence within the image. By evaluating the PLPD under such conditions, one can assess the accuracy of the model when key information is disrupted. The core idea is that if the model still assigns the same class to the input after noise or augmentation, it may have learned spurious correlations from the source, rather than focusing on the shape of the object in the input image. The noise or augmentation is, typically, designed to disrupt or destroy objects within the image. The desired scenario is a high PLPD, which indicates that after the object is corrupted, the model is unable to assign the same class. This scenario is crucial as it identifies samples that should be used later for adaptation, ensuring the model’s predictions are based on more meaningful features.

3. THE ETAGE METHOD

In this section, we introduce the proposed ETAGE, that leverages both a gradient norm threshold and PLPD to enhance sample filtering. We premise that a high gradient norm indicates that the model is highly sensitive to small perturbations in the input space. Such samples, even if they produce an acceptable PLPD, might be misleading because the underlying instability (captured by the high gradient norm) is not reflected in the PLPD calculation. To explore this premise, a mathematical counterexample is provided below.

ETAGE, first applies a gradient norm threshold in addition to an entropy threshold to ensure that only stable samples proceed to the PLPD calculation. The gradient norm is used to identify samples that may cause instability in the model’s predictions, and is defined as

$$\|\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{x}, \boldsymbol{\theta})\| \propto \left\| \frac{\partial P(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}} \right\|. \quad (2)$$

A high gradient norm indicates that the model’s output is highly sensitive to small changes in the input, which can be a sign of overfitting to noise or instability in the learned decision boundary. Now, consider a sample x where the gradient norm is very high, i.e.,

$$\|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\theta})\|_2 \gg \tau_{\text{Grad}},$$

Algorithm 1 Online Adaptation for ETAGE

```

1: Input: Model, Test samples
2: Output: Adapted Model
3: for each batch do
4:   Compute norm, and entropies
5:   Applying entropy and norm thresholds
6:   if filtered samples exist then
7:     Patch permutation on remaining samples
8:     Calculate PLPD
9:     Filter samples based on PLPD threshold
10:    if filtered samples exist then
11:      Compute loss using remaining entropies
12:      Backpropagate loss to update model
13:      Update optimizer for the model
14:    end if
15:  end if
16: end for

```

where τ_{Grad} is the gradient norm threshold. A high gradient norm suggests that even a small change in \mathbf{x} will cause a significant change in the model's prediction. Suppose we apply a small perturbation δ to \mathbf{x} , resulting in $\mathbf{x}' = \mathbf{x} + \delta$. Using the first-order Taylor expansion, the predicted probability is approximated as

$$P(\mathbf{y} | \mathbf{x}', \theta) \approx P(\mathbf{y} | \mathbf{x}, \theta) + \nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta) \cdot \delta. \quad (3)$$

Equation (1) can then be approximated as

$$\text{PLPD}_{\theta}(\mathbf{x}, \mathbf{x}') = (P(\mathbf{y} | \mathbf{x}, \theta) - (P(\mathbf{y} | \mathbf{x}, \theta) + \nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta) \cdot \delta))_{\hat{\mathbf{y}}}, \quad (4)$$

which simplifies to

$$\text{PLPD}_{\theta}(\mathbf{x}, \mathbf{x}') \approx -(\nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta) \cdot \delta)_{\hat{\mathbf{y}}}. \quad (5)$$

The value of PLPD can be further estimated as

$$\text{PLPD}_{\theta}(\mathbf{x}, \mathbf{x}') \approx \|\nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta)\|_2 \|\delta\|_2 \cos(\theta), \quad (6)$$

where θ is the angle between $\nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta)$ and δ . The PLPD values primarily depend on two components: the gradient norm $\|\nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta)\|_2$ and the perturbation magnitude $\|\delta\|_2$. Now, consider a scenario where the input is highly sensitive to small perturbations. In such cases, as illustrated in Equation (2), these samples will generate high gradient norms, meaning the PLPD is largely influenced by $\|\nabla_{\mathbf{x}} P(\mathbf{y} | \mathbf{x}, \theta)\|_2$. Consequently, even with low perturbations, the PLPD values will be disproportionately high. This undermines the reliability of the PLPD since it is meant to reflect the difference in predicted probabilities before and after perturbation. In other words, for noisy gradient samples, the reason the PLPD is high (and why such samples might pass the threshold) is due to the high gradient norm. To tackle this issue, we refine the sample selection to exclude those with high gradient norms, i.e.,

$$\mathbf{S}'_{\theta}(\mathbf{x}) = \{\mathbf{x} | \text{Ent}_{\theta}(\mathbf{x}) > \tau_{\text{Ent}}, \|\nabla_{\theta} \mathcal{L}(\mathbf{x}, \theta)\| < \tau_{\text{Grad}}, \text{PLPD}_{\theta}(\mathbf{x}, \mathbf{x}') > \tau_{\text{PLPD}}\}$$

We remove these noisy, high-gradient samples from the set, leading to a more stable and effective adaptation process. The final entropy loss to be minimized is given by

$$\mathcal{L}_{\text{final}}(\theta) = \sum_{\mathbf{x} \in \mathbf{S}'_{\theta}(\mathbf{x})} \text{Ent}_{\theta}(\mathbf{x}). \quad (7)$$

Algorithm 1 outlines steps to perform ETAGE method. Figure 2 also presents different samples of CIFAR-10-C for Gaussian noise

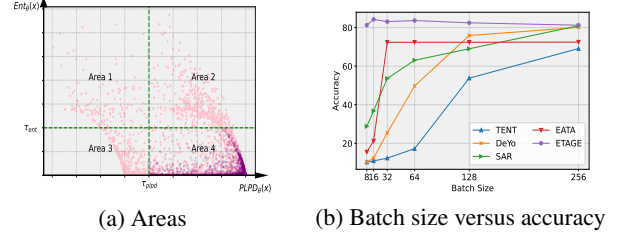


Fig. 2. The left figure classifies the test set data into four different areas using three thresholds: entropy, gradient norm, and PLPD. It should be noted that only the points/samples in purple, located in area 4, pass the thresholds and will be used for ETAGE. The right figure shows the performance of different methods under various batch sizes for Gaussian noise on CIFAR-10-C.

distortion based on the entropy and PLPD generated by the CLIP model. Following the common approach [11], a threshold, τ_{ent} , is defined for entropy selection to remove samples with low confidence (Areas 1 and 2 in Figure 2(a)). For the remaining samples, then the norm of the gradients are estimated to remove noisy samples, then the norm of the PLPD and keeping those samples with high PLPD (larger than the threshold) is the next step performed to remove samples that model tries to predict without considering the object in the image (i.e., model memorized such from the source domain).

As it is depicted in Area 4 of Figure 2(a), ETAGE further refines the sample selection process for adaptation by removing noisy gradient samples and only utilizing healthy samples (those shown in purple). This in turn leads to adapting the model with fewer samples. The proposed ETAGE method differs from SAR [14] in the sense that the latter fails to explicitly identify and remove noisy gradient samples. Different from our direct approach, SAR only indirectly mitigates the effects of noisy gradient samples using sharpness aware minimization. ETAGE also differs from DeYO [17] in the sense that the latter does not take into account the harmful effect of noisy samples with high gradient norm that can significantly degrade the adaptation performance.

4. SIMULATIONS AND RESULTS

To evaluate performance of the ETAGE method, we have used CLIP [18], which is a large pre-trained model that utilizes contrastive learning to connect image and text data. This approach was first introduced in the groundbreaking CLIP research [19]. CLIP is composed of two main components: a text encoder and an image encoder. The text encoder is based on a Transformer architecture, which effectively processes and understands natural language input. The image encoder, on the other hand, can be implemented using either a Vision Transformer (ViT) (the one that we have utilized is the ViT-B/32) or a Convolutional Neural Network (CNN). Together, these encoders are trained to align images with their corresponding text descriptions, enabling the model to perform various tasks involving both visual and textual information. For the CLIP model we modify LN layers while keeping the rest parameters freezed in the adaptation phase.

4.1. Datasets

We employed CIFAR-10-C and CIFAR-100-C datasets in our experiments to evaluate the ETAGE method. Below, we provide a brief yet comprehensive introduction to these datasets, highlighting their key characteristics and relevance to our research objectives.

4.1.1. CIFAR-10-C

This dataset shares its training samples with CIFAR-10 [20], but the test set, CIFAR-10-C [21], is created by applying various distortions to the original CIFAR-10 test images. These modifications include a range of corruptions and noises, such as blur, Gaussian

Table 1. Performance of different methods on various corruptions, each with severity of 5, on CIFAR-10-C dataset regarding accuracy (%). The best-performing results are in **bold**, the second-best in underline.

Method	CIFAR-10-C										
	Brightness	Contrast	Elastic Transform	Gaussian Blur	Gaussian Noise	JPEG Compression	Pixelate	Saturate	Shot Noise	Spatter	Speckle Noise
TENT	<u>96.61</u>	95.84	89.65	93.86	53.78	36.06	93.22	94.39	55.00	92.60	60.19
DEYO	94.25	94.70	87.47	91.45	<u>75.88</u>	79.41	91.76	93.24	81.39	91.16	81.17
SAR	96.65	95.84	91.50	<u>94.51</u>	69.01	<u>86.00</u>	94.26	<u>95.28</u>	<u>84.41</u>	94.29	79.06
EATA	96.58	91.78	91.14	91.38	72.38	84.76	86.66	95.35	79.00	94.78	<u>81.35</u>
ETAGE	96.18	<u>95.65</u>	<u>91.34</u>	94.52	82.47	86.45	<u>93.59</u>	95.17	86.08	<u>94.38</u>	86.68
											91.14

Table 2. Performance of different methods on various corruptions, each with severity of 5, on CIFAR-100-C dataset regarding accuracy (%). The best-performing results are in **bold**, the second-best in underline.

Method	CIFAR-100-C										
	Brightness	Contrast	Elastic Transform	Gaussian Blur	Gaussian Noise	JPEG Compression	Pixelate	Saturate	Shot Noise	Spatter	Speckle Noise
TENT	84.31	81.61	61.93	77.09	16.19	23.94	77.79	77.88	32.97	77.22	22.12
DEYO	84.89	83.03	73.64	79.94	<u>57.24</u>	64.65	79.65	80.01	65.63	79.07	66.86
SAR	85.76	<u>83.28</u>	74.90	80.91	41.28	66.32	80.14	80.72	51.46	79.69	64.75
EATA	1.70	1.33	1.56	1.49	1.37	1.47	1.51	1.53	1.40	1.64	1.26
ETAGE	<u>85.18</u>	83.43	<u>74.14</u>	80.36	61.03	<u>65.75</u>	<u>79.72</u>	<u>80.56</u>	<u>64.94</u>	<u>79.14</u>	<u>66.17</u>
											74.58

Table 3. comparison of calibration and discriminative metrics across different methods on CIFAR-10-C.

Method	ECE	MCE	Brier	AUROC
TENT	0.4389	0.4539	0.08972	0.7905
DEYO	0.1945	0.4144	<u>0.04251</u>	0.9618
SAR	0.2731	0.3699	0.05783	0.8861
EATA	<u>0.1888</u>	0.3981	0.04528	<u>0.9633</u>
ETAGE	0.1327	<u>0.3727</u>	0.03007	0.9793

noise, and shot noise. CIFAR-10-C presents these corruptions at five different levels of severity, resulting in a total of 50,000 test images for each type of noise. This dataset is particularly valuable for assessing the performance of image recognition models across diverse real-world scenarios, therefore, contributing to the enhancement of their reliability and robustness.

4.1.2. CIFAR-100-C

This dataset extends CIFAR-100 by applying a range of corruptions (similar to that of CIFAR-10-C) to its test set creating challenging scenarios for evaluating image recognition models. This dataset is commonly used in computer vision tasks to assess the performance of different models under distribution shifts.

4.2. Results

We estimate and compare ETAGE with some state-of-the-art methods in the literature namely, DeYo, EATA, SAR, and TENT. The results presented in Table 1 and Table 2 demonstrate that ETAGE consistently outperforms the other methods in the literature across a variety of corruption types. Specifically, on the CIFAR-10-C dataset (Table 1), ETAGE achieves the highest average accuracy, surpassing all other models. Similarly, on the CIFAR-100-C dataset (Table 2), ETAGE also leads with the highest average performance, indicating its robustness and effectiveness in handling diverse corruptions. These results validate the superiority of our approach, reinforcing ETAGE’s capability to adapt effectively under varying conditions compared to state-of-the-art methods. It is also worth noting that for TTA, lowering the batch size is challenging (lower batch sizes such as 1, 2, 4 are known as wild scenarios). This is because TTA relies on calculating statistics (mean and variance) over normaliza-

tion layers (BN, GN, LN). When the batch size is small, these statistics may be less representative of the data distribution, leading to instability in the model’s performance. The performance of different methods for different batch sizes are shown in Figure 2(b) for Gaussian noise corruption with severity 5 of CIFAR-10-C. ETAGE retains its accuracy even with low batch sizes, while its counterparts’ performance drop severely as batch size decreases. This underscores the effect of identifying and filtering out the noisy gradient samples which has been described in more details in Section 3.

Table 3 provides a comprehensive comparison of different methods under Gaussian noise on CIFAR-10-C, including TENT, DEYO, SAR, EATA, and ETAGE, evaluated across key performance metrics: Expected Calibration Error (ECE), Maximum Calibration Error (MCE), Brier score, and Area Under the Receiver Operating Characteristic curve (AUROC). Lower ECE and MCE values suggest that a model’s probability estimates are more dependable. Similarly, a lower Brier score indicates higher accuracy in probability predictions, while a higher AUROC value reflects better discrimination between classes. Notably, ETAGE outperforms others across most metrics, showcasing superior calibration and discriminative capability which in turn underscores the effectiveness of ETAGE in handling Gaussian noise while maintaining reliable predictions.

5. CONCLUSION

This study introduce ETAGE, an improved method for TTA, addressing the limitations of existing approaches that primarily rely on entropy as a confidence metric. By integrating gradient norms with the PLPD, our approach effectively filters out noisy samples, leading to more stable and reliable model adaptation. The application of this method to the CIFAR-10-C and CIFAR-100-C datasets demonstrated its effectiveness in handling various distribution shifts, with consistent performance improvements over baseline methods. Our findings emphasize the importance of considering gradient information alongside entropy in TTA, providing a pathway to more robust model adaptation. As part of future work, we plan to extend our method to wild setups involving lower batch sizes and evaluate its performance on the IMAGENET-C dataset. These efforts will further assess the scalability and generalization of our approach in more diverse and challenging environments.

6. REFERENCES

- [1] Afshar Shamsi, Hamzeh Asgharnezhad, Shirin Shamsi Jokandan, Abbas Khosravi, Parham M Kebria, Darius Nahavandi, Saeid Nahavandi, and Dipti Srinivasan, “An uncertainty-aware transfer learning-based framework for covid-19 diagnosis,” *IEEE transactions on neural networks and learning systems*, vol. 32, no. 4, pp. 1408–1417, 2021.
- [2] Maryam Habibpour, Hassan Gharoun, Mohammadreza Mehdipour, AmirReza Tajally, Hamzeh Asgharnezhad, Afshar Shamsi, Abbas Khosravi, and Saeid Nahavandi, “Uncertainty-aware credit card fraud detection using deep learning,” *Engineering Applications of Artificial Intelligence*, vol. 123, pp. 106248, 2023.
- [3] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [4] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al., “Wilds: A benchmark of in-the-wild distribution shifts,” in *International conference on machine learning*. PMLR, 2021, pp. 5637–5664.
- [5] Sinno Jialin Pan and Qiang Yang, “A survey on transfer learning,” *IEEE Transactions on knowledge and data engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.
- [6] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar, “Do cifar-10 classifiers generalize to cifar-10?,” *arXiv preprint arXiv:1806.00451*, 2018.
- [7] Yaroslav Ganin and Victor Lempitsky, “Unsupervised domain adaptation by backpropagation,” in *International conference on machine learning*. PMLR, 2015, pp. 1180–1189.
- [8] Changhwa Park, Jonghyun Lee, Jaeyoon Yoo, Minhoe Hur, and Sungroh Yoon, “Joint contrastive learning for unsupervised domain adaptation,” *arXiv preprint arXiv:2006.10297*, 2020.
- [9] Pedro O Pinheiro, “Unsupervised domain adaptation with similarity learning,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8004–8013.
- [10] Fatemeh Azimi, Sebastian Palacio, Federico Raue, Jörn Hees, Luca Bertinetto, and Andreas Dengel, “Self-supervised test-time adaptation on video data,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 3439–3448.
- [11] Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell, “Tent: Fully test-time adaptation by entropy minimization,” *arXiv preprint arXiv:2006.10726*, 2020.
- [12] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui Tan, “Efficient test-time model adaptation without forgetting,” in *International conference on machine learning*. PMLR, 2022, pp. 16888–16905.
- [13] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al., “Overcoming catastrophic forgetting in neural networks,” *Proceedings of the national academy of sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [14] Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Zhiquan Wen, Yaofo Chen, Peilin Zhao, and Mingkui Tan, “Towards stable test-time adaptation in dynamic wild world,” *arXiv preprint arXiv:2302.12400*, 2023.
- [15] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein, “Visualizing the loss landscape of neural nets,” *Advances in neural information processing systems*, vol. 31, 2018.
- [16] Sara Beery, Grant Van Horn, and Pietro Perona, “Recognition in terra incognita,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 456–473.
- [17] Jonghyun Lee, Dahuin Jung, Saehyung Lee, Junsung Park, Juhyeon Shin, Uiwon Hwang, and Sungroh Yoon, “Entropy is not enough for test-time adaptation: From the perspective of disentangled factors,” *arXiv preprint arXiv:2403.07366*, 2024.
- [18] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt, “Openclip,” 2021.
- [19] Alec Radford, Jong Wook Kim, Chris Hallacy, A. Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever, “Learning transferable visual models from natural language supervision,” in *International Conference on Machine Learning (ICML)*, 2021.
- [20] Alex Krizhevsky, Geoffrey Hinton, et al., “Learning multiple layers of features from tiny images,” 2009.
- [21] Bao Gia Doan, Afshar Shamsi, Xiao-Yu Guo, Arash Mohammadi, Hamid Alinejad-Rokny, Dino Sejdinovic, Damith C Ranasinghe, and Ehsan Abbasnejad, “Bayesian low-rank learning (bella): A practical approach to bayesian neural networks,” *arXiv preprint arXiv:2407.20891*, 2024.