

CONTRASTIVE RESIDUAL ENERGY TEST-TIME ADAPTATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Test-Time Adaptation (TTA) enhances model robustness by enabling adaptation to target distributions that differ from training distributions, improving real-world generalizability. However, most existing TTA approaches focus on adjusting the conditional distribution and therefore exhibit poor calibration, as they rely on uncertain predictions in the absence of labels. Energy-based TTA frameworks provide an alternative by modeling the marginal distribution of target data without depending on label predictions, but their reliance on costly sampling hinders scalability in real-world scenarios where decisions must be made without latency. In this work, we propose *Contrastive Residual Energy Test-time Adaptation (CRETTA)*, a practical solution for reliable adaptation. We first redefine the marginal distribution of target data using residual energy function and embed it into contrastive objective. This design prevents overfitting through adaptive gradient reweighting mechanism that leverages the relative residual energy while eliminating the sampling process. Extensive experiments demonstrate that CRETTA achieves scalable and well-calibrated adaptation under real-world computational constraints.

1 INTRODUCTION

Deep learning models can achieve high accuracy on training and testing data from the same distribution. However, when the distribution of the test data diverges from the original training dataset, the performance of the deep learning models deteriorates. This *distribution shift* refers to changes in the underlying data statistics, such as feature distributions or environmental conditions, between training and deployment. It is a major challenge in real-world scenarios, where test samples are often drawn from a distribution that deviates from the training data.

To address distribution shifts during testing, *test-time adaptation* (TTA) strategy aims to adapt trained model instantly, thereby maintaining robust performance on unexpected out-of-distribution samples. Since ground-truth labels are unavailable at test-time, existing approaches such as Pseudo-labeling (Lee et al., 2013) and SHOT (Liang et al., 2020) use the model’s own predictions as pseudo-labels. Likewise, TENT (Wang et al., 2020) operates similarly by using the model’s predicted probability distribution as a surrogate ground-truth distribution within an entropy-minimization objective.

Formally, the entropy minimization objective for a test sample x_t is expressed as $-\sum_{k=1}^C p_\theta(\hat{y}_k | x_t) \log p_\theta(\hat{y}_k | x_t)$, where C is the number of classes and $p_\theta(\hat{y}_k | x_t)$ is the predicted probability of class \hat{y}_k . While this approach demonstrates promising accuracy, it relies on uncertain predictions without ground-truth supervision. As a result, optimizing the entropy minimization objective often drives the predicted probabilities to collapse to extreme values of 0 or 1, leading to overconfident predictions (Press et al., 2024). This behavior increases calibration error as illustrated in Figure 10. In high-stakes real-world scenarios, such overconfidence can be detrimental, underscoring the need for alternative TTA strategies beyond simple entropy minimization. In such settings, *well-calibrated adaptation is critical to avoid overconfident errors and ensure safe, reliable model behavior in real-world scenarios*.

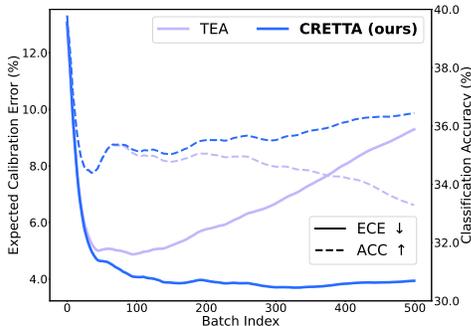
Adaptation with marginal distribution leads to better calibration. Rather than learning from the conditional distribution $p_\theta(\hat{y}|x)$ with unreliable model predictions, some approaches instead focus on modeling the marginal distribution $p_\theta(x)$ by directly maximizing its likelihood. This formulation avoids dependence on predicted labels and mitigates the overconfidence issues often associated with

054 entropy minimization. TEA (Yuan et al., 2024) applies maximum likelihood estimation (MLE) in
 055 the TTA setting and leverages contrastive divergence (CD) loss, an energy-based objective (Hinton,
 056 2002; LeCun et al., 2006; Song & Kingma, 2021), to align the model with unseen target distributions.
 057 Achieving better calibration while maintaining strong classification accuracy, TEA provides experi-
 058 mental evidence of stable and reliable adaptation, extending beyond the theoretical insights of earlier
 059 works (He et al., 2021; Wang et al., 2021; Schröder et al., 2024).

060 Despite its strengths, energy-based models with MLE methods for TTA still face two critical limita-
 061 tions: (i) **Unresolvable approximation error:** The normalization approximation relies on short-run
 062 Markov chains, which yield biased gradient estimates. As a result, parameter updates at test time
 063 can be unstable or converge to suboptimal solutions (Song & Kingma, 2021; Yair & Michaeli, 2021),
 064 particularly during few-shot adaptation or when dealing with high-dimensional data as shown in
 065 Figure 1. (ii) **High computational cost due to sampling-based approximation:** It requires repeated
 066 sampling to re-estimate the normalization constant for every incoming test batch, incurring substantial
 067 computational overhead during adaptation. This burden grows with model sizes, since each sampling
 068 step entails expensive gradient computations, thereby making real-time or resource-constrained TTA
 069 deployment infeasible, highlighting the need for methods that reduce overhead while maintaining
 070 strong performance.

071 To this end, we propose CRETТА, a novel residual
 072 energy-based test-time adaptation method that optimizes with marginal distribution while
 073 eliminating the normalization constant approximation, achieving high computational efficiency
 074 and scalability for real-world deployment.

075 First, we redesign the TTA task as learning the
 076 residual component of the distribution shift that the pretrained model has not yet captured and
 077 model the discrepancy with a residual energy
 078 function. Then, we embed the residual energy
 079 function into the contrastive learning objective,
 080 offsetting the normalization constant for the
 081 marginal distribution which typically requires
 082 extensive computations for approximation.



083 Figure 1: ImageNet-C (Sev 5) ECE(↓) and Acc(↑)
 084 over batch progress. CRETТА maintains stable
 085 calibration performance, while TEA experiences
 086 approximation error accumulation.

087 The primary contributions of our work can be summarized as follows:

- 088 • We introduce CRETТА, a novel sampling-free residual energy-based test-time adaptation (TTA)
 089 framework, offering computational efficiency and scalability for practical applications.
- 090 • Our method mathematically redefines the target marginal distribution with a residual energy
 091 function and optimizes it with contrastive learning objective. Consequently, this formulation
 092 enables TTA without costly normalization constant approximation and mitigates overfitting,
 093 thereby ensuring stable adaptation.
- 094 • CRETТА is well-calibrated and achieves strong performance across various distribution shifts.

095
 096
 097 **2 PRELIMINARIES**

098
 099 **Problem Setup** Let Q denote the marginal distribution of the source training data x_s . Consider a
 100 classifier $f_\phi(x)$ with parameters ϕ , which is pretrained on a labeled source dataset $\{(x_s^{(i)}, y_s^{(i)})\}_{i=1}^M$. Al-
 101 though this pretrained model performs well on in-distribution test data (i.e., $x_s \sim Q$), its performance
 102 can degrade substantially when tested on data from a different distribution $P(\neq Q)$, commonly
 103 referred to as out-of-distribution data.

104 Test-time adaptation (TTA) aims to mitigate this issue by adapting the pretrained parameters ϕ
 105 to better align with the target marginal distribution P . In this setting, we are given a set of N
 106 unlabeled target samples $\{x_t^{(i)}\}_{i=1}^N$ drawn from P , which arrive in online batches. To cope with the
 107 absence of label information, existing methods often rely on unsupervised objectives, particularly
 entropy minimization. EATA (Niu et al., 2022) and SAR (Niu et al., 2023) build on this foundation

by incorporating surrogate objectives and sample selection mechanisms to filter out unreliable predictions. However, these enhancements still fundamentally depend on uncertain model outputs, as dictated by the nature of entropy-based objectives.

Moving beyond entropy-based methods, recent work by TEA (Yuan et al., 2024) demonstrates the promise of energy-based modeling for TTA, where the central idea is to represent the test data marginal distribution through an energy function. Within this framework, TEA employs MLE on the marginal distribution of test samples $\{x_t^{(i)}\}_{i=1}^N$, so that the energy function is learned to assign lower energy (i.e., higher likelihood) to observed test inputs. By directly modeling and adapting to the marginal distribution, TEA mitigates distribution shifts without requiring labeled data or relying heavily on potentially unreliable model predictions. Building on this insight, our proposed method, CRETТА, approaches test-time adaptation not merely as an MLE procedure, but as effective learning of an energy function that represents an unnormalized marginal distribution.

Energy-based Models (EBMs) LeCun et al. (2006) express the marginal distribution in the energy-based model framework using Gibbs distribution, which can be formulated as $q_\phi(x) = \exp(-E_\phi(x))/Z(\phi)$, where $\phi \in \Phi$, with Φ representing the parameter space and $Z(\phi) = \int_{\mathbb{R}^d} \exp(-E_\phi(x))dx$ is the normalizing constant (partition function). The energy function $E_\phi(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ maps a d -dimensional data point to a scalar energy value, thereby defining an unnormalized density over the data space. The fundamental principle of EBMs is to represent the likelihood of a sample through this energy landscape: lower energy corresponds to higher likelihood and vice versa. A well-trained EBM thus learns to assign low energy values (i.e., high likelihood) to samples drawn from the in-distribution (source) Q , while assigning high energy (i.e., low likelihood) to out-of-distribution samples, such as those from a shifted target distribution P , where $P \neq Q$.

Grathwohl et al. (2019) and Yang & Ji (2021) present an innovative perspective on reframing discriminative models within the EBM framework. In their formulation, the energy function for a given input-label pair (x, y) is defined as $E_\phi(x, y) = -f_\phi(x)[y]$, where $f_\phi(x)[y]$ denotes the logit corresponding to label y in the discriminative model f_ϕ (i.e., classifier). Furthermore, the energy function derived from a discriminative model for a single input x can be expressed as the negative log-sum-exp of the logits across all classes in the final classifier layer:

$$E_\phi(x) = -T \cdot \log \sum_{k=1}^C e^{f_\phi(x)[k]/T}, \quad (1)$$

where T is a temperature parameter that controls the sharpness of the distribution (Liu et al., 2020). Finally, using a discriminative model within the EBM framework allows one to express the marginal probability of a data sample x . The main challenge in optimization stems from the normalization constant Z , which requires integrating over the entire input space, a task that is generally intractable. Consequently, EBMs often rely on specialized training methods such as contrastive divergence (Carreira-Perpinan & Hinton, 2005) or Markov chain Monte Carlo (MCMC) sampling to approximate or avoid the direct computation of Z .

3 METHODS

A Residual Perspective on Distribution Shift To characterize the distribution shift at test-time, we utilize a residual energy function that captures the discrepancy between the source and target distributions. Formally, let Q denote the source distribution and P the target distribution. We can express P in terms of Q via an exponential factor encoding the residual energy: $P = Q \exp(-R)/Z$, where R is a residual energy function, and Z is a normalization constant.

By analogy, the marginal distribution of the target data p_θ can be written as the product of the pretrained source model q_ϕ and an exponential residual term:

$$p_\theta(x) = \frac{1}{Z(\theta)} q_\phi(x) \exp\left(-\frac{1}{\beta} \tilde{E}_\theta(x)\right), \quad (2)$$

where the residual energy function \tilde{E}_θ is designed to model only the discrepancy between the fixed source model and the target distribution. Moreover, $\beta > 0$ is a temperature parameter and $\log Z(\theta)$ is constant across samples. During TTA, the source model q_ϕ remains frozen, and \tilde{E}_θ is learned to

capture only the distributional differences that arise under domain shift. In other words, \tilde{E}_θ focuses exclusively on the distribution-shift-induced residuals, refining the energy landscape so that the combined model aligns more closely with the true target distribution while preserving the original source knowledge.

Our residual interpretation of distribution shift can be viewed as an extension of the architectural constraints commonly employed in standard TTA setup. Building on the observation that updating only a subset of model parameters, such as the batch-normalization (BN) layers, enables efficient and stable adaptation by mitigating overfitting to severe distribution changes (Wang et al., 2020; Wu et al., 2024; Zhao et al., 2023), numerous TTA methods adopt such restricted update strategies. From a mathematical perspective, we extend this idea by recasting TTA as the problem of learning the residual component of the distribution shift that the pretrained model has not yet captured. This formulation serves as an implicit regularizer: it constrains the target model to learn only the unmodeled portion of the shift, thereby limiting deviation from the source distribution and preventing overfitting, as discussed in subsection 4.3.

Learning the Residual Energy via Contrastive Objective Energy-based models (EBMs) trained with maximum-likelihood estimation (MLE) often suffer from *biased gradients and prohibitive sampling costs*, primarily due to the need to approximate the intractable partition function Z (Song & Kingma, 2021). These limitations render conventional MLE approaches fundamentally ill-suited for practical TTA, where efficiency and stability are critical.

To overcome these challenges, we propose a *contrastive learning* framework that directly learns the residual energy function without any estimation or approximation of the partition function Z . Instead of optimizing likelihoods, we operate entirely on pairwise energy differences between source and target samples. This eliminates the need for sampling from the model distribution, making our method both tractable and scalable for TTA.

Our method shares conceptual similarities with Noise Contrastive Estimation (NCE) (Gutmann & Hyvärinen, 2010), in that both use contrastive objectives to bypass normalization. However, unlike NCE, which retains an implicit dependence on Z through the requirement of globally normalized densities, our formulation dispenses with Z entirely, as it only requires relative energies for learning.

Formally, we reinterpret the residual energy function $\tilde{E}_\theta(x)$ as arising from the density ratio between the target distribution $p_\theta(x)$ and the fixed source model $q_\phi(x)$:

$$\tilde{E}_\theta(x) = -\beta \left(\log \frac{p_\theta(x)}{q_\phi(x)} + \log Z(\theta) \right).$$

Assuming that the residual energy function \tilde{E}_θ should favor target samples x_t over source samples x_s (i.e., assigning lower energy to x_t than to x_s), we model the probability that the residual energy of a target sample is lower than that of a source sample, as:

$$P\left(\tilde{E}_\theta(x_t) < \tilde{E}_\theta(x_s)\right) = \frac{1}{1 + \exp\left(-\tilde{E}_\theta(x_s) + \tilde{E}_\theta(x_t)\right)} = \frac{1}{1 + \exp\left(\beta \log \frac{p_\theta(x_s)}{q_\phi(x_s)} - \beta \log \frac{p_\theta(x_t)}{q_\phi(x_t)}\right)},$$

During optimization, this objective drives the residual energy function \tilde{E}_θ to consistently reflect the distribution shift by lowering the relative energy of target samples with respect to source samples, thereby aligning the model with the target distribution while keeping the source model fixed.

Finally, we derive the optimization objective for learning the target model p_θ with the source model q_ϕ . Given a set of source and target pairs (x_s, x_t) , the objective can be formulated as minimizing the negative log-likelihood of the probability:

$$\mathcal{L}(\theta; \phi, \mathcal{B}) = -\frac{1}{|\mathcal{B}|} \sum_{(x_s, x_t) \sim \mathcal{B}} \log \sigma \left(\beta \underbrace{\left(\log \frac{p_\theta(x_t)}{q_\phi(x_t)} - \log \frac{p_\theta(x_s)}{q_\phi(x_s)} \right)}_l \right), \quad (3)$$

$$\text{where } l = -(E_\theta(x_t) - E_\theta(x_s)) + (E_\phi(x_t) - E_\phi(x_s)).$$

Crucially, as a consequence of our pairwise contrastive objective, the partition function Z cancels out completely, and no sampling is required unlike in MLE or NCE settings. Instead, we leverage a minimal buffer of source samples $\mathcal{B}_s = \{x_s^{(i)} \mid i = 1, \dots, |\mathcal{B}_s|\}$, to guide optimization and enable stable contrastive adaptation under test-time distribution shift. Despite its critical role in stabilizing optimization, the buffer size is negligibly small, imposing virtually no burden in modern memory settings and introducing no practical limitations in real-world deployment.

To optimize our objective (Equation 3), we construct a pairwise mini-batch $\mathcal{B} = \{(x_s^{(i)}, x_t^{(i)})\}$, where each pair consists of a target sample $x_t^{(i)} \in \mathcal{B}_t$ from the current target stream and a corresponding source sample $x_s^{(i)} \in \mathcal{B}_s$ randomly drawn from the source buffer \mathcal{B}_s . We demonstrate that our method maintains consistent performance even when the source buffer size $|\mathcal{B}_s|$ is reduced to as little as 1% of the source dataset, significantly lowering memory overhead. The robustness of our approach is further validated through an ablation study, as presented in Table 5.

Our proposed objective effectively aligns the model with the target distribution while avoiding the explicit estimation of the residual energy function. Furthermore, we reformulate both the source and target models as energy-based models, denoted as E_ϕ and E_θ , respectively, following Equation 1, with the target one initialized as $\theta = \phi$. This reformulation allows us to express the objective solely in terms of energy functions, eliminating the need for explicit normalization constants through algebraic simplifications. For a detailed derivation, we provide the full mathematical formulation in subsection A.1.1.

Why Does Contrastive Residual Learning Yield Stable Adaptation? The stable adaptation achieved by contrastive residual learning can be clarified through a gradient analysis. The gradient of our objective in Equation 3 is computed as follows:

$$\begin{aligned} \nabla_\theta \mathcal{L}(\theta; \phi, \mathcal{B}) &= -\frac{1}{|\mathcal{B}|} \sum_{(x_s, x_t) \sim \mathcal{B}} \beta \cdot w(x_t, x_s) \cdot (\nabla_\theta \log p_\theta(x_t) - \nabla_\theta \log p_\theta(x_s)), \\ \text{where } w(x_t, x_s) &= \sigma\left(\beta \log \frac{p_\theta(x_s)}{q_\phi(x_s)} - \beta \log \frac{p_\theta(x_t)}{q_\phi(x_t)}\right) \\ &= \sigma(\beta \cdot (E_\theta(x_t) - E_\phi(x_t)) - \beta \cdot (E_\theta(x_s) - E_\phi(x_s))). \end{aligned} \quad (4)$$

In this context, the term *contrastive* does not merely imply decreasing the energies of target samples or increasing those of source samples. Rather, the gradient weights are modulated by the relative energy levels of paired source-target samples, which promotes stable adaptation (see subsection 4.3 for further analysis).

If we remove the residual assumption, the pairwise contrastive learning objective reduces to the form in subsection C.1. In this case, there are no bias terms parameterized by ϕ , so the gradient magnitude depends solely on the target model’s energies, making the method more prone to overfitting, as illustrated in Figure 3. A more detailed mathematical discussion of the residual assumption and the pairwise contrastive approach is provided in Appendix B and Appendix C.

4 EXPERIMENT

In this section, we present a comprehensive analysis of our proposed method, CRETТА, and conduct a detailed comparison against state-of-the-art approaches using widely adopted benchmark datasets.

4.1 EXPERIMENTAL SETUP

Benchmark Datasets and Metrics To evaluate corruption robustness in test-time adaptation, we selected three benchmark datasets: (i) **CIFAR10-C**, (ii) **CIFAR100-C**, and (iii) **TinyImageNet-C** (Hendrycks & Dietterich, 2019). Each dataset contains 15 unique corruption types, categorized into 5 severity levels. In our evaluation, we reported the performance as the average across all 15 corruption types to provide a comprehensive measure of robustness. To rigorously evaluate the practical applicability of the proposed TTA method, we used three evaluation metrics: (i) **Accuracy (ACC)**, (ii) **Expected Calibration Error (ECE)**, and (iii) **Giga Floating-Point Operations (GFLOPs)**.

Table 1: Comparison of classification accuracy (Acc \uparrow) and expected calibration error (ECE \downarrow) on the CIFAR10-C, CIFAR100-C, and TinyImageNet-C datasets at corruption severity level 5 and the average across severity levels 1-5. The best results are emphasized in **BOLD**, while the second-best results are UNDERLINED.

Method	CIFAR-10-C				CIFAR-100-C				TinyImageNet-C			
	Severity L5		Severity Avg		Severity L5		Severity Avg		Severity L5		Severity Avg	
	Acc(\uparrow)	ECE(\downarrow)										
Source	81.73%	10.18%	88.82%	5.45%	53.25%	17.71%	64.11%	11.73%	35.12%	16.17%	43.16%	13.46%
<i>Normalization</i>												
BN Adapt	85.46%	4.85%	89.12%	3.15%	60.74%	8.32%	65.83%	6.88%	39.60%	<u>13.66%</u>	44.72%	<u>12.12%</u>
<i>Pseudo Labeling</i>												
PL	84.85%	10.10%	90.09%	6.20%	56.33%	23.81%	65.72%	16.66%	35.40%	30.95%	43.79%	23.47%
SHOT	87.91%	5.42%	90.78%	3.86%	<u>64.41%</u>	8.93%	<u>68.80%</u>	7.44%	39.84%	13.81%	44.95%	12.24%
<i>Entropy Minimization</i>												
TENT	87.84%	5.49%	90.74%	3.89%	64.31%	8.93%	68.73%	7.47%	39.83%	13.82%	44.94%	12.24%
ETA	85.46%	4.85%	89.12%	3.15%	61.77%	8.54%	66.66%	7.10%	39.67%	13.70%	44.79%	12.16%
EATA	85.46%	4.85%	89.12%	3.15%	61.79%	8.54%	66.65%	7.11%	39.68%	13.70%	44.79%	12.16%
SAR	86.54%	4.79%	89.80%	3.13%	62.71%	8.31%	67.36%	6.91%	39.66%	13.72%	44.77%	12.16%
AEA	<u>88.27%</u>	5.09%	<u>90.88%</u>	3.73%	64.40%	9.16%	68.75%	7.61%	39.87%	13.82%	44.97%	12.25%
<i>Energy-based Models</i>												
TEA	88.06%	3.83%	90.67%	2.68%	63.66%	7.68%	67.93%	6.33%	<u>39.96%</u>	13.84%	<u>45.08%</u>	12.24%
CRETTA	88.30%	<u>4.15%</u>	91.01%	<u>2.88%</u>	64.52%	<u>7.99%</u>	69.05%	<u>6.82%</u>	40.30%	13.52%	45.75%	11.85%

Baselines We compared our method with state-of-the-art approaches. (i) **Source:** The pre-trained classifier from the source data which performs inference on test data without adaptation. (ii) **Normalization-based:** BN-Adapt (Schneider et al., 2020) updates batch normalization statistics for test samples. (iii) **Pseudo-labeling-based:** Pseudo-Labeling (PL) (Lee et al., 2013) and SHOT (Liang et al., 2020) where test samples are filtered based on a confidence threshold, and the model is optimized using these pseudo-labels. (iv) **Entropy-based:** TENT (Wang et al., 2020), ETA, EATA (Niu et al., 2022), SAR (Niu et al., 2023), and AEA (Choi et al., 2025) aim to minimize entropy on test samples to achieve alignment with the target distribution. (v) **Energy-based:** TEA (Yuan et al., 2024) adapts to the marginal probability of the target distribution using energy-based learning with SGLD sampling.

Implementation Details In our experiments, we employed WRN-40-2 (Zagoruyko, 2016) for CIFAR10-C and CIFAR100-C datasets, and WRN-28-10 for TinyImageNet-C as backbones. Pre-trained weights were sourced from RobustBench (Croce et al., 2020). If unavailable, models were trained from scratch. We conduct online adaptation and evaluation following TENT (Wang et al., 2020) and TEA (Yuan et al., 2024) employing the Adam optimizer (Kingma, 2014) and reported results over three different random seeds. To further enhance robustness during adaptation, we incorporated data augmentation into the source buffer, which contained only 10% of the original source dataset. For more detailed information, please refer to the appendices.

4.2 PERFORMANCE COMPARISON

Classification Accuracy and Calibration Error Table 1 reports accuracy, focusing on the highest severity level 5 and the average across severity levels (1-5) across all datasets and corruption severities. Our proposed method consistently outperformed all other baselines under corrupted settings, notably achieving accuracy of 40.30% at the highest severity (level 5) on TinyImageNet-C. This consistent improvement highlights CRETTA’s adaptability and effectiveness in handling larger and more complex datasets, reinforcing its suitability for real-world test-time adaptation scenarios.

In TTA, model calibration is crucial for quantifying the prediction uncertainty, ensuring reliability under domain shifts and unlabeled data scenarios. We evaluated calibration performance using Expected Calibration Error (ECE) with 10 bins. While TEA performs well on CIFAR datasets, it fails to maintain the same level of superiority on TinyImageNet-C. In contrast, CRETTA demonstrates strong overall performance across all datasets (Table 1). Specifically, on TinyImageNet-C, CRETTA consistently outperforms other methods on most of corruption types in calibration as reported in Table 9.

Scalability We further evaluate CRETТА on PACS and ImageNet-C datasets. On PACS, CRETТА maintains competitive classification accuracy (Table 2) while achieving the lowest average ECE, outperforming the entropy-based method TENT and the existing energy-based method TEA by a significant margin (Table 3). On ImageNet-C, CRETТА achieved substantially lower ECE(2.69%) with higher accuracy, whereas TEA’s ECE was 7.21% (Table 12).

We interpret TEA’s weaker performance on both datasets as a consequence of approximation errors when estimating the normalization constant during sampling. These results underscore that CRETТА generalizes well to large-scale, style shifted datasets, achieving strong predictive performance and superior calibration.

Computational Efficiency A major challenge for previous energy-based TTA methods was the high computational cost for SGLD sampling. This makes them impractical for real-time TTA scenarios that demand rapid adaptation. This computational burden becomes even more pronounced as the input sample size increases. More precisely, TEA not only incurs approximately six times the computational cost (213K GFLOPs) compared to CRETТА (34K GFLOPs) but also struggles to maintain competitive performance. In contrast, CRETТА enables adaptation without explicitly tracking the normalization constant within a pair-wise contrastive learning framework. As shown in Figure 2, CRETТА consistently outperforms comparison methods, including TENT and BN-Adapt, while maintaining relatively low GFLOPs, demonstrating its efficiency for real-time TTA.

4.3 RESISTANCE MECHANISM TO OVERFITTING

In this subsection, we further analyze how the residual interpretation on distribution shift introduced in section 3 inherently provides a mechanism to mitigate overfitting and catastrophic forgetting.

Performance Under Gradual Distribution Shift To validate this mechanism, we conducted an experiment under a gradual distribution shift scenario. In this setting, the model continuously adapts from the source distribution Q through increasing shift intensities ($1 \rightarrow 5$), where severity 5 corresponds to the final target distribution P . After the model had fully adapted to P , we further froze the model and evaluated its classification accuracy on the original source distribution to observe if the target model remembers the original source distribution.

As summarized in Table 4, CRETТА demonstrates robust adaptation to progressively diverging target distributions,

Table 2: Classification accuracy (Acc \uparrow) on PACS.

Source Domain	Method	Target Domain				Avg
		Photo	Art	Cartoon	Sketch	
Photo	Source	-	55.39	22.61	23.17	33.72
	TENT	-	57.89	63.82	32.84	51.52
	TEA	-	53.16	50.33	33.17	45.55
	CRETТА	-	64.45	64.83	29.89	53.06
Art	Source	88.58	-	49.26	35.06	57.64
	TENT	92.10	-	63.18	35.23	63.50
	TEA	83.95	-	53.38	36.52	57.95
	CRETТА	91.52	-	65.90	40.99	66.13
Cartoon	Source	70.52	62.52	-	50.44	61.16
	TENT	82.34	64.63	-	41.47	62.81
	TEA	76.01	53.40	-	33.30	54.24
	CRETТА	83.65	68.86	-	40.04	64.18
Sketch	Source	13.89	14.50	19.03	-	15.81
	TENT	31.88	30.99	49.60	-	37.49
	TEA	18.60	25.47	48.34	-	30.80
	CRETТА	26.63	30.19	51.35	-	36.06

Table 3: ECE (\downarrow) on each source domain of PACS.

Method	P	A	C	S	AVG
TENT	44.41	35.20	34.69	58.63	43.23
TEA	41.71	34.62	35.02	50.26	40.40
CRETТА	37.42	28.22	26.65	51.68	35.99

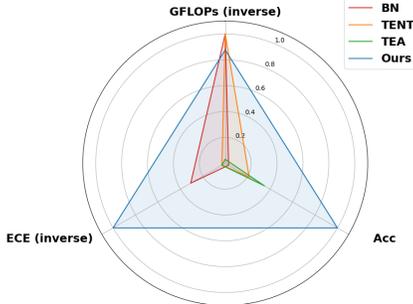


Figure 2: Comparison of GFLOPs, ECE and Acc against competitive baselines on TinyImageNet-C at the average across severity levels 1-5.

Table 4: Gradual distribution shift on CIFAR10(-C) and CIFAR100(-C).

Domain	CIFAR10		CIFAR100	
	OURS	TEA	OURS	TEA
Source (Q)	93.46	93.45	73.97	73.88
1	92.88	92.80	71.90	71.41
2	92.03	91.92	71.57	70.40
3	91.63	91.29	69.99	67.71
4	90.25	89.81	67.99	65.23
5 (P)	89.47	88.78	65.47	60.26
Source (Q)	94.03	93.58	75.70	69.25

378 outperforming TEA with classification accuracy on both datasets. Notably, the model’s accuracy on
 379 Q improved after adaptation to P (+1.73%). This improvement provides empirical evidence that
 380 CRETТА’s residual-energy formulation acts as a regularizer that prevents forgetting and facilitates
 381 robust adaptation. In contrast, MLE-based method TEA lacks this regularization mechanism and
 382 sometimes suffers from forgetting, with performance degradation (-4.63%) after adaptation.
 383

384 **Benefits of Contrastive Residual Energy Learning** CRETТА tends to reduce target-sample
 385 energy during adaptation across severities as shown in Figure 3, enhancing robustness to strong
 386 distribution shifts (Yuan et al., 2024). Notably, CRETТА achieves this with up to **8× reduction of**
 387 **the computational cost** compared to existing energy-based method TEA.

388 While the observed energy reduction
 389 is meaningful, the core strengths of
 390 CRETТА lie in *contrastive residual*
 391 *energy learning* to prevent conver-
 392 gence to trivial solutions as discussed
 393 in section 3. The bias terms with re-
 394 spect to E_ϕ in Equation 3 prevent E_θ
 395 from becoming excessively small or
 396 large, thereby stabilizing the adapta-
 397 tion process. As shown in Figure 3,
 398 the energy of target samples increases
 399 drastically after adaptation without
 400 bias terms, resulting in performance
 401 degradation.

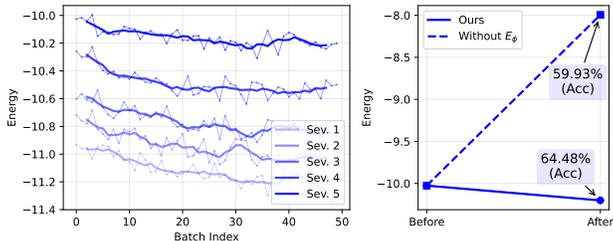


Figure 3: Energy trajectories of target samples across severities (left) and the effect of relative residual energy learning on stable adaptation (right) in CIFAR100-C.

402 Consequently, optimization proceeds
 403 adaptively based on the relative energy difference between source and target samples. When the
 404 energy difference between source and target is already aligned with the desired preference (i.e.,
 405 $E_\theta(x_t) < E_\theta(x_s)$), the gradient weight $w(x_t, x_s)$ in Equation 4 decreases, leading to weaker updates.
 406 On the other hand, when the energy difference exists in the opposite direction of the target preference
 407 (i.e., $E_\theta(x_t) > E_\theta(x_s)$), the gradient weight $w(x_t, x_s)$ increases to enforce stronger corrections.
 408 By letting gradient updates depend on these relative energy differences, CRETТА automatically
 409 modulates learning strengths through the weighting term, enabling dynamic and stable adaptation.

410 **4.4 ABLATION STUDIES**

411
 412 Other critical concerns regarding the replay buffer
 413 can be summarized in two folds: (1) Can the source
 414 data in the replay buffer be replaced with samples
 415 unseen during the pretraining phase? and (2) Does
 416 the model maintain its performance regardless of the
 417 quality of the samples included in the buffer?

418 The first question becomes particularly important
 419 from a data privacy perspective when the original
 420 source data is unavailable. To address this, we eval-
 421 uate the performance of CRETТА on CIFAR10-C
 422 with the replay buffer composed of CIFAR100 that
 423 we assume is distributionally similar but was not seen
 424 during the pretraining phase. As shown in Table 6,
 425 no performance drop is observed when the buffer is
 426 constructed from either CIFAR100 training or validation set. These findings suggest that CRETТА
 427 can operate effectively even without access to the original source data.

428 We further examined the performance of CRETТА in extreme cases where the source buffer com-
 429 position is biased. Specifically, the buffer was constructed using high confidence (top 10%) and
 430 low confidence (bottom 10%) samples respectively based on the source model’s confidence score.
 431 The results summarized in Table 7, indicate that adaptation performance remains unaffected by such
 variation in buffer content.

Table 5: Effectiveness of buffer size.

Buffer Ratio	CIFAR10-C	CIFAR100-C	TinyImageNet-C
1%	88.00%	64.21%	40.18%
2%	88.17%	64.42%	40.24%
10%	88.30%	64.52%	40.30%

Table 6: Effectiveness of source buffer content on CIFAR10-C.

Buffer Type	Sev 5	Sev 1-5 Avg.
Default(Random)	88.30%	91.01%
CIFAR100-trainset	88.20%	91.02%
CIFAR100-valset	88.21%	91.03%

432 These comprehensive findings highlight that
 433 a small and randomly sampled buffer suffices
 434 for effective adaptation, regardless of its size
 435 or composition. This insensitivity to buffer
 436 configuration underscores the practicality of
 437 CRETТА, enabling deployment in real-world
 438 scenarios with minimal memory overhead and
 439 flexible buffer sourcing.

Table 7: Effectiveness of source buffer confidence.

Buffer Type	CIFAR10-C	CIFAR100-C	TinyImageNet-C
Default (Random)	88.30%	64.52%	40.30%
High Confidence	86.92%	64.10%	40.18%
Low Confidence	88.02%	64.65%	40.92%

440 5 RELATED WORKS

441 **Test-time Adaptation** TTA is an emerging paradigm that has demonstrated immense potential in
 442 adapting pretrained models to unlabeled test data during testing phase. Early methods such as batch
 443 normalization adaptations (BN Adapt) (Schneider et al., 2020) leveraged test-batch statistics, while
 444 techniques like TTT (Sun et al., 2020) and TTT++ (Liu et al., 2021) utilized image augmentations.
 445 TENT (Wang et al., 2020), minimizes entropy to update BN layers, aiming for improved adaptation
 446 but often resulting in overconfident ‘that impair model calibration. EATA (Niu et al., 2022) and SAR
 447 (Niu et al., 2023) incorporates instance selection to filter unreliable samples, preserving performance,
 448 especially in continual test settings.

449 **Energy-based Models** Energy-based models (EBMs) are non-normalized probabilistic models
 450 widely used in classification and generative tasks (Grathwohl et al., 2019; Guo et al., 2023; Kim &
 451 Bengio, 2016). Energy provides a non-probabilistic scalar value capturing the density of the data
 452 distribution, making EBMs effective for capturing distribution shifts (Du & Mordatch, 2019). Due
 453 to this property, energy-based approaches are utilized in out-of-distribution (OOD) detection and
 454 unsupervised domain adaptation (Herath et al., 2023). Recent works such as AEA (Choi et al., 2025)
 455 and TEA (Yuan et al., 2024) extend energy to test-time adaptation scenario.

456 **Learning by Comparison** Noise-Contrastive Estimation (NCE) (Gutmann & Hyvärinen, 2010)per-
 457 forms maximum-likelihood estimation through nonlinear logistic regression, distinguishing real data
 458 from artificially generated noise. Although it provides insightful ideas as an optimization strategy,
 459 NCE still relies on the normalization constant implicitly, which can be difficult to handle in practice.
 460 By contrast, pairwise comparison removes the need for this constant by using linear logistic regression.
 461 Methods such as RLHF (Ouyang et al., 2022) and DPO (Rafailov et al., 2024) adopt this idea within
 462 autoregressive text-generation models to better align generated responses with human preferences.

463 6 CONCLUSION

464 In test-time adaptation, the entropy minimization objective often suffers from poor calibration due to
 465 the overconfidence problem, while existing energy-based methods incur significant computational
 466 overhead from extensive sampling to approximate the normalization constant for the marginal target
 467 distribution. In contrast, CRETТА achieves reliable and efficient adaptation by redefining the
 468 distribution shift with a residual energy function while optimizing a contrastive objective that avoids
 469 the sampling.

470 CRETТА provides two key benefits. First, it inherently mitigates overfitting by adaptively reweight-
 471 ing gradient signals based on relative energy differences, thereby ensuring stable adaptation. Second,
 472 by embedding the residual energy function into the contrastive objective, CRETТА eliminates the
 473 need for normalization constant approximation. Through comprehensive experiments and abla-
 474 tions CRETТА confirms that it bridges the gap between calibration-aware adaptation and practical
 475 feasibility, offering a scalable solution previously unattainable with conventional TTA frameworks.

REFERENCES

- 486
487
488 Ralph Allan Bradley and Milton E Terry. Rank analysis of incomplete block designs: I. the method
489 of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952.
- 490 Miguel A Carreira-Perpinan and Geoffrey Hinton. On contrastive divergence learning. In *International
491 workshop on artificial intelligence and statistics*, pp. 33–40. PMLR, 2005.
- 492
493 Wonjeong Choi, Do-Yeon Kim, Jungwuk Park, Jungmoon Lee, Younghyun Park, Dong-Jun Han, and
494 Jaekyun Moon. Adaptive energy alignment for accelerating test-time adaptation. In *The Thirteenth
495 International Conference on Learning Representations*, 2025.
- 496 Francesco Croce, Maksym Andriushchenko, Vikash Sehwal, Edoardo DeBenedetti, Nicolas Flam-
497 marion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial
498 robustness benchmark. *arXiv preprint arXiv:2010.09670*, 2020.
- 499 Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances
500 in Neural Information Processing Systems*, 32, 2019.
- 501
502 Taesik Gong, Jongheon Jeong, Taewon Kim, Yewon Kim, Jinwoo Shin, and Sung-Ju Lee. Note:
503 Robust continual test-time adaptation against temporal correlation. *Advances in Neural Information
504 Processing Systems*, 35:27253–27266, 2022.
- 505 Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi,
506 and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like
507 one. *arXiv preprint arXiv:1912.03263*, 2019.
- 508
509 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
510 networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- 511 Qiushan Guo, Chuofan Ma, Yi Jiang, Zehuan Yuan, Yizhou Yu, and Ping Luo. Egc: Image
512 generation and classification via a diffusion energy-based model. In *Proceedings of the IEEE/CVF
513 International Conference on Computer Vision*, pp. 22952–22962, 2023.
- 514 Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle
515 for unnormalized statistical models. In *Proceedings of the thirteenth international conference on
516 artificial intelligence and statistics*, pp. 297–304. JMLR Workshop and Conference Proceedings,
517 2010.
- 518
519 Tianxing He, Bryan McCann, Caiming Xiong, and Ehsan Hosseini-Asl. Joint energy-based
520 model training for better calibrated natural language understanding models. *arXiv preprint
521 arXiv:2101.06829*, 2021.
- 522 Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common
523 corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.
- 524
525 Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
526 examples in neural networks. *arXiv preprint arXiv:1610.02136*, 2016.
- 527 Samitha Herath, Basura Fernando, Ehsan Abbasnejad, Munawar Hayat, Shahram Khadivi, Mehrtash
528 Harandi, Hamid Reza Tofighi, and Gholamreza Haffari. Energy-based self-training and normal-
529 ization for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International
530 Conference on Computer Vision*, pp. 11653–11662, 2023.
- 531 Geoffrey E Hinton. Training products of experts by minimizing contrastive divergence. *Neural
532 computation*, 14(8):1771–1800, 2002.
- 533
534 Taesup Kim and Yoshua Bengio. Deep directed generative models with energy-based probability
535 estimation. *arXiv preprint arXiv:1606.03439*, 2016.
- 536
537 Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
538 2014.
- 539 Yann LeCun, Sumit Chopra, Raia Hadsell, M Ranzato, Fugie Huang, et al. A tutorial on energy-based
learning. *Predicting structured data*, 1(0), 2006.

- 540 Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for
541 deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, pp.
542 896. Atlanta, 2013.
- 543 Jian Liang, Dapeng Hu, and Jiashi Feng. Do we really need to access the source data? source
544 hypothesis transfer for unsupervised domain adaptation. In *International conference on machine*
545 *learning*, pp. 6028–6039. PMLR, 2020.
- 546 Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection.
547 *Advances in neural information processing systems*, 33:21464–21475, 2020.
- 548 Yuejiang Liu, Parth Kothari, Bastien Van Delft, Baptiste Bellot-Gurlet, Taylor Mordan, and Alexandre
549 Alahi. Ttt++: When does self-supervised test-time training fail or thrive? *Advances in Neural*
550 *Information Processing Systems*, 34:21808–21820, 2021.
- 551 Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Yafo Chen, Shijian Zheng, Peilin Zhao, and Mingkui
552 Tan. Efficient test-time model adaptation without forgetting. In *International conference on*
553 *machine learning*, pp. 16888–16905. PMLR, 2022.
- 554 Shuaicheng Niu, Jiaxiang Wu, Yifan Zhang, Zhiquan Wen, Yafo Chen, Peilin Zhao, and Mingkui
555 Tan. Towards stable test-time adaptation in dynamic wild world. *arXiv preprint arXiv:2302.12400*,
556 2023.
- 557 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
558 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
559 instructions with human feedback. *Advances in neural information processing systems*, 35:27730–
560 27744, 2022.
- 561 Ori Press, Ravid Shwartz-Ziv, Yann LeCun, and Matthias Bethge. The entropy enigma: Success and
562 failure of entropy minimization, 2024. URL <https://arxiv.org/abs/2405.05012>.
- 563 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
564 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
565 *in Neural Information Processing Systems*, 36, 2024.
- 566 Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias
567 Bethge. Improving robustness against common corruptions by covariate shift adaptation. *Advances*
568 *in neural information processing systems*, 33:11539–11551, 2020.
- 569 Tobias Schröder, Zijing Ou, Jen Lim, Yingzhen Li, Sebastian Vollmer, and Andrew Duncan. Energy
570 discrepancies: a score-independent loss for energy-based models. *Advances in Neural Information*
571 *Processing Systems*, 36, 2024.
- 572 Yang Song and Diederik P Kingma. How to train your energy-based models. *arXiv preprint*
573 *arXiv:2101.03288*, 2021.
- 574 Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei Efros, and Moritz Hardt. Test-time training
575 with self-supervision for generalization under distribution shifts. In *International conference on*
576 *machine learning*, pp. 9229–9248. PMLR, 2020.
- 577 Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully
578 test-time adaptation by entropy minimization. *arXiv preprint arXiv:2006.10726*, 2020.
- 579 Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation.
580 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
581 7201–7211, 2022.
- 582 Yezhen Wang, Bo Li, Tong Che, Kaiyang Zhou, Ziwei Liu, and Dongsheng Li. Energy-based
583 open-world uncertainty modeling for confidence calibration. In *Proceedings of the IEEE/CVF*
584 *International Conference on Computer Vision*, pp. 9302–9311, 2021.
- 585 Yanan Wu, Zhixiang Chi, Yang Wang, Konstantinos N. Plataniotis, and Songhe Feng. Test-time
586 domain adaptation by learning domain-aware batch normalization, 2024. URL <https://arxiv.org/abs/2312.10165>.

594 Omer Yair and Tomer Michaeli. Contrastive divergence learning is a time reversal adversarial game,
595 2021. URL <https://arxiv.org/abs/2012.03295>.
596

597 Xiulong Yang and Shihao Ji. Jem++: Improved techniques for training jem. In *Proceedings of the*
598 *IEEE/CVF International Conference on Computer Vision*, pp. 6494–6503, 2021.

599 Longhui Yuan, Binhui Xie, and Shuang Li. Robust test-time adaptation in dynamic scenarios.
600 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
601 15922–15932, 2023.

602

603 Yige Yuan, Bingbing Xu, Liang Hou, Fei Sun, Huawei Shen, and Xueqi Cheng. Tea: Test-time
604 energy adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
605 *Recognition*, pp. 23901–23911, 2024.

606 Sergey Zagoruyko. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
607

608 Hao Zhao, Yuejiang Liu, Alexandre Alahi, and Tao Lin. On pitfalls of test-time adaptation, 2023.
609 URL <https://arxiv.org/abs/2306.03536>.
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A TECHNICAL APPENDICES

A.1 DERIVATION AND FUNCTION OF CRETТА

A.1.1 DERIVATION OF CRETТА

The marginal distribution of target data p_θ can be written as the product of the pretrained source model q_ϕ and an exponential residual term:

$$p_\theta(x) = \frac{1}{Z} q_\phi(x) \exp\left(-\frac{1}{\beta} \tilde{E}_\theta(x)\right),$$

where \tilde{E}_θ represents the residual energy function encoding the distribution shift. During TTA, the source model q_ϕ remains fixed, and the objective is to learn \tilde{E}_θ so as to align the source model more closely with the target distribution. By expanding the equation with respect to the energy function \tilde{E}_θ , we can compute the residual energy score of an image sample x .

$$\tilde{E}_\theta(x) = -\beta \left(\log \frac{p_\theta(x)}{q_\phi(x)} + \log Z \right).$$

Next, we substitute the ground-truth residual energy function \tilde{E}_θ^* into the Bradley-Terry (BT) model (Bradley & Terry, 1952), which only depends on the difference in energy values between source and target pairs:

$$P(\tilde{E}(x_t) < \tilde{E}(x_s)) = \frac{1}{1 + \exp(-\tilde{E}_\theta(x_s) + \tilde{E}_\theta(x_t))} = \frac{1}{1 + \exp\left(\beta \log \frac{p_\theta(x_s)}{q_\phi(x_s)} - \beta \log \frac{p_\theta(x_t)}{q_\phi(x_t)}\right)},$$

where x_t and x_s denote the target and source samples, respectively. Here, for pairwise comparison, we use the negative residual energy.

Having derived the probability of the target distribution data in terms of the optimal energy function, which can further be expressed using ϕ and θ , our objective for the target model is as follows:

$$\mathcal{L}(\theta; \phi) = -\mathbb{E}_{(x_s, x_t) \sim B} \left[\log \sigma \left(\beta \log \frac{p_\theta(x_t)}{q_\phi(x_t)} - \beta \log \frac{p_\theta(x_s)}{q_\phi(x_s)} \right) \right] \quad (5)$$

In section 3, we emphasize that the key advantage of CRETТА is that it avoids the costly Stochastic Gradient Langevin Dynamics (SGLD) sampling required to compute the normalization constant as required in TEA (Yuan et al., 2024). However, the objective Equation 5 still includes the normalization constant for both target and source model.

To eliminate both normalization constants, we first redefine the target and source models using the Gibbs distribution as follows:

$$p_\theta(x) = \frac{\exp(-E_\theta(x))}{Z(\theta)}, \quad q_\phi(x) = \frac{\exp(-E_\phi(x))}{Z(\phi)}$$

By integrating p_θ and q_ϕ into the Equation 5 and applying the logarithm, the normalization constants for both target and source model, i.e., $Z(\theta)$ and $Z(\phi)$, are canceled out as shown in below:

$$\mathcal{L}(\theta; \phi) = -\mathbb{E}_{(x_s, x_t) \sim B} \left[\ln \sigma \left(\beta \left(-E_\theta(x_t) - \ln Z(\theta) + E_\phi(x_t) + \ln Z(\phi) \right) - \beta \left(-E_\theta(x_s) - \ln Z(\theta) + E_\phi(x_s) + \ln Z(\phi) \right) \right) \right] \quad (6)$$

Therefore, the final learning objective is expressed as follows:

$$\mathcal{L}(\theta; \phi) = -\mathbb{E}_{(x_s, x_t) \sim B} [\ln \sigma (\beta (-E_\theta(x_t) + E_\theta(x_s) + E_\phi(x_t) - E_\phi(x_s)))]$$

A.1.2 FUNCTION OF CRETТА

In this section, we provide a detailed explanation of how each component of CRETТА contributes to adaptation, as well as the expected behavior during early and late stages of online adaptation.

If the target model θ successfully optimizes this objective, then residual energy function $\tilde{E}_\theta(x)$ in $p_\theta(x) = \frac{1}{Z} q_\phi(x) \exp(-\frac{1}{\beta} \tilde{E}_\theta(x))$ models the residual component of the distribution shift between the source and target domains.

At the beginning of adaptation, the model has not yet encoded the distribution shift. Therefore, the residual energy $E_\theta(x)$ is close to zero for both source samples x_s and target samples x_t . This results in: $p_\theta(x) \approx q_\phi(x)$ meaning that predictions for both source and target data remain similar to the source model outputs.

As training progresses, the residual energy function learns the discrepancy between target and source distributions. For source samples x_s , $\tilde{E}_\theta(x_s)$ remains small, leading to $p_\theta(x_s) \approx q(x_s)$, preserving source performance. For target samples x_t , the residual energy adjusts predictions reflecting the learned domain shift and improving performance on the target domain. By progressively learning the residual while maintaining alignment with the source model, CRETТА achieves better generalization.

A.1.3 BUFFER MANAGEMENT OF CRETТА

CRETТА initializes the source buffer \mathcal{B}_s at model initialized, prior to adaptation, by randomly sampling source data up to a fixed buffer size, with an equal number of samples per class. During adaptation, the samples in the buffer are used sequentially in batches without any additional sampling or refresh, unlike TEA, thereby incurring no additional computational overhead.

Table 8: Comparison of classification accuracy (Acc \uparrow) and expected calibration error (ECE \downarrow) on the CIFAR10-C, CIFAR100-C, and TinyImageNet-C datasets at corruption severity level 5, the average across severity levels 1-5, and on clean data. The best adaptation results are emphasized in **BOLD**, while the second-best results are UNDERLINED.

Method	CIFAR-10-C					CIFAR-100-C					TinyImageNet-C				
	Clean Acc(\uparrow)	Corr Severity 5 Acc(\uparrow)	ECE(\downarrow)	Corr Severity 1-5 Avg Acc(\uparrow)	ECE(\downarrow)	Clean Acc(\uparrow)	Corr Severity 5 Acc(\uparrow)	ECE(\downarrow)	Corr Severity 1-5 Avg Acc(\uparrow)	ECE(\downarrow)	Clean Acc(\uparrow)	Corr Severity 5 Acc(\uparrow)	ECE(\downarrow)	Corr Severity 1-5 Avg Acc(\uparrow)	ECE(\downarrow)
Source	95.08%	81.73%	10.18%	88.82%	5.45%	76.28%	53.25%	17.71%	64.11%	11.73%	59.60%	35.12%	16.17%	43.16%	13.46%
<i>Normalization</i>															
BN Adapt	93.59%	85.46%	4.85%	89.12%	3.15%	72.84%	60.74%	8.32%	65.83%	6.88%	56.72%	39.60%	<u>13.66%</u>	44.72%	<u>12.12%</u>
<i>Pseudo Labeling</i>															
PL	94.85%	84.85%	10.10%	90.09%	6.20%	75.98%	56.33%	23.81%	65.72%	16.66%	58.95%	35.40%	30.95%	43.79%	23.47%
SHOT	94.38%	87.91%	5.42%	90.78%	3.86%	75.00%	<u>64.41%</u>	8.93%	<u>68.80%</u>	7.44%	56.90%	39.84%	13.81%	44.95%	12.24%
<i>Entropy Minimization</i>															
TENT	94.35%	87.84%	5.49%	90.74%	3.89%	74.95%	64.31%	8.93%	68.73%	7.47%	56.92%	39.83%	13.82%	44.94%	12.24%
ETA	93.72%	85.46%	4.85%	89.12%	3.15%	73.71%	61.77%	8.54%	66.66%	7.10%	56.82%	39.67%	13.70%	44.79%	12.16%
EATA	93.72%	85.46%	4.85%	89.12%	3.15%	73.66%	61.79%	8.54%	66.65%	7.11%	56.86%	39.68%	13.70%	44.79%	12.16%
SAR	93.61%	86.54%	4.79%	89.80%	3.13%	73.73%	62.71%	8.31%	67.36%	6.91%	56.77%	39.66%	13.72%	44.77%	12.16%
AEA	94.21%	<u>88.27%</u>	5.09%	<u>90.88%</u>	3.73%	75.17%	64.40%	9.16%	68.75%	7.61%	56.97%	39.87%	13.82%	44.97%	12.25%
<i>Energy-based Models</i>															
TEA	94.06%	88.06%	3.83%	90.67%	2.68%	74.18%	63.66%	7.68%	67.93%	6.33%	57.17%	<u>39.96%</u>	13.84%	<u>45.08%</u>	12.24%
CRETТА (Ours)	94.43%	88.30%	<u>4.15%</u>	91.01%	<u>2.88%</u>	75.26%	64.52%	7.99%	69.05%	<u>6.82%</u>	<u>58.23%</u>	40.30%	13.52%	45.75%	11.85%

A.2 ADDITIONAL EXPERIMENTS AND ANALYSIS

A.2.1 DETAILED PERFORMANCE COMPARISON

Detailed Performance Comparison on Accuracy Table 8 reports accuracy on the highest severity level 5, the average across severity levels (1-5), and performance on the clean dataset (i.e., without corruption) for CIFAR10-C, CIFAR100-C, and TinyImageNet-C. This table extends the results of Table 1 by additionally reporting accuracy on the clean dataset, providing a more complete view of model performance.

While CRETТА achieves the second-best accuracy on clean data among all methods, with the PL method performing the best. However, this can lead to overfitting and significant degradation in

Table 9: Comparison of expected calibration error (ECE \downarrow) on TinyImageNet-C datasets across all corruptions at the average across severity level 1-5. (Values are reported to one decimal place for space efficiency.)

Method	Noise			Blur				Weather				Digital				Avg
	Gauss.	Shot	Impul.	Defoc.	Glass	Motion	Zoom	Snow	Frost	Fog	Brit.	Contr.	Elastic	Pixel	JPEG	
Source	12.5%	11.6%	13.1%	<u>10.8%</u>	13.3%	<u>10.8%</u>	10.5%	14.7%	14.6%	18.4%	13.9%	25.9%	11.0%	10.1%	<u>10.6%</u>	13.5%
BN Adapt	<u>12.1%</u>	11.9%	12.4%	11.0%	<u>11.9%</u>	11.0%	<u>10.3%</u>	<u>13.1%</u>	<u>12.7%</u>	<u>14.0%</u>	<u>12.1%</u>	<u>16.7%</u>	<u>10.9%</u>	10.7%	10.9%	<u>12.1%</u>
PL	20.9%	19.5%	26.7%	18.4%	20.4%	18.1%	17.7%	22.8%	20.6%	38.4%	19.7%	54.7%	18.3%	17.6%	18.2%	23.5%
SHOT	12.2%	12.1%	12.5%	11.1%	12.1%	11.2%	10.5%	13.2%	12.8%	14.2%	12.2%	16.9%	11.0%	10.7%	11.0%	12.2%
TENT	12.2%	12.1%	12.5%	11.1%	12.1%	11.1%	10.5%	13.2%	12.8%	14.2%	12.2%	16.9%	11.0%	10.8%	11.0%	12.2%
ETA	12.1%	12.0%	12.5%	11.1%	12.0%	11.1%	10.4%	13.1%	12.7%	14.1%	12.2%	16.7%	10.9%	10.7%	10.9%	12.2%
EATA	12.1%	12.0%	12.4%	11.1%	12.0%	11.1%	10.4%	13.2%	12.7%	14.1%	12.2%	16.8%	10.9%	10.7%	10.9%	12.2%
SAR	12.1%	12.0%	12.5%	11.1%	12.0%	11.1%	10.4%	13.2%	12.7%	14.1%	12.2%	16.8%	10.9%	10.7%	10.9%	12.2%
AEA	12.2%	12.1%	12.5%	11.2%	12.1%	11.2%	10.4%	13.3%	12.8%	14.2%	12.2%	16.9%	11.0%	10.8%	11.0%	12.3%
TEA	12.1%	12.0%	12.6%	11.2%	12.1%	11.1%	10.5%	13.2%	12.7%	14.1%	12.2%	16.9%	11.0%	10.8%	11.0%	12.2%
CRETTA	12.0%	<u>11.7%</u>	<u>12.4%</u>	10.7%	11.9%	10.5%	10.0%	12.7%	12.1%	13.6%	11.9%	16.6%	10.7%	<u>10.3%</u>	10.6%	11.9%

performance under severe corruptions. Notably, while PL exhibits substantial drops in performance under corruption, CRETTA remains robust and effective across both clean and corrupted settings, demonstrating its reliability in both distributions.

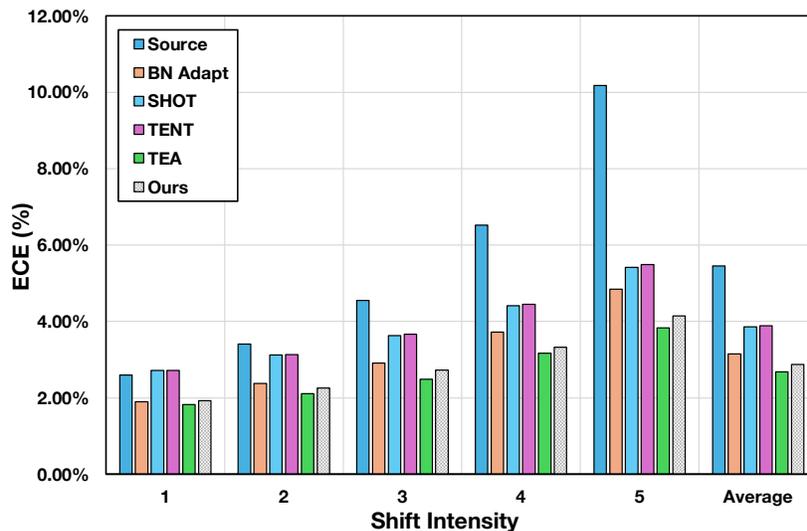


Figure 4: Comparison of Expected Calibration Error (ECE \downarrow) on the CIFAR10-C dataset across different corruption severity levels.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

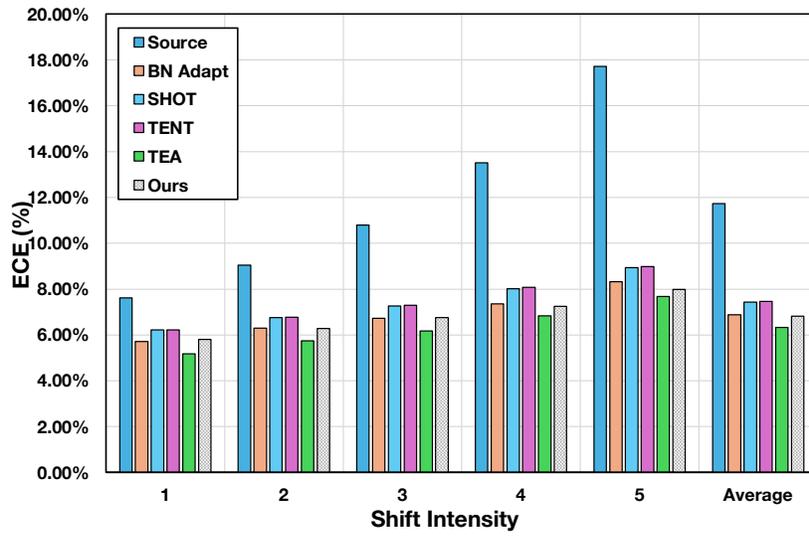


Figure 5: Comparison of Expected Calibration Error (ECE \downarrow) on the CIFAR100-C dataset across different corruption severity levels.

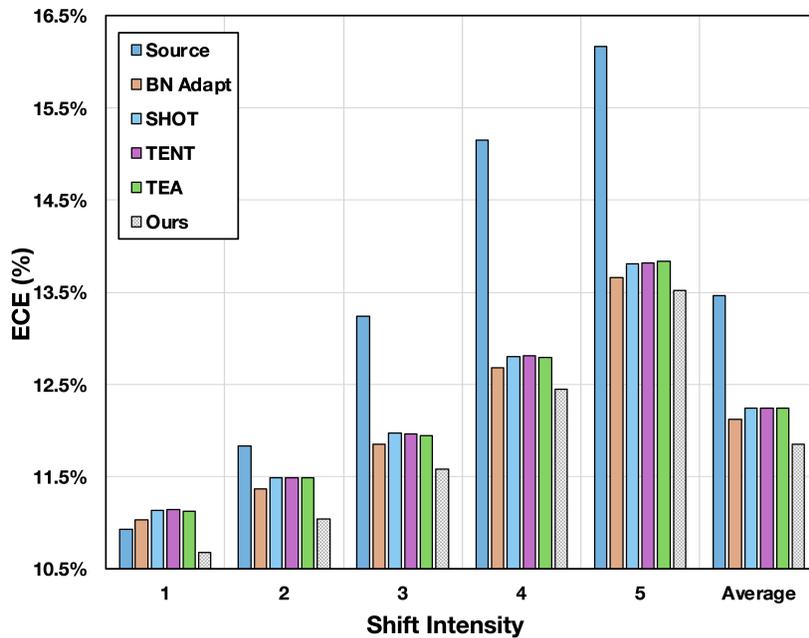


Figure 6: Comparison of Expected Calibration Error (ECE \downarrow) on the TinyImageNet-C dataset across different corruption severity levels.

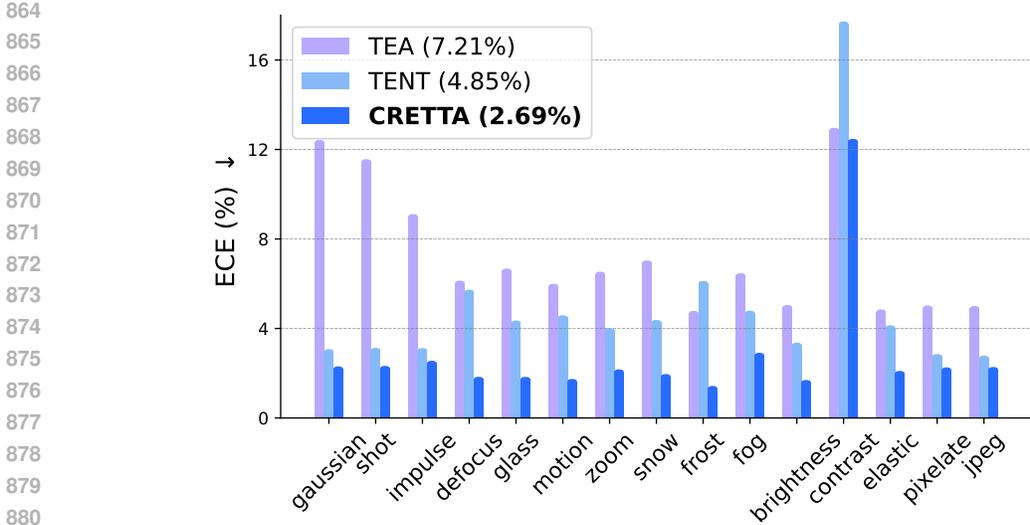


Figure 7: Comparison of Expected Calibration Error (ECE ↓) on ImageNet-C across various corruption types, with results averaged over severities 1–5.

Detailed Performance Comparison on Calibration Error In this section, we provide a detailed analysis of the Expected Calibration Error (ECE) for CIFAR10-C and CIFAR100-C. This expands upon the results shown in Table 1.

As seen in Figure 4 and Figure 5, energy-based methods such as TEA and CRETТА consistently outperform baseline approaches like TENT, which suffers from overconfidence issues. Furthermore, our method maintains computational advantage over TEA, making it more efficient while achieving comparable or superior performance.

On TinyImageNet-C dataset, shown in Figure 6, CRETТА outperforms all competing methods across all severity levels. This consistent superiority over all baseline methods demonstrates the robustness and adaptability of our approach in high-complexity datasets.

Table 10: Comparison of computational cost (GFLOPs), Memory Cost (Peak Memory Usage), and performance metrics (ECE and Acc) for baselines on the CIFAR10-C

	GFLOPs(↓)	Memory Cost(↓)	Acc(↑)	ECE(↓)
Source	131.53	443.98 MB	88.82%	5.45%
BN Adapt	131.53	452.61 MB	89.12%	3.15%
TENT	132.59	1546.05 MB	90.74%	3.89%
TEA	4335.82	3464.78 MB	90.67%	2.68%
Ours	527.40	2651.83 MB	91.01%	2.88%

Table 11: Comparison of computational cost (GFLOPs), Memory Cost (Peak Memory Usage), and performance metrics (ECE and Acc) for baselines on the CIFAR100-C

	GFLOPs(↓)	Memory Cost(↓)	Acc(↑)	ECE(↓)
Source	131.53	443.03 MB	64.11%	11.73%
BN Adapt	131.53	452.70 MB	65.83%	6.88%
TENT	132.59	1546.21 MB	68.73%	7.47%
TEA	4335.82	3465.00 MB	67.93%	6.33%
Ours	527.40	2651.85 MB	69.05%	6.82%

Detailed Performance Comparison on Computational Efficiency To further demonstrate the computational advantages of our proposed method, we present a comprehensive comparison of com-

Table 12: Comparison of classification accuracy (Acc \uparrow) and expected calibration error (ECE \downarrow) on the ImageNet-C dataset.

Method	Severity L5		Severity Avg	
	Acc \uparrow	ECE \downarrow	Acc \uparrow	ECE \downarrow
TENT	37.39%	7.75%	43.78%	4.85%
TEA	31.60%	8.39%	38.72%	7.21%
CRETТА	37.05%	4.43%	43.54%	2.69%

putational cost (GFLOPs), peak GPU memory usage (MB) with performance metrics (Accuracy and ECE) across CIFAR10-C, CIFAR100-C, and TinyImageNet-C, as summarized in Table 10, Table 11. Compared to TEA, which incurs substantial computational overhead due to SGLD-based sampling, CRETТА reduces GFLOPs by more than sevenfold across datasets. Furthermore, despite incorporating a source buffer, CRETТА maintains a modest peak GPU memory usage, significantly lower than TEA. The peak GPU memory usage is measured as the maximum allocated GPU memory during adaptation. Consequently, CRETТА offers a practical balance between performance, computational cost, and memory efficiency, making it well-suited for deployment in real-world, resource-constrained environments.

Scalability In this section, we provide a detailed results on ImageNet-C.

As shown in Table 12, CRETТА achieves performance by a significant margin, outperforming the entropy-based method TENT and the existing energy-based method TEA in ECE.

Entropy minimizations’s overconfidence and MLE-based approach’s approximation error introduced when estimating its normalization constant term leads to poor calibration which is inappropriate in real-world TTA scenarios. In contrast, CRETТА generalizes well to large-scale datasets, achieving strong predictive performance with superior calibration.

Table 13: Comparison of classification accuracy on CIFAR10(-C), CIFAR100(-C) under gradual distribution shift

Domain	CIFAR10			CIFAR100		
	OURS	TEA	TENT	OURS	TEA	TENT
Source (Q)	93.46	93.45	93.43	73.97	73.88	73.57
1	92.88	92.80	92.77	71.90	71.41	71.70
2	92.03	91.92	91.92	71.57	70.40	71.36
3	91.63	91.29	91.35	69.99	67.71	70.04
4	90.25	89.81	90.03	67.99	65.23	68.28
5 (P)	89.47	88.78	88.58	65.47	60.26	65.23

Detailed Performance Comparison Under Gradual Shift scenario In subsection 4.3, we demonstrated that our contrastive residual energy-based learning shows superior performance over CD MLE-based adaptation method TEA. This tendency was consistently observed under the gradual distribution shift setting in Table 13, and here we additionally report comparisons with TENT.

For CIFAR10-C, CRETТА maintains the best performance throughout the shift. For CIFAR100-C, CRETТА shows clear gains under stronger shifts. At severity 5, it achieves 65.47%, notably higher than TEA(60.26%) and TENT (65.23%). While TENT is competitive at mid-level severities, it degrades more under severe shifts. Overall, CRETТА provides robust adaptation across gradual shifts while preventing forgetting, outperforming both TEA and TENT.

Test-time Adaptation for Non-IID Settings

Our previous experiments are conducted under the assumption of i.i.d. test samples which is a widely adopted setting in prior work. Nonetheless, real-world applications can also encounter non-i.i.d. samples (Gong et al., 2022; Yuan et al., 2023; Wang et al., 2022). To further examine the robustness and generalizability of our method beyond the i.i.d. assumption, we constructed a non-i.i.d. test-time adaptation scenario. Specifically, we simulated non i.i.d. data stream by leveraging a Dirichlet distribution to control the class allocation ratio within batch, denoted as δ . A higher δ value brings the distribution closer to i.i.d., whereas a lower δ value results in a more non-i.i.d. distribution, where a specific class might dominate the batch. We conducted our experiment on CIFAR100-C using the WRN-28-10 backbone.

As Table 14 shows, our method consistently outperforms entropy minimization and instance selection approaches across all δ values. Specifically, CRETТА achieves the highest average accuracy of 60.99%, surpassing TENT’s 58.85% by 2.14%p. Also, even at the most imbalanced setting where $\delta = 0.01$, our method achieves a competitive accuracy of 48.33%. These findings demonstrate that our method not only excels in i.i.d. scenarios but also is effective in dynamic real-world environments.

A.3 ABLATION STUDY

A.3.1 DETAILED ABLATION STUDY

Table 15: Comparison of classification accuracy(Acc) and expected calibration error(ECE) on benchmark datasets between CRETТА(Default) and CRETТА(Loss Term without Source Model) at severity level 5.

Method	CIFAR10-C		CIFAR100-C		TinyImageNet-C	
	Acc(\uparrow)	ECE(\downarrow)	Acc(\uparrow)	ECE(\downarrow)	Acc(\uparrow)	ECE(\downarrow)
CRETТА	88.30	4.15	64.52	7.99	40.30	13.52
w.o Source Model Term	88.09	4.66	60.02	5.93	37.46	14.55

Loss Ablation We observed that eliminating the source model consistently degraded both accuracy and calibration (ECE) in most cases across our benchmark datasets. These results collectively demonstrate that incorporating source model related terms into our contrastive residual learning is essential for stable adaptation.

Gradient Ablation The gradient coefficient $w(x_t, x_s)$ is the key mechanism that turns relative energy into stable updates. To verify this role, we conducted an ablation study that disrupts the proposed weighting scheme by replacing $w(x_t, x_s)$ with values randomly sampled from a uniform distribution $[0, 1)$. As shown in Table 16, this replacement lead to lower accuracy and higher calibration error, confirming that gradient coefficient is critical for stable optimization and robust adaptation under noisy target data.

Table 14: Test-time adaptation in dynamic scenarios using CIFAR100-C at severity 5. Our method demonstrates higher robustness compared to baselines across varying the allocation ratio δ .

Method	$\delta = 10$	$\delta = 1$	$\delta = 0.1$	$\delta = 0.01$	Avg Acc.
BN Adapt	61.44%	61.11%	59.02%	45.61%	56.79%
PL	44.03%	37.27%	39.06%	43.38%	40.93%
SHOT	<u>63.94%</u>	<u>63.60%</u>	<u>61.20%</u>	46.54%	58.82%
TENT	63.91%	63.56%	<u>61.20%</u>	<u>46.72%</u>	58.85%
ETA	62.31%	62.04%	59.89%	46.04%	57.57%
EATA	62.35%	62.04%	59.84%	46.04%	57.57%
SAR	61.54%	61.22%	59.12%	45.66%	56.89%
TEA	62.58%	62.29%	60.08%	46.22%	57.79%
Ours	66.20%	65.95%	63.47%	48.33%	60.99%

1026 Table 17: Comparison of classification accuracy(Acc) and expected calibration error(ECE) on
 1027 benchmark datasets between CRETТА(Default) and CRETТА(Single Source Sample in Buffer) at
 1028 severity level 5.
 1029

Method	CIFAR10-C		CIFAR100-C		TinyImageNet-C	
	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)
CRETТА	88.30	4.15	64.52	7.99	40.30	13.52
CRETТА with single source sample	87.62	5.39	62.67	8.93	40.30	14.13

1036
 1037 **Extended Buffer Ablation** While the specific content of the buffer has less impact on performance,
 1038 as shown in Table 5, this does not imply that the source buffer itself plays a trivial role. To further
 1039 verify this, we additionally conducted an experiment where the buffer consists of only a single source
 1040 sample. As shown in Table 17, accuracy dropped by up to 1.7% and ECE increased by up to 1.2%
 1041 across datasets.

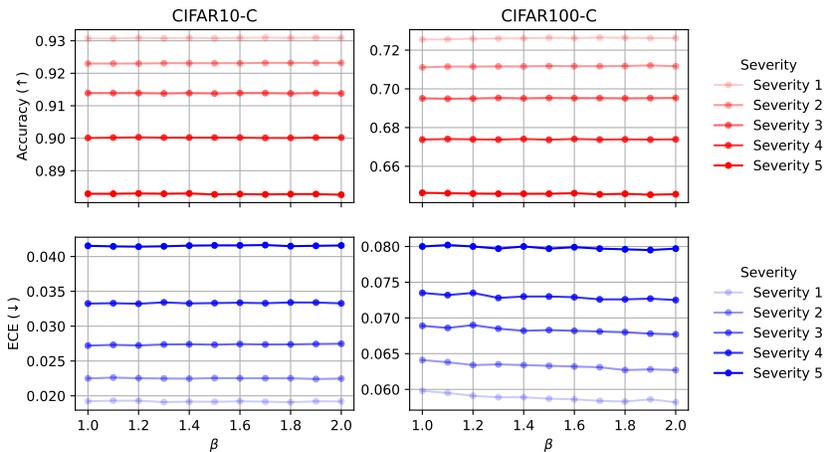
1042 Table 18: Effectiveness of preference pair size on CIFAR10-C, CIFAR100-C, and TinyImageNet-C.
 1043

	CIFAR10-C	CIFAR100-C	TinyImageNet-C
CRETТА wo/ CP	88.30%	64.52%	40.30%
CRETТА w/ CP	88.24%	64.69%	40.44%

1044
 1045
 1046
 1047
 1048
 1049 **Pair Size Ablation** In CRETТА, we assume that the samples in a test batch represent the target
 1050 distribution, while the source replay buffer represents the source distribution. The loss is computed
 1051 by forming pairs between target and source samples within each batch, enabling a direct comparison
 1052 between the two distributions.
 1053

1054 To demonstrate the assumption is valid, we examined the impact of increasing the number of pair
 1055 combinations using a Cartesian Product (CP) to generate all possible combinations of target and
 1056 source data within each batch. For example, we use 200 pairs for each adaptation in CIFAR10-C,
 1057 while the Cartesian Product results in 200×200 pairs.

1058 Our results across three datasets summarized in Table 18 indicate that generating more pairs does not
 1059 necessarily lead to performance gain. With only a few pairs, CRETТА can efficiently adapt to the
 1060 target distribution.
 1061



1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077 Figure 8: Ablation on varying β values on CIFAR10-C and CIFAR100-C at severity 1-5.
 1078
 1079

Hyperparameter β Ablation The hyperparameter β in Equation 3 controls the deviation from the pretrained source model, serving as a scaling parameter. To evaluate the robustness of our method, we experiment its performance across varying values of β , assessing both accuracy and expected calibration error (ECE) on CIFAR10-C, CIFAR100-C and TinyImageNet-C. As shown in Figure 8, our method consistently demonstrates stable performance across all corruption severity levels (1-5), validating its robustness.

In addition, we further examine the effectiveness of CRETТА across varying values of the hyperparameter β on TinyImageNet-C, averaging results over severity levels 1 to 5, and compare its performance against competitive baselines (see Figure 9). These results confirm that the strong adaptation performance of CRETТА is not reliant on a specific setting of the temperature parameter β , but rather stems from our contrastive residual learning objective itself.

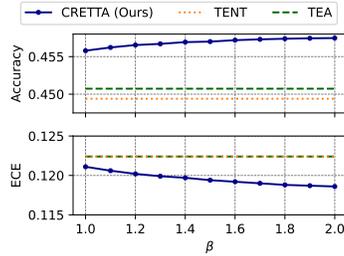


Figure 9: Ablation on varying values of β .

A.3.2 DETAILED SETTING OF CRETТА

Table 19: Detailed hyperparameters settings for each dataset.

Dataset	LR	β	Batch Size	Transformation Type (probability)
CIFAR10-C	1e-3	1.0	200	rotate(1.0)
CIFAR100-C	2e-3	2.0	200	flip, rotate, affine, perspective, crop(0.2)
TinyImageNet-C	1e-3	2.0	1000	None

Hyperparameters This section details the hyperparameter settings for CRETТА. To optimize performance, minimal hyperparameter tuning was conducted, focusing solely on learning rate, β and type and probability of random transformations for source buffer. With only slight adjustments, CRETТА achieved significantly better performance than the current state-of-the-art (SOTA). The batch sizes were aligned with the default settings used in TENT and TEA, which are 200 for CIFAR10-C and CIFAR100-C, 1000 for TinyImageNet-C. For ImageNet-C we follow TENT default settings, using a batch size of 64 and learning rate of $2.5e-4$. These settings ensured consistency across experiments while highlighting the robustness and effectiveness of CRETТА. For the PACS domain-generalization task, we used a learning rate of $1e-3$, a batch size of 100, applying source-sample augmentation in the same way as for CIFAR100-C. All experiments were conducted using a single NVIDIA RTX A6000 GPU (48GB).

Evaluation Metrics Expected Calibration Error (ECE) (Guo et al., 2017) is a metric used to measure the calibration quality of a probabilistic model. Calibration refers to how closely the predicted probabilities of a model match the actual probabilities. ECE quantifies the discrepancy between predicted confidence and actual accuracy. ECE is calculated as shown in Equation 7:

$$\text{ECE} = \sum_{m=1}^M \frac{|\text{bin}_m|}{N} \cdot |\text{confidence}_m - \text{accuracy}_m| \quad (7)$$

where M is the number of bins, N is the total number of data points, bin_m is the number of predictions in m -th bin, and confidence_m and accuracy_m are the confidence and accuracy of bin m , respectively.

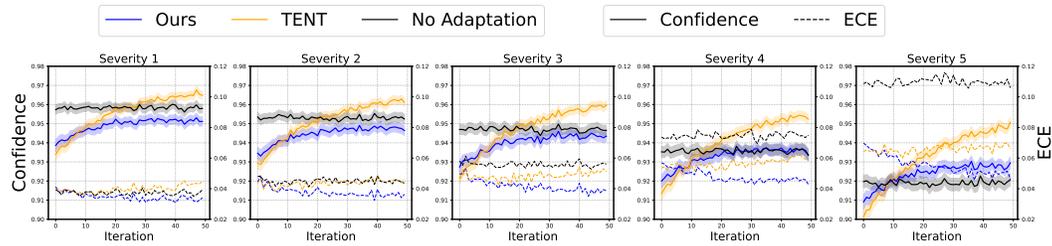


Figure 10: The overconfidence problem of entropy minimization in test-time adaptation on CIFAR10-C. TENT tends to increase a model’s confidence in uncertain predictions as adaptation progresses, often leading to worse calibration due to overconfidence. In contrast, CRETТА (Ours) stabilizes the adaptation process by gradually reducing the expected calibration error.

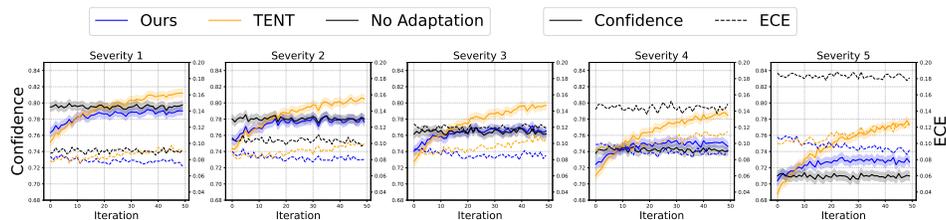


Figure 11: The overconfidence problem of entropy minimization in test-time adaptation on CIFAR100-C.

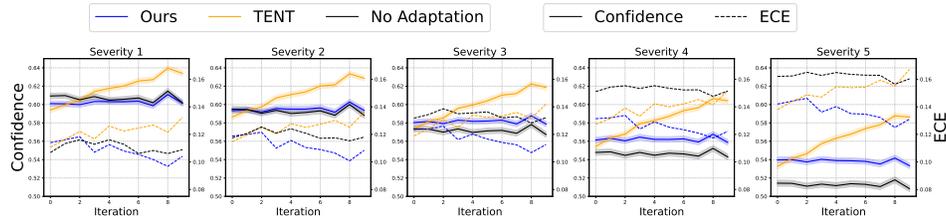


Figure 12: The overconfidence problem of entropy minimization in test-time adaptation on TinyImageNet-C.

A.3.3 DETAILED RESULTS FOR OVERCONFIDENCE PROBLEM OF ENTROPY MINIMIZATION

The overconfidence issue inherent in entropy minimization has been thoroughly investigated in prior works (Liu et al., 2020; Hendrycks & Gimpel, 2016; Guo et al., 2017). Building on this, we explored that increasing a model’s prediction confidence especially when the label information is unavailable can lead to bad calibration as shown in Figure 4. The trend consistently appears in other benchmark datasets including CIFAR100-C and TinyImageNet-C as illustrated in Figure 11 and Figure 12. Entropy minimization raises the model’s confidence across all severity levels, with the rate of increase becoming steeper as corruption severity intensifies, thereby exacerbating error accumulation.

On the other hand, CRETТА maintains stable confidence managing uncertainty during test-time adaptation and even reduces calibration error as adaptation progresses. These results suggest that maximizing the marginal likelihood of target samples provides a safer and more effective strategy compared to relying on uncertain predicted probabilities $p_{\theta}(\hat{y}|x)$ in the test-time learning objective.

B NOISE CONTRASTIVE ESTIMATION

We first define a reward function $r(\cdot)$ to properly compare samples from two different sets or distributions.

$$r(x; \theta, \phi) = \log P_{\theta}(x) - \log P_{\phi}(x)$$

1188 where P_θ is the target distribution and P_ϕ is the source distribution.

1189

1190 B.1 NON-RESIDUAL

1191

1192 If we define energy functions for each of them by utilizing gibbs distribution,

1193

$$E_\theta(x) = -\log P_\theta(x) - \log Z(\theta)$$

1194

$$E_\phi(x) = -\log P_\phi(x) - \log Z(\phi)$$

1195

1196 Then the reward function becomes

1197

$$r(x; \theta, \phi) = -(E_\theta(x) - E_\phi(x)) + C$$

1198

1199 Then the loss function becomes

1200

$$\mathcal{L}(\theta; \phi) = -\mathbb{E}_{x_t} [\log \sigma(r(x; \theta, \phi))] - \mathbb{E}_{x_s} [\log(1 - \sigma(r(x; \theta, \phi)))]$$

1201

1202

1203

1204

1205

1206 B.2 RESIDUAL

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1218 C PAIR-WISE CONTRASTIVE ESTIMATION

1220 We first define a reward function $r(\cdot)$ to properly compare samples from two different sets or

1221 distributions.

1222

1223

$$r(x_t, x_s) = \tilde{r}(x_t) - \tilde{r}(x_s)$$

1224

1225 where \tilde{r} is a reward function that assigns higher values to target samples than source samples

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1227 If we define energy functions for each of them,

1228

1229

$$E_\theta(x) = -\log P_\theta(x) - \log Z(\theta)$$

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1230 Then the reward function becomes

1231

1232

$$r(x_t, x_s) = \log P_\theta(x_t) - \log P_\theta(x_s) = -(E_\theta(x_t) - E_\theta(x_s))$$

1233

1234

1235

1236

1237

1238

1239

1240

1241

1233 Then the loss function becomes

1234

1235

1236

1237

1238

1239

1240

1241

$$\begin{aligned} \mathcal{L}(\theta; \phi) &= -\mathbb{E}_{x_t, x_s} [\log \sigma(r(x_t, x_s))] \\ &= -\mathbb{E}_{x_t, x_s} [\log \sigma(-(E_\theta(x_t) - E_\theta(x_s)))] \end{aligned}$$

1237 The gradient becomes

1238

1239

1240

1241

$$\begin{aligned} \nabla_\theta \mathcal{L}(\theta; \phi) &= -\mathbb{E}_{x_t, x_s} [\sigma(-r(x_t, x_s)) \nabla_\theta r(x_t, x_s)] \\ &= \mathbb{E}_{x_t} [\sigma(E_\theta(x_t) - E_\theta(x_s)) (\nabla_\theta E_\theta(x_t) - \nabla_\theta E_\theta(x_s))] \end{aligned}$$

1242 D USE OF LLMs

1243
1244 We used a large language model (ChatGPT) only as a general purpose assistive tool for minor editing
1245 tasks such as polishing sentences, correcting grammar and spelling and making small LaTeX table
1246 formatting adjustments. The LLM was not involved in research ideation, experimental design, data
1247 analysis, or substantive writing. All technical decisions, interpretations, and the writing of the core
1248 content were carried out entirely by the authors, who take full responsibility for the originality of the
1249 manuscript.

1251 E REBUTTAL

1252
1253 **Clarifying Why Energy Minimization Improves Generalization and How It Differs from En-**
1254 **trophy Minimization** In this section, we clarify (1) why energy minimization improves generaliza-
1255 tion capability and (2) how this approach differs, especially from an optimization perspective, from
1256 entropy minimization, which is known to suffer from overconfidence issues.

1257 First, to understand why minimizing the energy of target samples leads to improved generalization
1258 capability (i.e., higher classification accuracy) on the target distribution, it is essential to examine
1259 how optimizing the energy function reshapes the representations.

1260 We assume a classifier $f_\theta(x) = g(h_\theta(x))$ with feature extractor h_θ and a frozen linear classifier
1261 $g(z) = Wz + b$, where $W = [w_1, \dots, w_C]^\top \in \mathbb{R}^{C \times d}$. The logits can be expressed as

$$1262 a_k(x) = w_k^\top h_\theta(x) + b_k.$$

1263
1264 Since target samples contain no labels, we optimize the marginal energy of the input,

$$1265 E_\theta(x) = -\log \sum_{k=1}^C \exp(a_k(x)),$$

1266
1267 which is the standard unnormalized negative log-density in energy-based models. The gradient of
1268 this energy with respect to the logit $a_k(x)$ is

$$1269 \frac{\partial E_\theta(x)}{\partial a_k(x)} = -\frac{\exp(a_k(x))}{\sum_{k'=1}^C \exp(a_{k'}(x))} = -p_\theta(y = k | x).$$

1270
1271 Using $\frac{\partial a_k}{\partial z} = w_k$, the gradient of the energy with respect to the feature representation becomes

$$1272 \frac{\partial E_\theta(x)}{\partial z} = \sum_{k=1}^C \frac{\partial E_\theta(x)}{\partial a_k} \frac{\partial a_k}{\partial z} = -\sum_{k=1}^C p_\theta(y = k | x) w_k = -\mathbb{E}_{p_\theta(y|x)}[w_y].$$

1273
1274 This expression shows that the gradient descent updates shift the feature z toward the expected
1275 classifier weight vector. Although the classifier remains frozen, this shift alters the logits and thus
1276 reshapes the conditional distribution, enabling the model to align its predictions to the target domain
1277 solely through feature-level adjustments. From a representation perspective, energy minimization
1278 adapts the feature extractor to better capture the target-domain distribution, producing stronger
1279 representations that improve classification accuracy under distribution shift even without labels.

1280 Furthermore, to understand how energy minimization and entropy minimization objectives behave
1281 differently during optimization, it is essential to compare their gradients. Since we already expressed
1282 the gradient of the energy objective, we now present the gradient expressions for the entropy objective.

1283 For unlabeled target data, the entropy of model prediction is expressed as

$$1284 \mathcal{L}_{\text{ent}}(x) = -\sum_{k=1}^C p_k \log p_k,$$

1285 where $p_k = p_\theta(y = k | x)$. The corresponding gradient is

$$1286 \frac{\partial \mathcal{L}_{\text{ent}}}{\partial z} = -\sum_{y=1}^C (\log p_y + 1) p_y \left(w_y - \sum_{k=1}^C p_k w_k \right).$$

Letting $\mathbb{E}[w] = \sum_k p_k w_k$ denote the expectation of classifier weights under the current predictive distribution, we obtain the final compact form:

$$\frac{\partial \mathcal{L}_{\text{ent}}}{\partial z} = - \sum_{y=1}^C (\log p_y + 1) p_y (w_y - \mathbb{E}[w]).$$

In the gradient, the term $(\log p_y + 1) p_y$ heavily weights classes for which p_y is already large and $\log p_y$ is less negative. Thus, the gradient of entropy moves the feature z in the direction that increases confidence for the most likely classes, effectively reducing prediction entropy and behaving similarly to pseudo-labeling. By directly modifying the conditional distribution, entropy minimization mainly pushes the model to become more confident, often excessively.

In contrast, the gradient of the energy objective depends only on $\mathbb{E}_{p_{\theta}(y|x)}[w_y]$, which is a smooth expectation over classifier weights. It does not contain the entropy term’s confidence-amplifying multiplier. This trend is also empirically confirmed in Figure 10, Figure 11, and Figure 12.

Overall, energy minimization improves generalization by altering the logits and reshaping the conditional distribution, allowing the model to align its predictions to the target domain. While achieving strong classification performance, it is also more robust than entropy minimization because it optimizes a smooth expectation over classifier weights, avoiding the entropy objective’s confidence-amplifying multiplier. Together, these properties make energy-based adaptation a more balanced and principled optimization approach, enhancing representation quality rather than merely increasing confidence, which explains its superior robustness in test-time adaptation.

Extended Buffer Ablation While the specific content of the buffer has less impact on performance, as shown in Table 5, this does not imply that the source buffer itself plays a trivial role. To further verify this, we additionally conducted an experiment where the buffer consists of only a single source sample. As shown in Table 17, accuracy dropped by up to 1.7% and ECE increased by up to 1.2% across datasets.

To understand why this is the case, recall that source energy (i.e., $E(x_s)$) is required not merely because of residual learning, but because it is a crucial component of the pairwise contrastive objective. Using an arbitrary reference energy as source energy would likely cause training failure. From a gradient perspective, Equation 4 shows that during early adaptation, a high source energy drives the gradient weight w to collapse to zero, resulting in a trivial solution where learning cannot proceed—mirroring the well-known constraint in Noise Contrastive Estimation (NCE) regarding the choice of the noise distribution. To enable effective early adaptation, the buffer must therefore contain samples drawn from a distribution similar to that learned by the pretrained source model, ensuring that these samples receive low energy and provide meaningful contrastive learning signals.

Table 20: Performance Comparison of Source Buffer Contents on CIFAR100-C

Buffer Type	Severity 5	Severity 1–5
CRETTA (ours)	64.52%	69.05%
CIFAR-10 (train)	64.97%	69.37%
CIFAR-10 (val)	64.95%	69.37%
PACS (sketch)	63.74%	68.25%

This behavior is confirmed empirically. As shown in Table 20, CRETTA maintains strong performance on CIFAR100-C when the buffer contains CIFAR-10, which is not original source data but is distributionally similar. In contrast, performance deteriorates when the buffer contains samples with a very different distribution, such as PACS.

Thus, Table 6 should not be interpreted as suggesting that the absolute source-energy distribution plays a minor role. Rather, it demonstrates that CRETTA remains robust and consistently effective as long as the buffer contains data that are distributionally similar to the source distribution—even without access to the original source dataset. In practice, while true source data may be inaccessible due to privacy constraints, obtaining similar samples is usually far more feasible.

Overall, CRETТА’s superior performance does not arise merely from lowering target energy; instead, it results from the interplay between the residual formulation, the pairwise contrastive objective, and their integration, which together provide a stable and powerful adaptation mechanism.

Contrastive Component Ablation In this section, we clarify why the contrastive component is necessary and how it contributes to stable adaptation.

- **(1) Why — Contrastive learning is necessary for well-calibrated and stable adaptation.** We demonstrate this by *removing the contrastive term*, showing that direct minimization of target energy leads to unstable energy collapse and degraded calibration across benchmarks.
- **(2) How — Contrastive learning provides a more informative gradient signal.** We validate this by *replacing source-buffer samples with target samples*, showing that target-only regularization produces weak gradients and yields only marginal adaptation.

First, to empirically validate whether contrastive learning is indeed essential for stable adaptation, we first conducted an additional experiment in which we removed the contrastive term and measured the resulting calibration error across the three benchmark datasets.

Table 21: Ablations on the contrastive component across benchmark datasets (ECE).

	CIFAR10-C Sev5	CIFAR10-C Sev1-5	CIFAR100-C Sev5	CIFAR100-C Sev1-5	TinyIN-C Sev5	TinyIN-C Sev1-5
W Contrastive Terms (CRETТА)	4.15%	2.88%	7.99%	6.82%	13.52%	11.85%
W/O Contrastive Terms	5.57%	4.08%	11.61%	9.65%	16.21%	14.12%

As shown in Table 21, removing the contrastive terms (i.e., $E(x_s)$) and directly minimizing the target energy led to a consistent degradation in calibration across all benchmarks. The effect was particularly pronounced on CIFAR100-C, where the calibration error deteriorated to 11.61% (+4%p), which is notably worse than TENT (8.93%), an entropy minimization-based method that is prone to overconfidence. These results clearly demonstrate that the stability of CRETТА’s adaptation does not arise simply from reducing target energies, but instead stems from the contrastive learning mechanism.

To further analyze how the contrastive term contributes to stable adaptation, we also examined the behavior of energy levels and calibration error on CIFAR100-C (Severity 5) during adaptation, comparing CRETТА against the variant where the contrastive terms are ablated.

Table 22: Target energy and ECE of CRETТА vs. without contrastive terms during adaptation on CIFAR100-C (Severity 5).

Batch Idx	Target Energy		ECE	
	W Contrastive (OURS)	W/O Contrastive	W Contrastive (OURS)	W/O Contrastive
0	-9.9781	-9.9781	11.11%	11.56%
9	-10.1208	-10.9536	10.48%	11.64%
19	-10.2124	-11.5247	9.48%	12.97%
29	-10.1896	-11.7863	9.85%	12.35%
39	-10.1856	-12.0177	8.88%	13.35%
49	-10.1904	-12.2531	9.15%	12.92%
Δ (Last-First)	-0.21	-2.28	-1.96%	+1.35%

As shown in Table 22, we observe that when the contrastive terms are removed and the model directly minimizes the target energy, the energy level drops rapidly during the early stages of adaptation. While this may facilitate fast initial adaptation, it poses a critical risk to stability since aggressively lowering target energies early on sharpens the energy landscape around target samples, which can lead to overfitting. Empirically, we indeed find that removing the contrastive terms results in an overall increase in calibration error.

In contrast, CRETТА reduces the energy level progressively within the contrastive learning framework, enabling a more stable adaptation trajectory. This gradual reduction improves calibration error over time, demonstrating that the contrastive mechanism plays a key role in stabilizing adaptation dynamics.

Table 23: Gradient coefficient w of CRETТА and Target-as-Source Buffer Data during adaptation on CIFAR100-C (Severity 5).

Batch Idx	CRETТА	W TRG as SRC
0	0.490	0.304
9	0.535	0.406
19	0.580	0.456
29	0.580	0.465
39	0.627	0.518
49	0.608	0.522
AVG	0.572	0.454

Table 24: Performance comparison of Target-as-Source buffer setting on CIFAR100-C (Severity 5).

Method	Sev5 Acc	Sev5 ECE
BN Adapt	60.74%	8.32%
W TRG as SRC	61.39% (+0.65%p)	8.19%
CRETТА	64.52% (+3.78%p)	7.99%

Furthermore, the contrastive learning mechanism proposed in CRETТА that utilizes a small buffer of source (or distributionally similar) samples provides a more meaningful gradient signal than using target samples for regularization. In our previous comment, we illustrated this using an extreme scenario where the source energy, driven by high-energy target samples, becomes sufficiently large that the gradient effectively vanishes. Here, we provide a more realistic explanation focusing on the relative magnitudes of the energy level differences.

Target samples within the same batch are drawn from the same underlying distribution, and so thus it is unlikely for their energy levels to differ significantly. In contrast, source (or distributionally similar) samples originate from distribution that the pretrained model has already learned, making them more likely to exhibit consistently lower energy values than newly encountered target samples. This creates a meaningful energy gap from the target sample energies, which in turn provides a strong gradient signal during adaptation. Consequently, CRETТА’s contrastive learning framework yields notable gains in classification performance while achieving better calibration.

CRETТА provides a substantially more informative learning signal than simply replacing source samples with low-energy target samples. To validate this, we constructed an experimental setup where the 50% of samples with the lowest model-computed energy values in each target batch served as source-buffer data to compute the source energy $E(x_s)$. We then compared this setup with CRETТА’s adaptation process and performance.

More concretely, we first compared the magnitude of the gradient coefficient w throughout the adaptation process. As shown in Table 23, replacing source samples with low energy target samples results in consistently smaller gradient coefficients than CRETТА across the entire adaptation trajectory. This indicates that the model receives weaker learning signals and therefore fails to sufficiently adapt to the target distribution. Consequently, as shown in Table 24, the accuracy improvement is only marginal amounting to just +0.65 percentage points compared to BN adapt, which performs adaptation solely through normalization without any learning. In contrast, CRETТА maintains a relatively meaningful gradient coefficient while gradually increasing the learning signal from the early, high-uncertainty stages of adaptation toward later stages. This leads to both improved classification performance and better calibration, ultimately achieving effective and stable adaptation.

Overall, CRETТА’s contrastive learning is an essential component for achieving well-calibrated and stable test-time adaptation. By progressively lowering target energy and thereby reducing calibration error, it enables a stable adaptation process. Moreover, CRETТА’s contrastive learning methodology, which leverages buffer data, is distinctly more effective than approaches that apply only target-sample-based regularization. Notably, CRETТА’s contrastive framework is also robust to variations in buffer-data content and quality, making it highly practical for real-world deployment.