CONTRASTIVE RESIDUAL ENERGY TEST-TIME ADAPTATION

Anonymous authorsPaper 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 $-\sum_{k=1}^C p_\theta(\hat{y}_k \mid x_t) \log p_\theta(\hat{y}_k \mid x_t)$, where C is the number of classes and $p_\theta(\hat{y}_k \mid x_t)$ 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

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

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

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

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

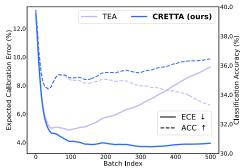


Figure 1: ImageNet-C (Sev 5) ECE(\downarrow) and Acc(\uparrow) over batch progress. CRETTA maintains stable calibration performance, while TEA experiences approximation error accumulation.

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

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

2 Preliminaries

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

Test-time adaptation (TTA) aims to mitigate this issue by adapting the pretrained parameters ϕ to better align with the target marginal distribution P. In this setting, we are given a set of N unlabeled target samples $\{x_i^{(i)}\}_{i=1}^N$ drawn from P, which arrive in online batches. To cope with the 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, CRETTA, 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 $Z(\phi) = \int_x \exp(-E_{\phi}(x))dx$ is the normalizing constant (partition function). The energy function $E_{\phi}(x): \mathbb{R}^d \to \mathbb{R}$ 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 $\log -\sup f$ of the logits across all classes in the final classifier layer:

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

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

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:

$$\mathcal{L}(\theta; \phi, \mathcal{B}) = -\frac{1}{|\mathcal{B}|} \sum_{(x_s, x_t) \sim \mathcal{B}} \log \sigma \left(\beta \underbrace{\left(\log \frac{p_{\theta}(x_t)}{q_{\phi}(x_t)} - \log \frac{p_{\theta}(x_s)}{q_{\phi}(x_s)} \right)}_{l} \right),$$
where $l = -(E_{\theta}(x_t) - E_{\theta}(x_s)) + (E_{\phi}(x_t) - E_{\phi}(x_s)).$ (3)

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 $\mathcal{B}_s = \{x_s^{(i)} \mid i=1,\ldots,|\mathcal{B}_s|\}$, 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 subsubsection 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:

$$\nabla_{\theta} \mathcal{L}(\theta; \phi, \mathcal{B}) = -\frac{1}{|\mathcal{B}|} \sum_{(x_s, x_t) \sim \mathcal{B}} \beta \cdot w(x_t, x_s) \cdot \left(\nabla_{\theta} \log p_{\theta}(x_t) - \nabla_{\theta} \log p_{\theta}(x_s) \right),$$
where $w(x_t, x_s) = \sigma \left(\beta \log \frac{p_{\theta}(x_s)}{q_{\phi}(x_s)} - \beta \log \frac{p_{\theta}(x_t)}{q_{\phi}(x_t)} \right)$

$$= \sigma \left(\beta \cdot \left(E_{\theta}(x_t) - E_{\phi}(x_t) \right) - \beta \cdot \left(E_{\theta}(x_s) - E_{\phi}(x_s) \right) \right).$$
(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 <u>subsection 4.3</u> 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, CRETTA, 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.

	CIFAR-10-C			CIFAR-100-C			TinyImageNet-C					
	Sever	ity L5	Sever	ity Avg	Sever	ity L5	Sever	ity Avg	Sever	ity L5	Severi	ty Avg
Method	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)	Acc(↑)	ECE(↓)
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%
Normaliza	ıtion											
BN Adapt	85.46%	4.85%	89.12%	3.15%	60.74%	8.32%	65.83%	6.88%	39.60%	13.66%	44.72%	12.12%
Pseudo Lo	ibeling											
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%	64.41%	8.93%	68.80%	7.44%	39.84%	13.81%	44.95%	12.24%
Entropy N	1inimizatio	n										
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	88.27%	5.09%	90.88%	3.73%	64.40%	9.16%	68.75%	7.61%	39.87%	13.82%	44.97%	12.25%
Energy-be	ised Model	ls .										
TEA	88.06%	3.83%	90.67%	2.68%	63.66%	7.68%	67.93%	6.33%	39.96%	13.84%	45.08%	12.24%
CRETTA	88.30%	4.15%	91.01%	2.88%	64.52%	7.99%	69.05%	6.82%	40.30%	13.52%	45.75%	11.85%

Baselines We compared our method with state-of-the-art approaches. (i) Source: The pretrained 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. Pretrained 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 CRETTA on PACS and ImageNet-C datasets. On PACS, CRETTA 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, CRETTA 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 CRETTA generalizes well to large-scale, style shifted datasets, achieving strong predictive performance and superior calibration.

Table 2: Classification accuracy (Acc ↑) on PACS.

Source	Method		Avg			
Domain		Photo	Art	Cartoon	Sketch	
	Source	-	55.39	22.61	23.17	33.72
Photo	TENT TEA	-	57.89 53.16	63.82 50.33	32.84 33.17	51.52 45.55
	CRETTA	-	64.45	64.83	29.89	53.06
	Source	88.58	-	49.26	35.06	57.64
Art	TENT TEA	92.10 83.95	-	63.18 53.38	35.23 36.52	63.50 57.95
	CRETTA	91.52	-	65.90	40.99	66.13
	Source	70.52	62.52	-	50.44	61.16
Cartoon	TENT TEA	82.34 76.01	64.63 53.40	-	41.47 33.30	62.81 54.24
	CRETTA	83.65	68.86	-	40.04	64.18
	Source	13.89	14.50	19.03	-	15.81
Sketch	TENT TEA	31.88 18.60	30.99 25.47	49.60 48.34	-	37.49 30.80
	CRETTA	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
CRETTA	37.42	28.22	26.65	51.68	35.99

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 CRETTA (34K GFLOPs) but also struggles to maintain competitive performance. In contrast, CRETTA enables adaptation without explicitly tracking the normalization constant within a pair-wise contrastive learning framework. As shown in Figure 2, CRETTA consistently outperforms comparison methods, including TENT and BN-Adapt, while maintaining relatively low GFLOPs, demonstrating its efficiency for real-time TTA.

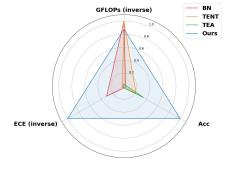


Figure 2: Comparison of GFLOPs, ECE and Acc against competitive baselines on TinyImageNet-C at the average across severity levels 1-5.

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, CRETTA demonstrates robust adaptation to progressively diverging target distributions,

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

Domain	CIFA	R10	CIFAR100		
20114111	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	

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

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

While the observed energy reduction is meaningful, the core strengths of CRETTA lie in contrastive residual energy learning to prevent convergence to trivial solutions as discussed in section 3. The bias terms with respect to E_{ϕ} in Equation 3 prevent E_{θ} from becoming excessively small or large, thereby stabilizing the adaptation process. As shown in Figure 3, the energy of target samples increases drastically after adaptation without bias terms, resulting in performance degradation.

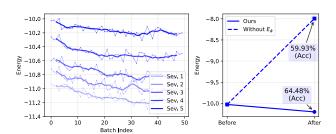


Figure 3: Energy trajectories of target samples across severities (left) and the effect of relative residual energy learning on stable adaptation (right) in CIFAR100-C.

Consequently, optimization proceeds

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

4.4 ABLATION STUDIES

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

The first question becomes particularly important from a data privacy perspective when the original source data is unavailable. To address this, we evaluate the performance of CRETTA on CIFAR10-C with the replay buffer composed of CIFAR100 that we assume is distributionally similar but was not seen during the pretraining phase. As shown in Table 6, no performance drop is observed when the buffer is

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%

constructed from either CIFAR100 training or validation set. These findings suggest that CRETTA can operate effectively even without access to the original source data.

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

These comprehensive findings highlight that a small and randomly sampled buffer suffices for effective adaptation, regardless of its size or composition. This insensitivity to buffer configuration underscores the practicality of CRETTA, enabling deployment in real-world scenarios with minimal memory overhead and 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%

5 RELATED WORKS

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

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

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

6 Conclusion

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

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

REFERENCES

486

487

488

489

491

492

493

494

495 496

497

498 499

500

501

502

504 505

506

507

508

509

510 511

512

513 514

515

516

517

518

519

520

521 522

523

524

525

526 527

528

529

530 531

532

533

534

535

537

538

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

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

Omer Yair and Tomer Michaeli. Contrastive divergence learning is a time reversal adversarial game, 2021. URL https://arxiv.org/abs/2012.03295. Xiulong Yang and Shihao Ji. Jem++: Improved techniques for training jem. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6494–6503, 2021. Longhui Yuan, Binhui Xie, and Shuang Li. Robust test-time adaptation in dynamic scenarios. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 15922-15932, 2023. Yige Yuan, Bingbing Xu, Liang Hou, Fei Sun, Huawei Shen, and Xueqi Cheng. Tea: Test-time energy adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 23901–23911, 2024. Sergey Zagoruyko. Wide residual networks. arXiv preprint arXiv:1605.07146, 2016. Hao Zhao, Yuejiang Liu, Alexandre Alahi, and Tao Lin. On pitfalls of test-time adaptation, 2023. URL https://arxiv.org/abs/2306.03536.