# TEA: Test-time Energy Adaptation
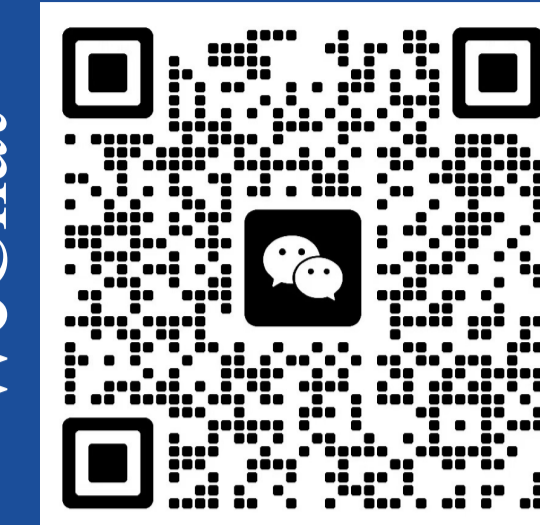
**Yige Yuan, Bingbing Xu, Liang Hou, Fei Sun, Huawei Shen, Xueqi Cheng**

CAS Key Laboratory of AI Safety, Institute of Computing Technology, Chinese Academy of Sciences
University of Chinese Academy of Sciences
Kuaishou Technology

Email: yuanyige20z@ict.ac.cn

Paper | Code | Homepage | WeChat
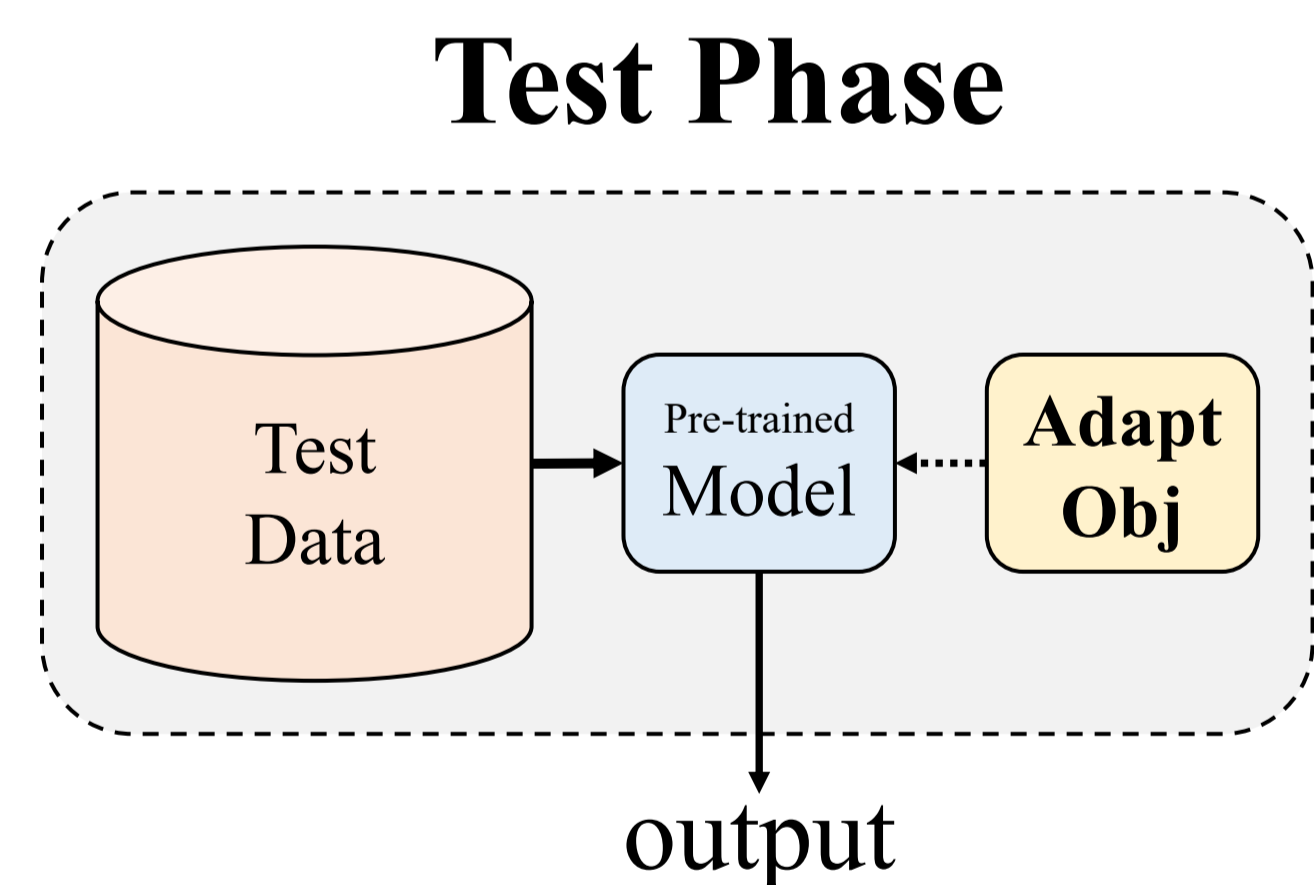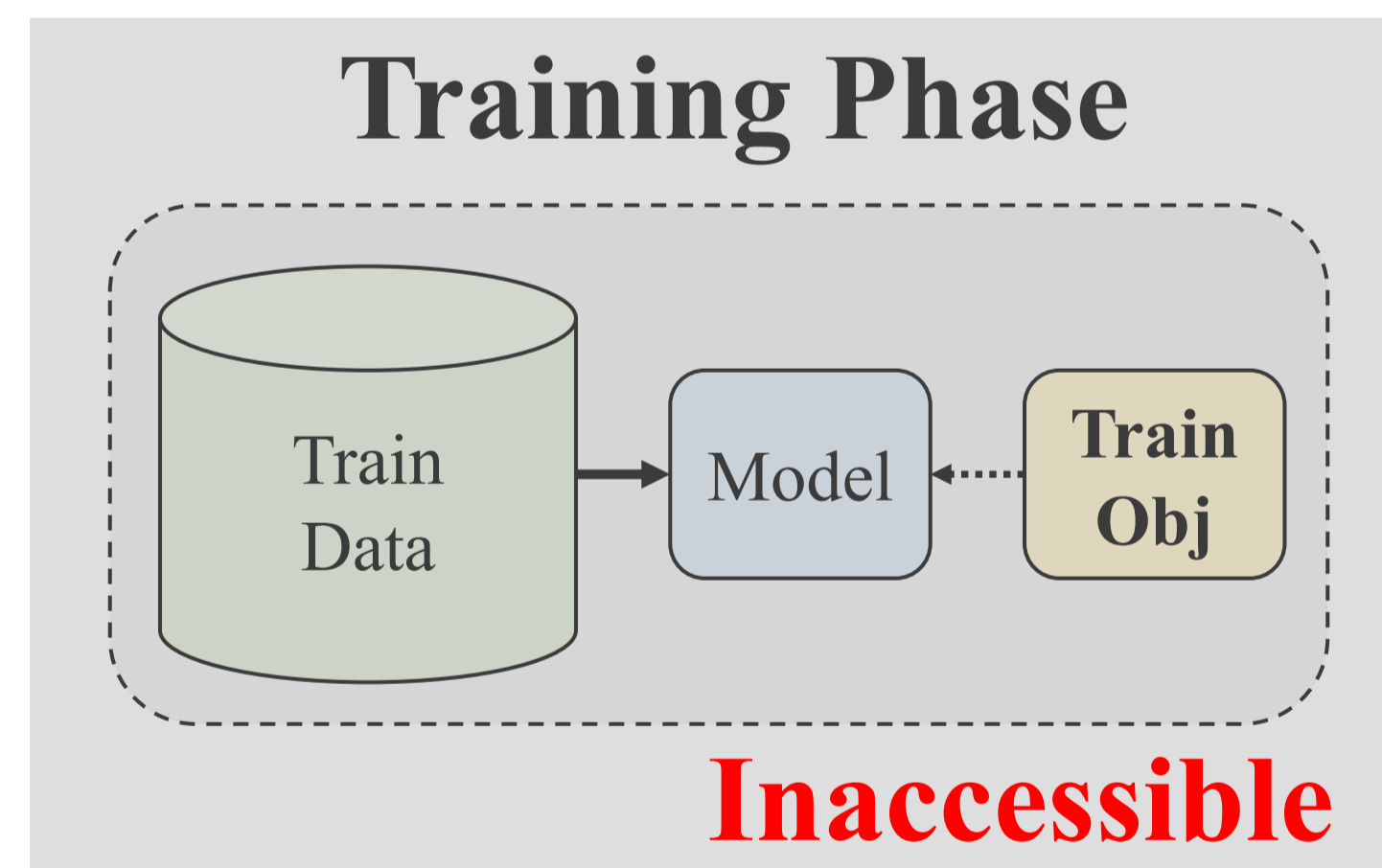
CVPR JUNE 17-21, 2024 SEATTLE, WA

## INTRODUCTION

**Objective**: Improving model **generalizability** when test data diverges from training distribution, without requiring access to training data and processes.

**Weakness of existing methods**: Current TTA methods fail to address the fundamental issue: **covariate shift**, i.e., the decreased generalizability can be attributed to the model's reliance on the marginal distribution of the training data, which may impair model calibration and introduce confirmation bias.
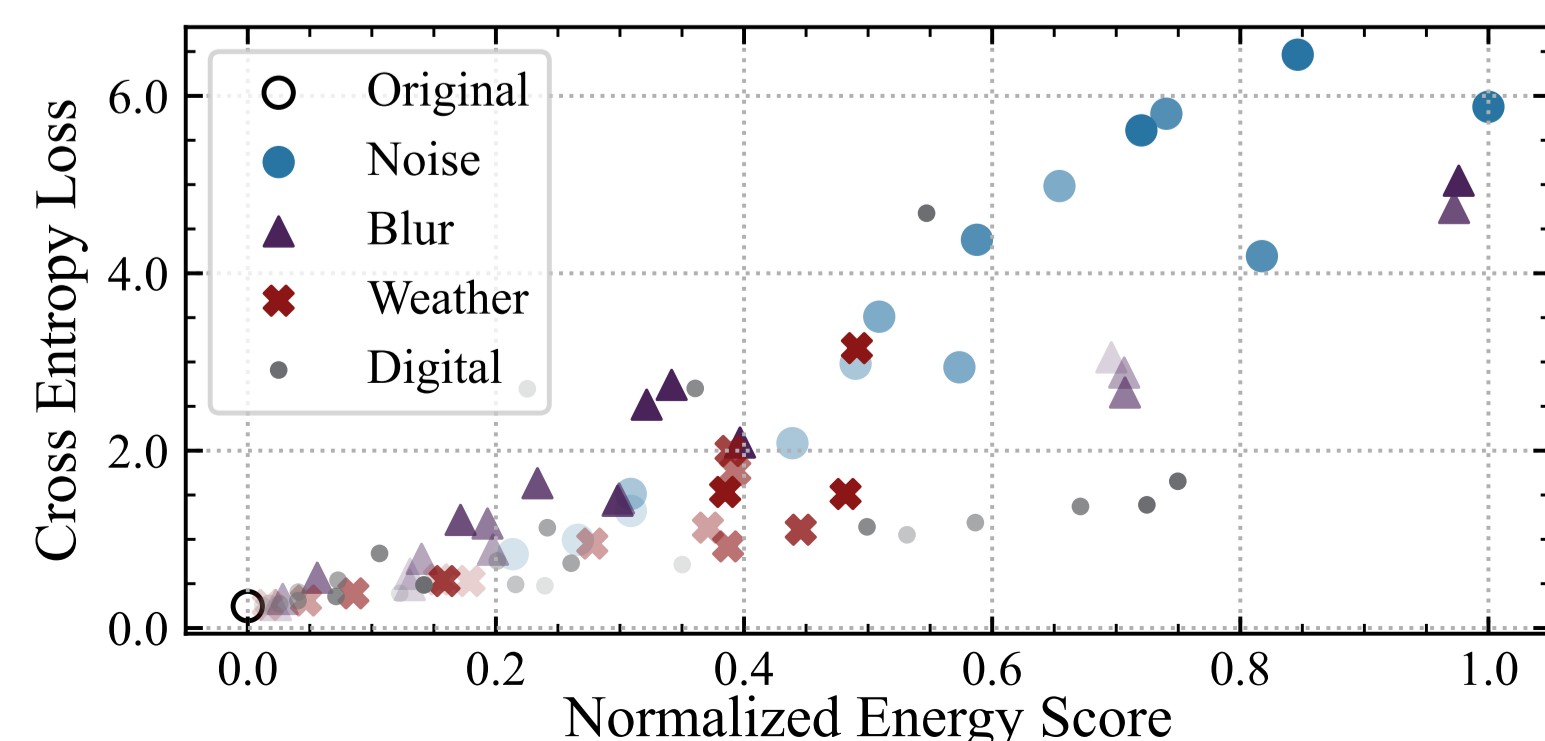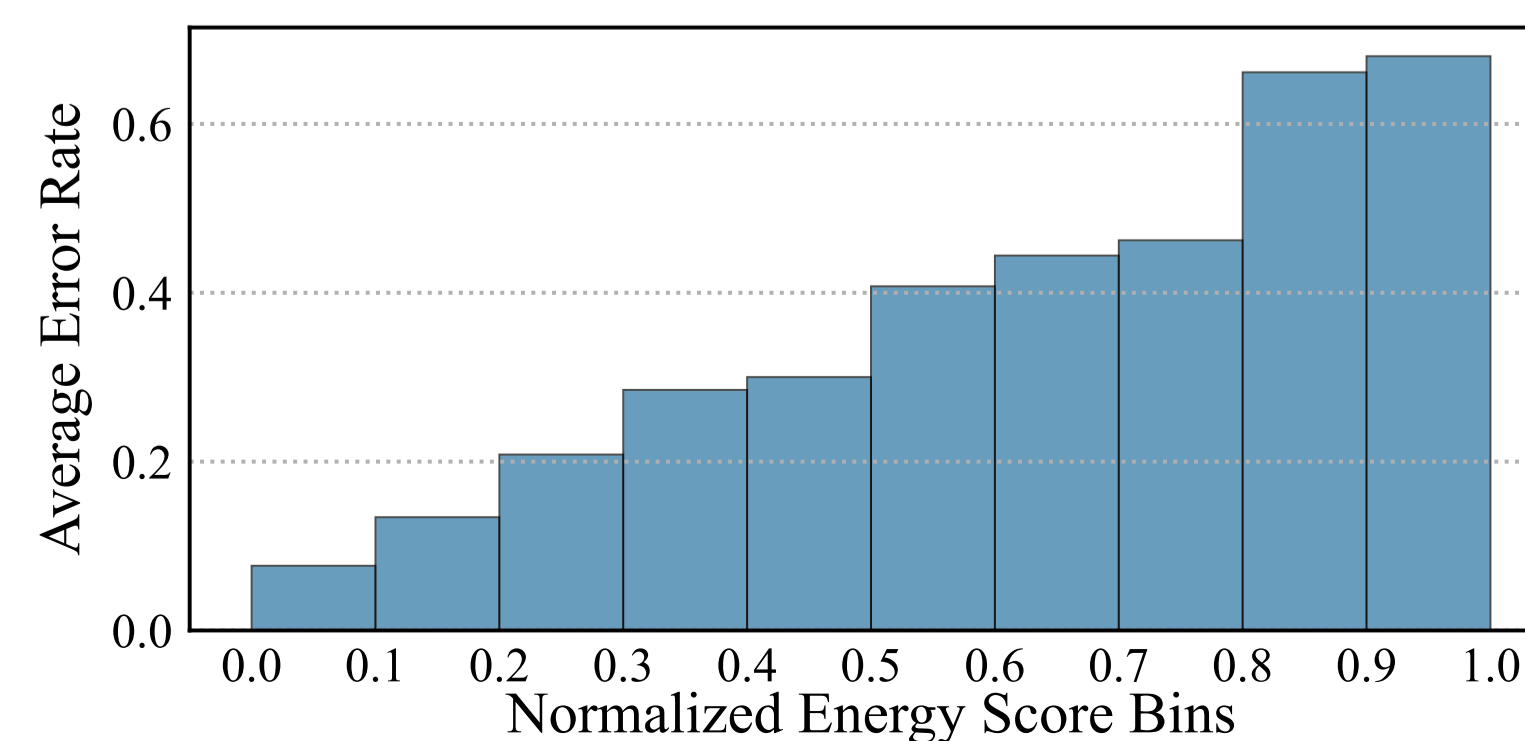
**Motivation**: Transforming the trained classifier into an **energy-based** model and aligning the model's distribution with the test data's, enhancing its ability to perceive test distributions and thus improving overall generalizability.
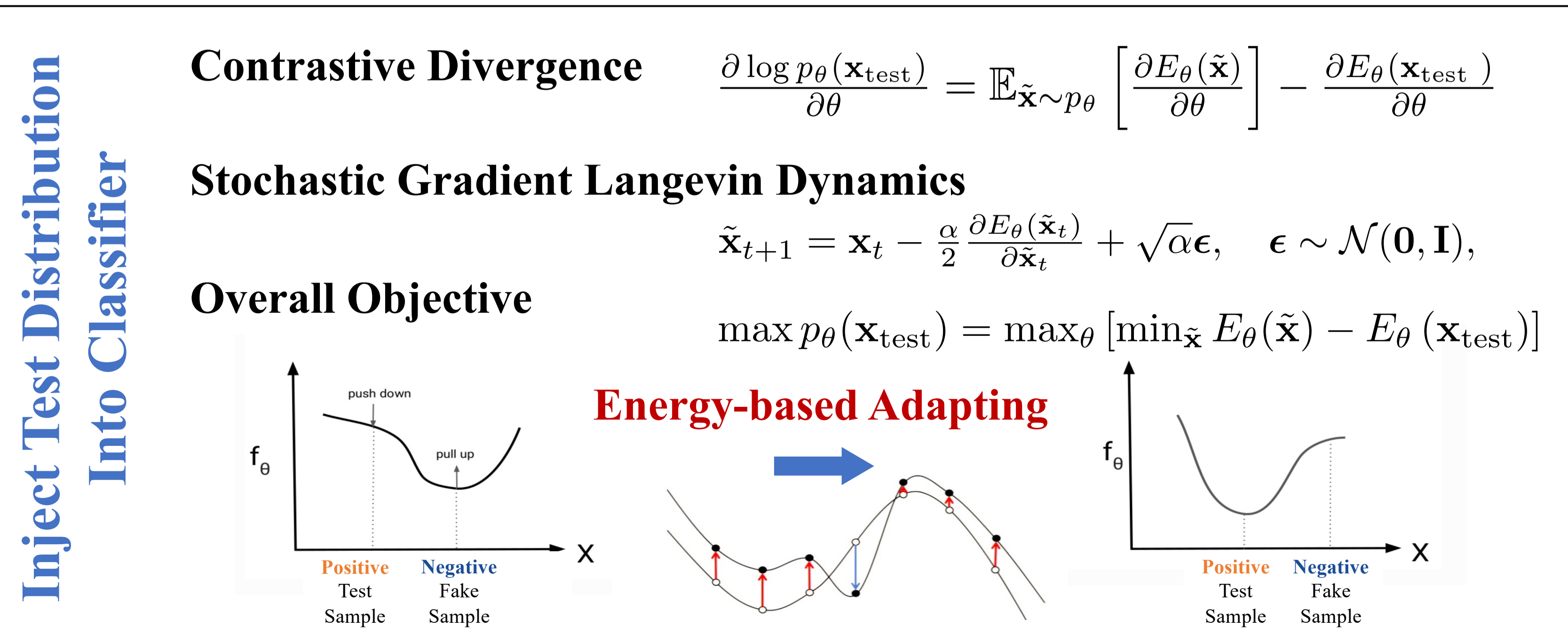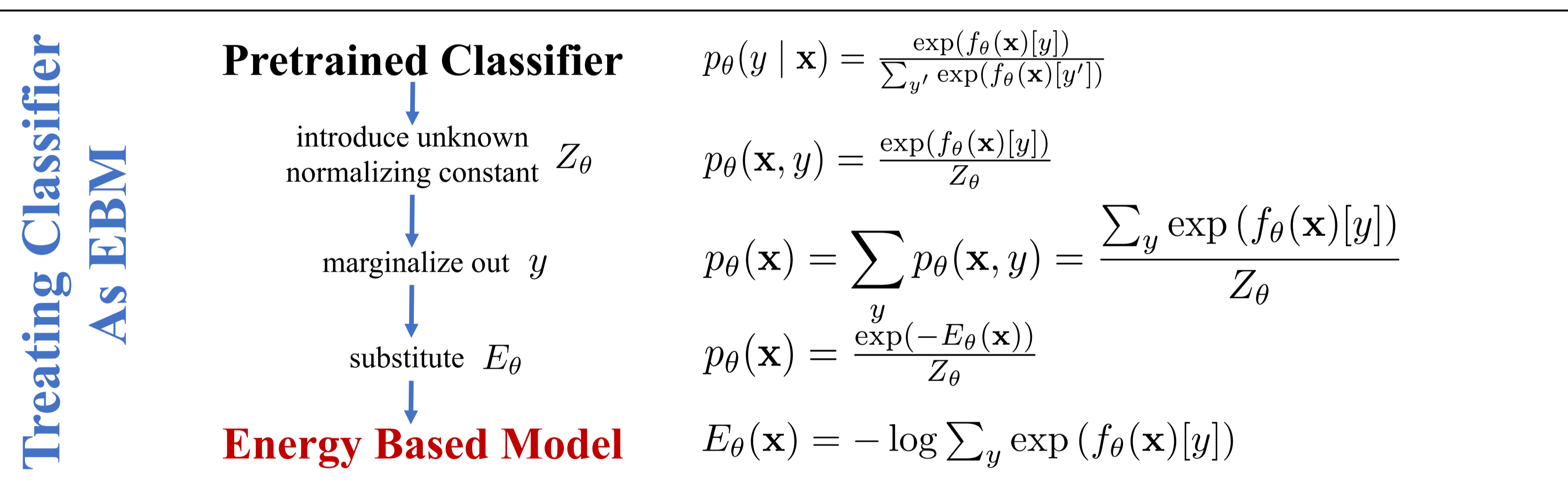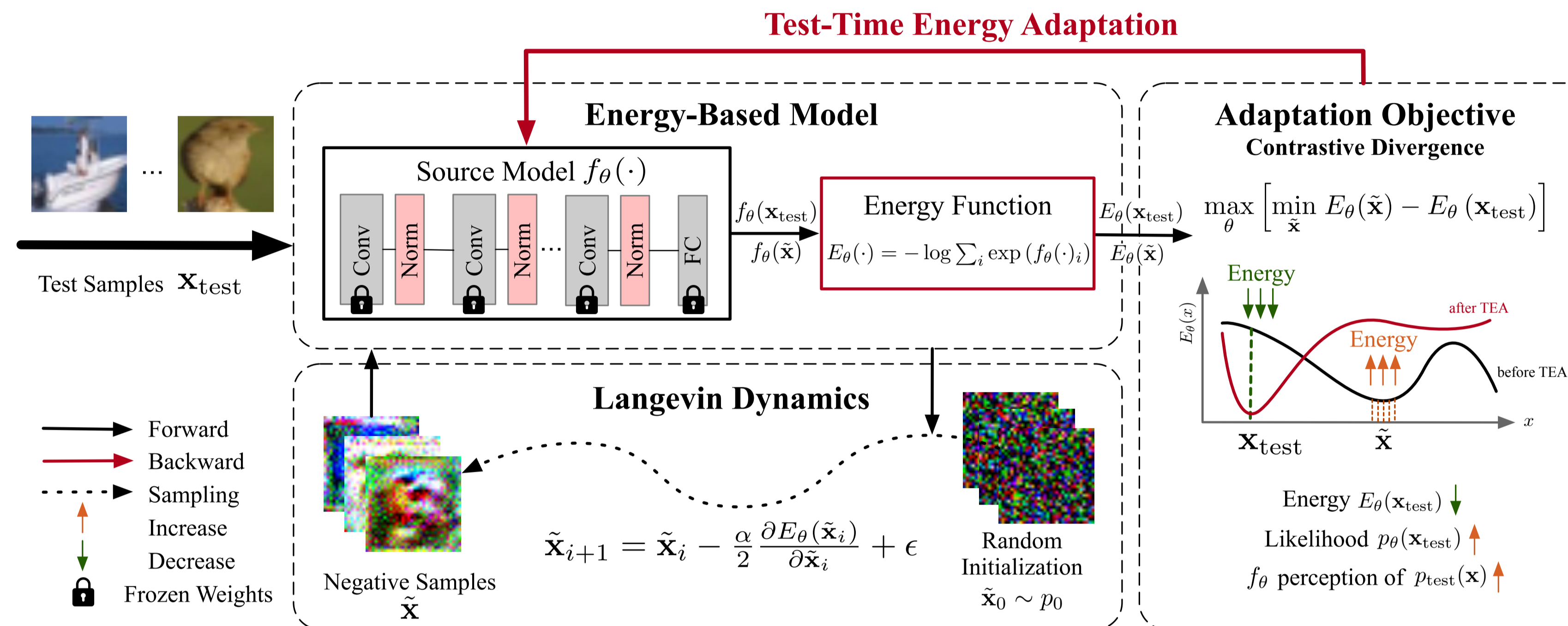
### BACKGROUND

**Training Phase** — Train Data → Model ← Train Obj

**Test Phase** — Test Data → Pre-trained Model ← Adapt Obj → output

**Inaccessible**

### MOTIVATION

**Low Energy** ⟹ **High Probability    High Performance**
**High Energy** ⟹ **Low Probability    Low Performance**



## METHOD

**Test-Time Energy Adaptation**

**Energy-Based Model**

Source Model $f_\theta(\cdot)$ — Conv, Norm, Conv, Norm, Conv, Norm, FC

$f_\theta(\mathbf{x}_{test})$, $f_\theta(\tilde{\mathbf{x}})$

Energy Function $E_\theta(\cdot) = -\log \sum_i \exp(f_\theta(\cdot)_i)$

$E_\theta(\mathbf{x}_{test})$, $\tilde{E}_\theta(\tilde{\mathbf{x}})$

**Adaptation Objective — Contrastive Divergence**

$\max_\theta \left[ \min_{\tilde{\mathbf{x}}} E_\theta(\tilde{\mathbf{x}}) - E_\theta(\mathbf{x}_{test}) \right]$

Energy $E_\theta(\mathbf{x}_{test})$ ↓
Likelihood $p_\theta(\mathbf{x}_{test})$ ↑
$f_\theta$ perception of $p_{test}(\mathbf{x})$ ↑

Test Samples $\mathbf{x}_{test}$

**Langevin Dynamics**

$\tilde{\mathbf{x}}_{i+1} = \tilde{\mathbf{x}}_i - \frac{\alpha}{2}\frac{\partial E_\theta(\tilde{\mathbf{x}}_i)}{\partial \tilde{\mathbf{x}}_i} + \epsilon$

Negative Samples $\tilde{\mathbf{x}}$ ← Random Initialization $\tilde{\mathbf{x}}_0 \sim p_0$

Forward / Backward / Sampling / Increase / Decrease / Frozen Weights

### Treating Classifier As EBM

**Pretrained Classifier**
$p_\theta(y \mid \mathbf{x}) = \frac{\exp(f_\theta(\mathbf{x})[y])}{\sum_{y'}\exp(f_\theta(\mathbf{x})[y'])}$

introduce unknown normalizing constant $Z_\theta$
$p_\theta(\mathbf{x}, y) = \frac{\exp(f_\theta(\mathbf{x})[y])}{Z_\theta}$

marginalize out $y$
$p_\theta(\mathbf{x}) = \sum_y p_\theta(\mathbf{x}, y) = \frac{\sum_y \exp(f_\theta(\mathbf{x})[y])}{Z_\theta}$

substitute $E_\theta$
$p_\theta(\mathbf{x}) = \frac{\exp(-E_\theta(\mathbf{x}))}{Z_\theta}$

**Energy Based Model**
$E_\theta(\mathbf{x}) = -\log \sum_y \exp(f_\theta(\mathbf{x})[y])$

### Inject Test Distribution Into Classifier

**Contrastive Divergence**
$\frac{\partial \log p_\theta(\mathbf{x}_{test})}{\partial \theta} = \mathbb{E}_{\tilde{\mathbf{x}} \sim p_\theta}\left[\frac{\partial E_\theta(\tilde{\mathbf{x}})}{\partial \theta}\right] - \frac{\partial E_\theta(\mathbf{x}_{test})}{\partial \theta}$

**Stochastic Gradient Langevin Dynamics**
$\tilde{\mathbf{x}}_{t+1} = \mathbf{x}_t - \frac{\alpha}{2}\frac{\partial E_\theta(\tilde{\mathbf{x}}_t)}{\partial \tilde{\mathbf{x}}_t} + \sqrt{\alpha}\epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$

**Overall Objective**
$\max p_\theta(\mathbf{x}_{test}) = \max_\theta \left[\min_{\tilde{\mathbf{x}}} E_\theta(\tilde{\mathbf{x}}) - E_\theta(\mathbf{x}_{test})\right]$

**Energy-based Adapting**

Positive Test Sample / Negative Fake Sample

## EXPERIMENTS

**TEA's Adaptation Performance**

| WRN-28-10 BatchNorm | | CIFAR-10(C) | | | | CIFAR-100(C) | | | | Tiny-ImageNet(C) | | | |
| | | Clean | Corr Severity 5 | | Corr Severity 1-5 | | Clean | Corr Severity 5 | | Corr Severity 1-5 | | Clean | Corr Severity 5 | | Corr Severity 1-5 |
| | | Acc (↑) | Acc (↑) | mCE (↓) | Acc (↑) | mCE (↓) | Acc (↑) | Acc (↑) | mCE (↓) | Acc (↑) | mCE (↓) | Acc (↑) | Acc (↑) | mCE (↓) | Acc (↑) | mCE (↓) |
| Source | | 94.77 | 56.47 | 100.00 | 73.45 | 100.00 | 35.39 | 100.00 | | | | 52.12 | 100.00 | | 63.19 | 21.21 | 19.50 | 34.13 | 100.00 |
| Norm | BN [52] | 93.97 | 79.56 | 52.65 | 85.63 | 60.00 | 80.83 | | | | | | | |
| | DUA* [41] | - | 80.10 | 50.78 | - | - | - | | | | | | | |
| Pseudo | PL [34] | 93.75 | 51.42 | 106.98 | 72.62 | 99.37 | 80.52 | | | | | | | |
| | SHOT [36] | 93.25 | 74.77 | 63.19 | 82.35 | 72.61 | 80.52 | | | | | | | |
| Entropy | TENT [60] | 93.66 | 81.41 | 48.13 | 86.75 | 56.17 | 80.14 | | | | | | | |
| | ETA [45] | 93.96 | 79.58 | 52.64 | 85.63 | 59.99 | 80.65 | | | | | | | |
| | EATA [45] | 93.96 | 79.59 | 52.62 | 85.64 | 59.98 | 80.68 | | | | | | | |
| | SAR [46] | 93.97 | 79.77 | 51.94 | 85.83 | 58.97 | 80.84 | | | | | | | |
| Energy | TEA | 94.09 | 83.34 | 43.69 | 87.88 | 52.00 | 80.88 | | | | | | | |

**TEA's Improvements in [calibrat]ion and Generalizability Enhancement**

Source: ECE(↓)=4.11% MCE(↓)=57.99%
TENT: ECE(↓)=5.50% MCE(↓)=59.73%
SHOT: ECE(↓)=6.35% MCE(↓)=58.53%
TEA: ECE(↓)=4.02% MCE(↓)=47.37%



**TEA's Distribution Perception and Generation**