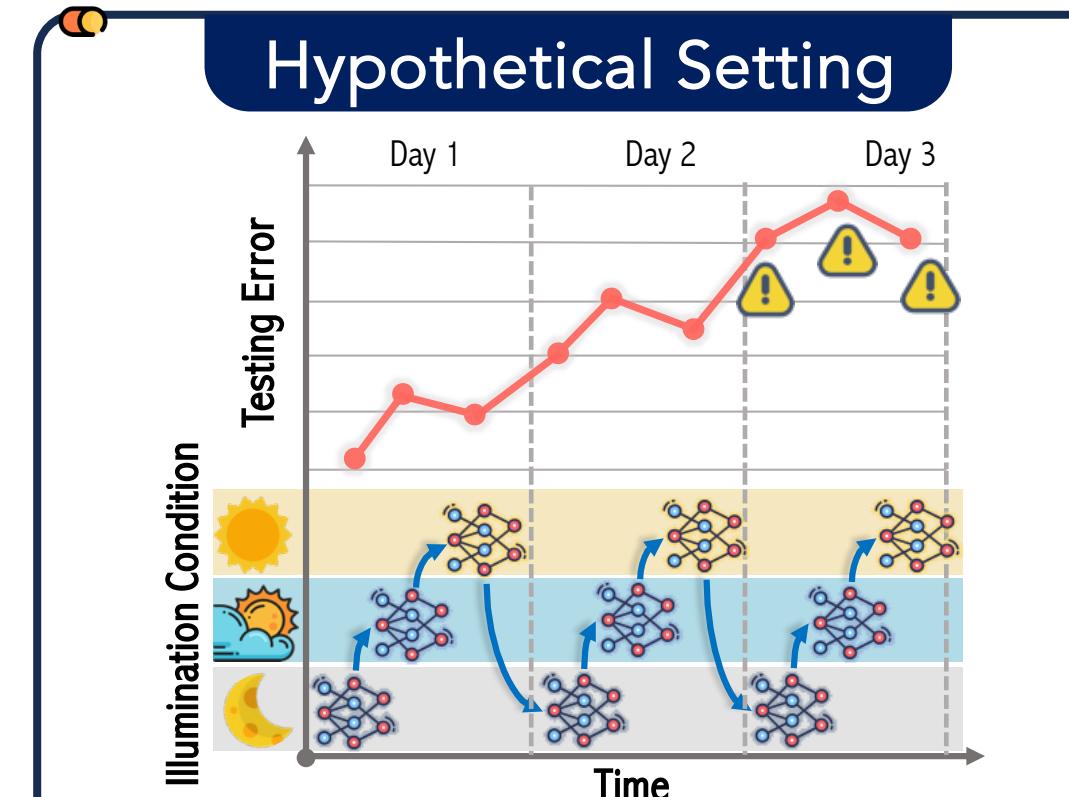


## INTRODUCTION

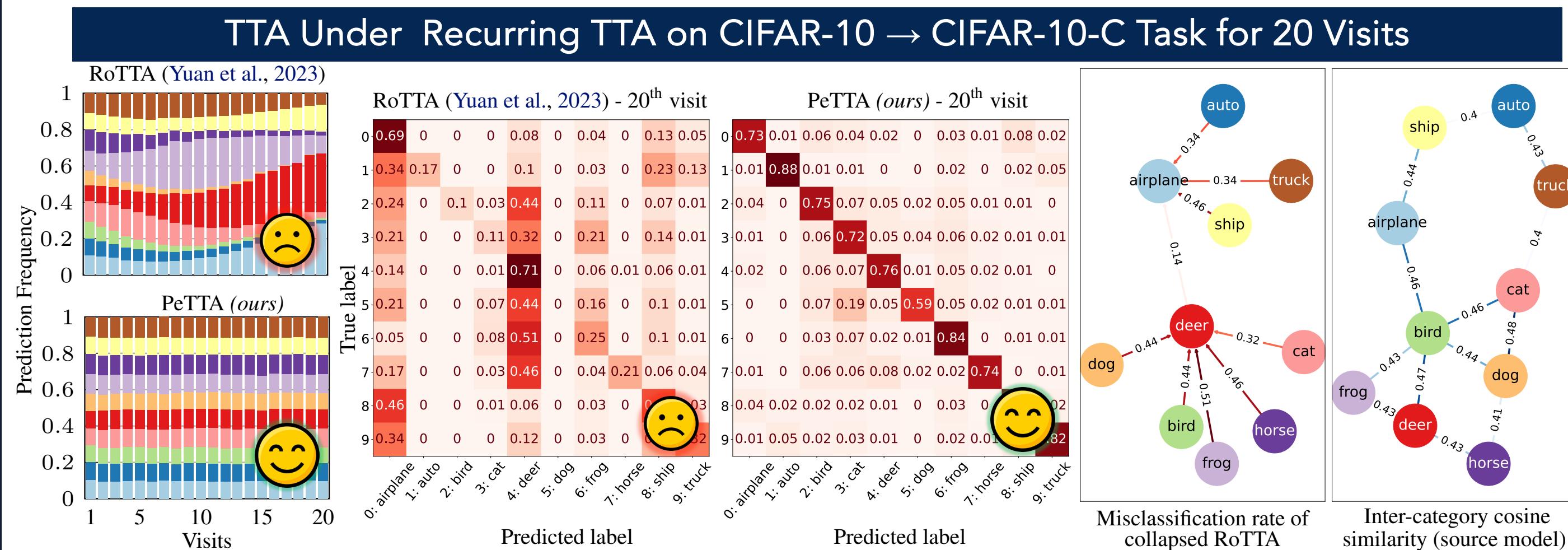
Continual Test-time Adaptation (TTA) operates on an ML classifier  $f_t: \mathcal{X} \rightarrow \mathcal{Y}$ , parameterized by  $\theta_t \in \Theta$  gradually changing over time. The model explores an online stream of testing data  $X_t \sim P_t$  for adapting itself  $f_{t-1} \rightarrow f_t$  (self-supervised learning) before predicting  $\hat{Y}_t = f_t(X_t)$ .

Does the model adaptability persist after a long time adapting to multiple data shifts?

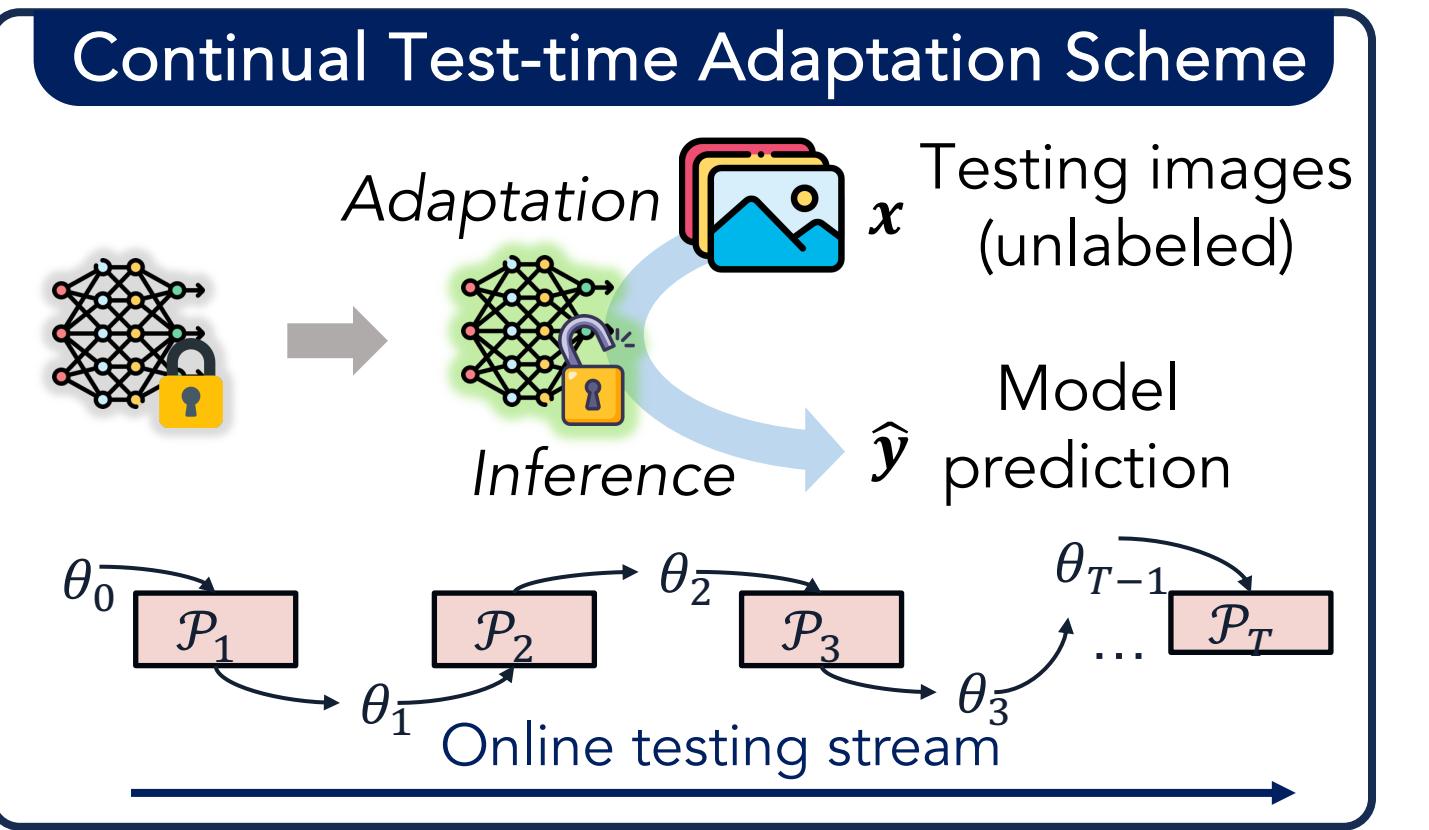


- In practice, testing environments may *change recurrently*.
- Preserving adaptability when visiting the same testing condition is *not guaranteed*.

**Recurring Test-time Scenario:**  $\mathcal{P}_1 \rightarrow \mathcal{P}_2 \rightarrow \dots \rightarrow \mathcal{P}_D \rightarrow \dots \rightarrow \mathcal{P}_1 \rightarrow \mathcal{P}_2 \rightarrow \dots \rightarrow \mathcal{P}_D$



(a) Histogram of model predictions: PeTTA achieves a persisting performance while RoTTA degrades. (b) Confusion matrix at the last visit (c) Force-directed graph showing (left) the most prone to misclassification; (right) similar categories tend to be easily collapsed.



## ε-PERTURBED GAUSSIAN MIXTURE MODEL CLASSIFIER (ε-GMMC)

ε-GMMC - a simple yet representative failure case of TTA for theoretical analysis

**Setting:** A simplified continual TTA process

- Let  $p_{y,t} = \Pr(Y_t = y)$ ;  $\hat{p}_{y,t} = \Pr(\hat{Y}_t = y)$ .
- Binary classification  $\mathcal{X} \times \mathcal{Y} = \mathbb{R} \times \{0,1\}$ .
- Underlying distribution follows a mixture of 2 Gaussian:  $P_t(x, y) = p_{y,t} \mathcal{N}(x; \mu_y, \sigma_y^2)$ .

**Main Task:** predicting  $X_t$  was sampled from cluster 0 or 1 (negative or positive).

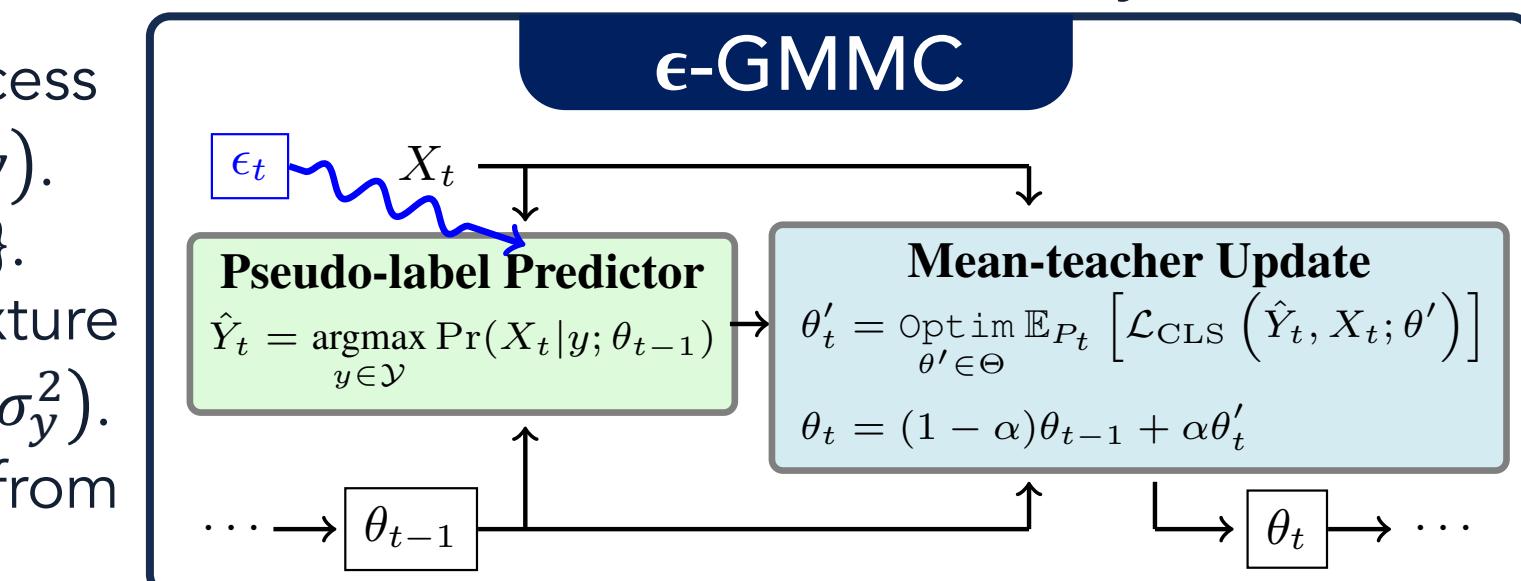
## A Mathematical Definition of Model Collapse

**Definition 1 (Model Collapse).** A model is said to be collapsed from step  $\tau \in \mathcal{T}$ ,  $\tau < \infty$  if there exists a non-empty subset of categories  $\tilde{\mathcal{Y}} \subseteq \mathcal{Y}$  such that  $\Pr\{\hat{Y}_t \in \tilde{\mathcal{Y}}\} > 0$  but the marginal  $\Pr\{\hat{Y}_t \in \tilde{\mathcal{Y}}\}$  converges to zero in probability:

$$\lim_{t \rightarrow \tau} \Pr\{\hat{Y}_t \in \tilde{\mathcal{Y}}\} = 0.$$

## Factors contributing to the model collapse:

- Data-dependent factors:** the prior data distribution ( $p_0$ ), the nature difference between two categories ( $|\mu_0 - \mu_1|$ ) from the dataset.
- Algorithm-dependent factors:** update rate ( $\alpha$ ), the false negative rate at each step ( $\varepsilon_t$ ).



ε-GMMC performs 2 main steps:

- Predicting pseudo-labels ( $\hat{Y}_t$ ).
- Updating with mean teacher model.

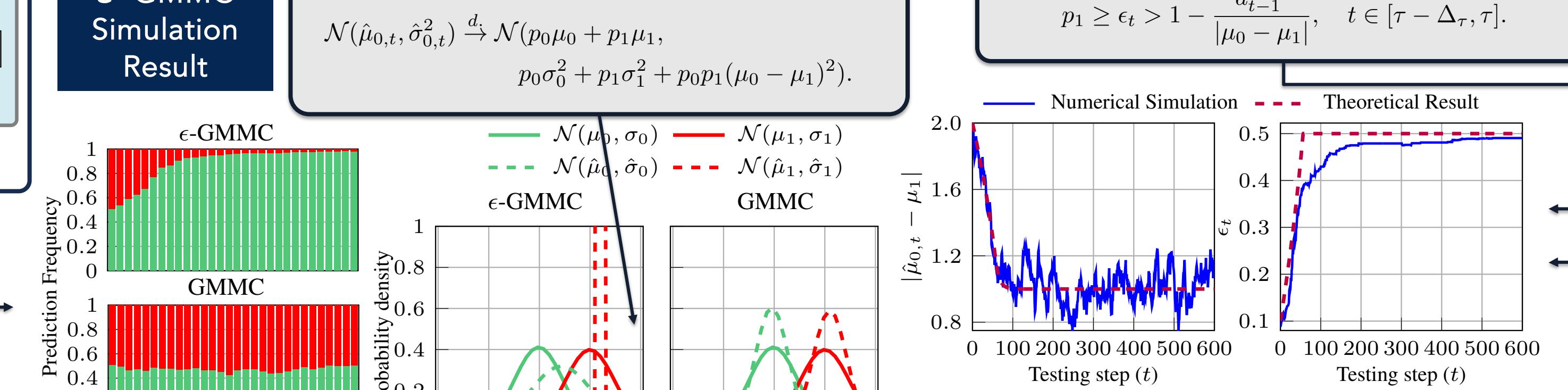
**Key Idea:** The predictor is perturbed for retaining a **false negative rate (FNR)** of  $\varepsilon_t = \Pr\{\hat{Y}_t = 1 | \hat{Y}_t = 0\}$  to simulate undesirable effects of the testing stream in TTA, making model prone to collapse.

**Assumption 1 (Static Data Stream).** The marginal distribution of the true label follows the same Bernoulli distribution  $\text{Ber}(p_0)$ :  $p_{0,t} = p_0$ ,  $(p_{1,t} = 1 - p_0)$ ,  $\forall t \in \mathcal{T}$ .

**Lemma 2 (ε-GMMC After Collapsing).** For a binary ε-GMMC model, with Assumption 1, if  $\lim_{t \rightarrow \tau} \hat{p}_{1,t} = 0$  (collapsing), the cluster 0 in GMMC converges in distribution to a single-cluster GMMC with parameters:

$$\mathcal{N}(\hat{\mu}_{0,t}, \hat{\sigma}_{0,t}^2) \xrightarrow{d} \mathcal{N}(p_0 \mu_0 + p_1 \mu_1, p_0 p_0^2 + p_1 p_1^2 + p_0 p_1 (\mu_0 - \mu_1)^2).$$

**Corollary 1 (A Condition for ε-GMMC Collapse).** With fixed  $p_0, \alpha, \mu_0, \mu_1$ , ε-GMMC is collapsed if there exists a sequence of  $\{\varepsilon_t\}_{t=\tau}^{\tau + \Delta_\tau}$  ( $\tau \geq \Delta_\tau > 0$ ) such that:

$$p_1 \geq \varepsilon_t > 1 - \frac{d_{t-1}^{0 \rightarrow 1}}{|\mu_0 - \mu_1|}, \quad t \in [\tau - \Delta_\tau, \tau].$$


**Lemma 1 (Increasing FNR).** Under Assumption 1, a binary ε-GMMC would collapsed (Def. 1) with  $\lim_{t \rightarrow \tau} \hat{p}_{1,t} = 0$  (or  $\lim_{t \rightarrow \tau} \hat{p}_{0,t} = 1$ , equivalently) if and only if  $\lim_{t \rightarrow \tau} \varepsilon_t = p_1$ .

**Theorem 1 (Convergence of ε-GMMC).** For a binary ε-GMMC model, with Assumption 1, let the distance from  $\hat{\mu}_{0,t}$  toward  $\mu_1$  is  $d_t^{0 \rightarrow 1} = |\mathbb{E}_{P_t}[\hat{\mu}_{0,t}] - \mu_1|$ , then:

$$d_t^{0 \rightarrow 1} - d_{t-1}^{0 \rightarrow 1} \leq \alpha \cdot p_0 \cdot \left( |\mu_0 - \mu_1| - \frac{d_{t-1}^{0 \rightarrow 1}}{1 - \varepsilon_t} \right).$$

(a) Histogram of model predictions. (b) The probability density function of the two clusters after convergence (dashed line) versus the true data distribution. (c) Distance toward  $\mu_1$  and false-negative rate ( $\varepsilon_t$ ) coincides with the theoretical analysis.

## PERSISTENT TEST-TIME ADAPTATION (PeTTA)

**Key Idea:** Striking a balance between **adaptation** and **preventing model collapse**

With  $\phi_{\theta_t}$  is the deep feature extractor of  $f_t$ , let  $\mathbf{z} = \phi_{\theta_t}(\mathbf{x})$ . Keeping track of a collection of the running mean of feature vector  $\mathbf{z}$ :  $\{\hat{\mu}_t^y\}_{y \in \mathcal{Y}}$  in which  $\hat{\mu}_t^y$  is exponential moving average updated with vector  $\mathbf{z}$  if  $f_t(\mathbf{x}) = y$ .

## Persistent TTA

(2) Adaptive Learning Rate  $\alpha_t$  and Regularization  $\lambda_t$ (1) Sensing the divergence from  $\theta_0$ 

$$\hat{\gamma}_t^y = \frac{1}{|\hat{\mathcal{Y}}_t|} \sum_{y \in \hat{\mathcal{Y}}_t} \gamma_t^y, \quad \hat{\mathcal{Y}}_t = \{\hat{Y}_t^{(i)} | i = 1, \dots, N_t\}$$

$\mu_t^y, \Sigma_t^y$  are pre-computed on the source distribution

$$\lambda_t = \hat{\gamma}_t \cdot \lambda_0, \quad \alpha_t = (1 - \hat{\gamma}_t) \cdot \alpha_0,$$

## (3) Anchor Loss

$$\theta'_t = \text{Optim} \mathbb{E}_{\theta' \in \Theta} [\mathcal{L}_{\text{CLS}}(\hat{Y}_t, X_t; \theta') + \mathcal{L}_{\text{AL}}(X_t; \theta')] + \lambda_t \mathcal{R}(\theta')$$

$$\theta_t = (1 - \alpha_t) \theta_{t-1} + \alpha_t \theta'_t.$$

$$\mathcal{L}_{\text{AL}}(X_t; \theta) = - \sum_{y \in \mathcal{Y}} \Pr(y|X_t; \theta_0) \log \Pr(y|X_t; \theta)$$

## EXPERIMENTAL RESULTS

Average classification error on the task ImageNet → ImageNet-C for 20 recurring TTA visits.

Method	Recurring TTA visit																				Avg
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	
LAME (Boudiaf et al., 2022)	82.0																				82.0
CoTTA (Wang et al., 2022)	98.6	99.1	99.4	99.4	99.5	99.5	99.5	99.5	99.6	99.6	99.6	99.6	99.6	99.6	99.6	99.6	99.6	99.6	99.7	99.5	
RoTTA (Niu et al., 2022)	60.4	59.3	65.4	72.6	79.1	84.2	88.7	92.7	95.2	96.9	97.7	98.1	98.4	98.6	98.7	98.8	98.9	98.9	99.0	89.0	
RMT (Döbler et al., 2022)	72.3	71.0	69.9	69.1	68.8	68.5	68.4	68.3	70.0	70.2	70.1	70.2	72.8	76.8	75.6	75.1	75.1	75.2	74.8	74.7	71.8
MECTA (Hong et al., 2023)	77.2	82.8	86.1	87.9	88.9	89.8	89.8	89.9	90.0	90.0	90.6	90.7	90.7	90.9	90.8	90.8	90.8	90.8	90.7	90.8	89.0
RoTTA (Yuan et al., 2023)	68.3	62.1	61.8	64.5	68.4	75.4	82.7	95.1	95.8	96.4	97.1	97.9	98.3	98.7	99.0	99.1	99.3	99.4	99.5	99.6	87.9
RDumb (Press et al., 2023)	72.2	73.0	73.2	72.8	72.8	73.2	73.7	72.3	73.1	73.2	73.2	73.1	72.1	72.6	73.3	73.1	72.8	73.2	73.3	72.8	
ROID (Marsden et al., 2024)	62.7	62.3	62.3	62.3	62.5	62.3	62.4	62.4	62.5	62.6	62.4	62.4	62.5	62.5	62.4	62.4	62.5	62.5	62.4	62.4	
TRIBE (Su et al., 2024)	63.6	64.0	64.9	67.8	69.6	71.7	73.5	75.5	77.4	79.8	85.0	96.5	99.4	99.9	99.8	99.9	99.9	99.9	99.9	99.9	
PeTTA (ours) <sup>(*)</sup>	65.3	61.7	59.8	59.1	59.4	59.6	59.8	59.3	59.4	60.0	60.3	61.0	60.7	60.4	60.6	60.7	60.8	60.7	60.4	60.2	

Does model reset help? A comparison with a reset-based approach at different frequencies.

Reset Every	Recurring TTA visit																				Avg
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20		




<tbl\_r cells="22