

# Resurrecting Old Classes with New Data for Exemplar-Free Continual Learning

Dipam Goswami<sup>1,2</sup> Albin Soutif-Cormerais<sup>1,2</sup> Yuyang Liu<sup>3,4,†</sup> Sandesh Kamath<sup>1,2</sup>  
Bartłomiej Twardowski<sup>1,2,5</sup> Joost van de Weijer<sup>1,2</sup>



# Resurrecting Old Classes with New Data for Exemplar-Free Continual Learning

Dipam Goswami<sup>1,2</sup> Albin Soutif-Cormerais<sup>1,2</sup> Yuyang Liu<sup>3,4,†</sup> Sandesh Kamath<sup>1,2</sup>  
Bartłomiej Twardowski<sup>1,2,5</sup> Joost van de Weijer<sup>1,2</sup>

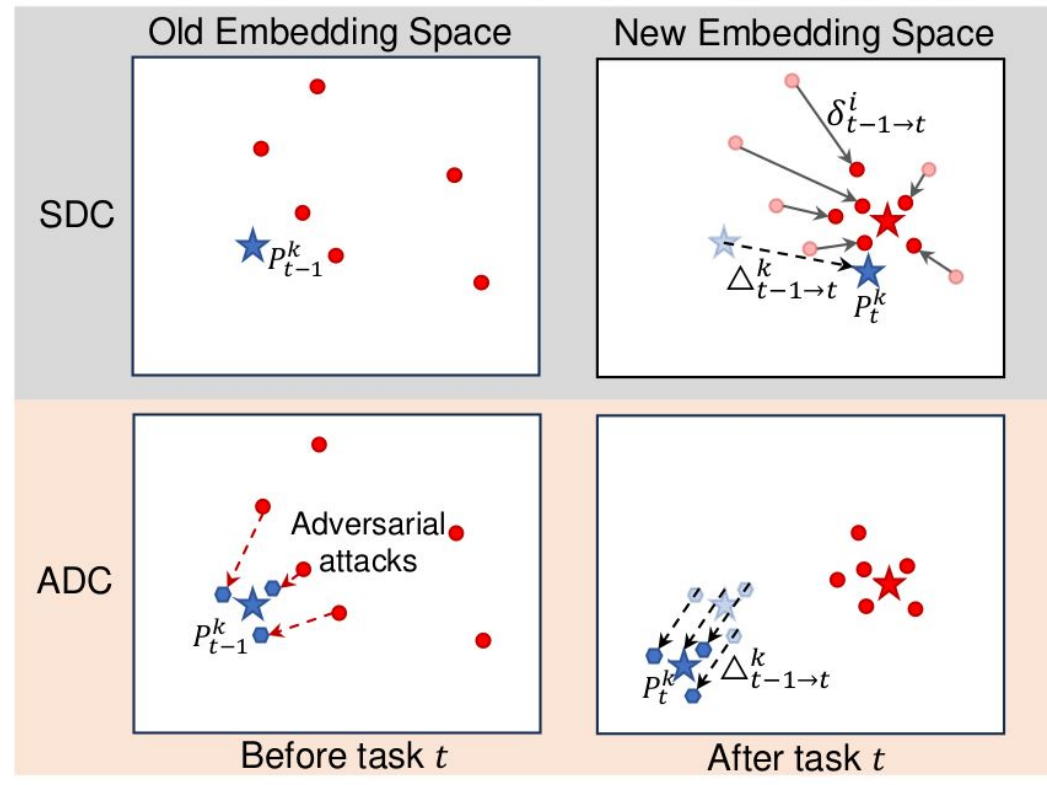
- When the feature extractor is continually trained on new tasks, old class prototypes drift in the embedding space. In the absence of data from old classes, we cannot estimate the new position of old prototypes in the embedding space.
- We address the problem of feature drift compensation for exemplar-free methods. We propose to generate adversarial samples using the new task samples and then use these adversarial samples to track the movement of old class prototypes in the feature space.

# Introduction

- Why exemplar-free continual Learning ?
  - Multiple challenges in practical settings with exemplar-based methods such as legal concerns with new regulations (e.g. European GDPR where users can request to delete personal data), and privacy issues when dealing with sensitive data like in medical imaging
- Small-start vs Big-start settings -
  - To reduce potential drift in the feature extractor, existing exemplar-free methods are typically evaluated in settings where the first task is significantly larger than subsequent tasks (half of the dataset in first task) and the model learns high-quality feature representations in the first task.
  - Usually, this drift is minimized with heavy functional regularization, which consequently restricts the plasticity of the network.
  - Small-start settings (equally splitting the dataset into tasks) starts from a smaller initial task and thus is a more challenging setting.

# Introduction

- New Samples
- Adversarial Samples
- ☆ Prototype → Semantic Drift



# Introduction

## Traditional Buffer (Exemplars)



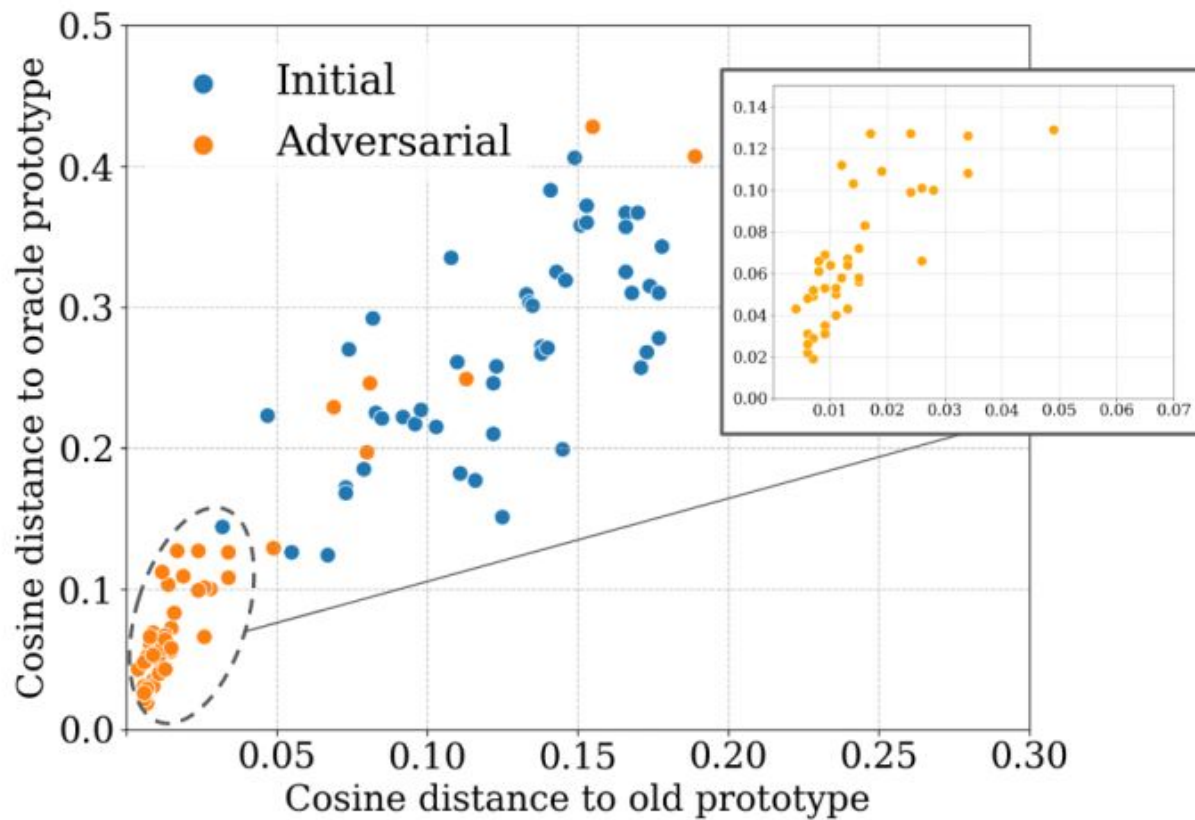
## Current Data



## Adversarial Samples (Pseudo-exemplars)



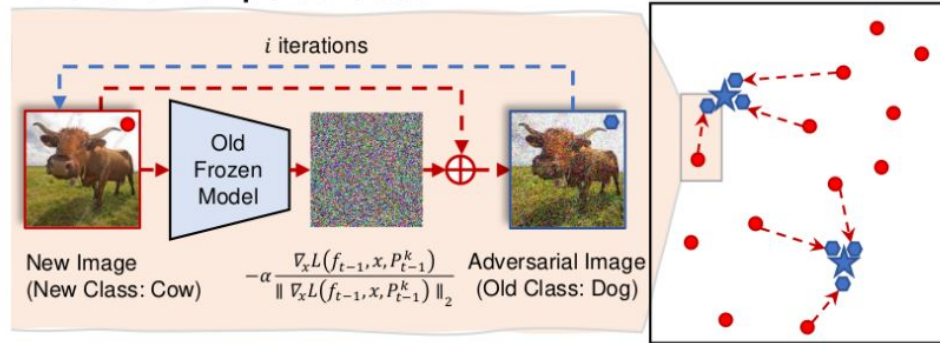
# Motivation





# Adversarial Drift Compensation

## Adversarial Sample Generation



## Adversarial Drift Compensation

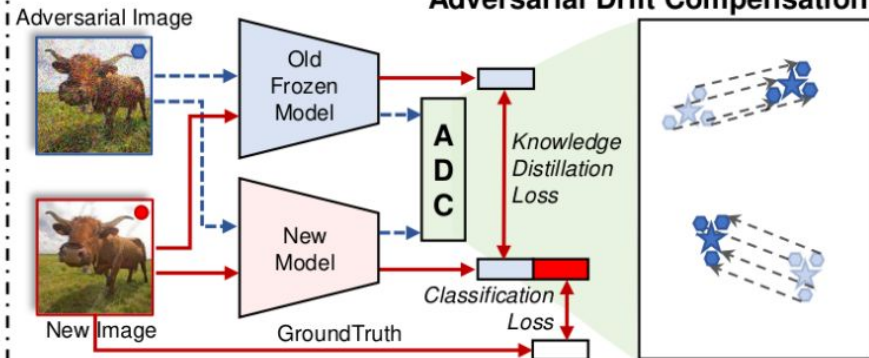


Figure 3. (a) Adversarial Sample Generation: On the old model feature space, the new samples closest to the old prototype are selected and iteratively perturbed in the direction of the target old prototype to generate adversarial samples which are now misclassified as the target old class resulting in embeddings closer to the old prototype. We perform this for every old class (we show 2 classes here for demonstration). (b) Model Training with Drift Compensation: The new model is trained using the classification loss for learning new classes and knowledge distillation loss to prevent forgetting of old classes. After the new model is trained, the adversarial samples generated using the old model are passed through both the models and the drift from old to new feature space is estimated. This is then used to update the old prototypes.

## Adversarial Sample Generation

We propose the following optimization objective by computing the mean squared error between the features  $f_{t-1}(x)$  and the prototype  $P_{t-1}^k$  as:

$$L(f_{t-1}, \mathcal{X}^k, P_{t-1}^k) = \frac{1}{|\mathcal{X}^k|} \sum_{x \in \mathcal{X}^k} \|f_{t-1}(x) - P_{t-1}^k\|_2^2. \quad (2)$$

In order to move the feature embedding in the direction of the target prototype  $P_{t-1}^k$ , we obtain the gradient of the loss with respect to the data  $x \in \mathcal{X}^k$ , normalize it to get the unit attack vector and scale it by  $\alpha$  as follows:

$$x_{adv} \leftarrow x - \alpha \frac{\nabla_x L(f_{t-1}, x, P_{t-1}^k)}{\|\nabla_x L(f_{t-1}, x, P_{t-1}^k)\|_2} \quad \forall x \in \mathcal{X}^k \quad (3)$$

where  $\nabla_x L(f_{t-1}, x, P_{t-1}^k)$  is the gradient of the objective function with respect to the data  $x$  and  $\alpha$  refers to the step size. We perform the attack for  $i$  iterations.

## Drift Compensation

The adversarial samples when passed through the new feature extractor  $f_t$  are expected to lie close to the drifted prototype and hence are used to compute the drift. After generating the adversarial samples for each target class  $k$ , we measure the prototype drift as:

$$\Delta_{t-1 \rightarrow t}^k = \frac{1}{|\mathcal{X}_{adv}^k|} \sum_{x_{adv} \in \mathcal{X}_{adv}^k} (f_t(x_{adv}) - f_{t-1}(x_{adv})) \quad (4)$$

where  $x_{adv} \in \mathcal{X}_{adv}^k$  is the set of only those adversarial samples which are classified as the target class  $k$  using the NCM classifier. We resurrect the old prototypes by compensating the drift as follows:

$$P_t^k = P_{t-1}^k + \Delta_{t-1 \rightarrow t}^k \quad (5)$$



---

**Algorithm 1** Adversarial Drift Compensation

---

**input:** Images  $X_t$ , feature extractors  $f_t$  and  $f_{t-1}$ , old prototypes  $P_{t-1}^{Y_{1:t-1}}$ , step size  $\alpha$ , number of generated samples  $m$ , number of iterations for perturbation  $i$ .

**output:** Resurrected prototypes  $P_t^{Y_{1:t-1}}$ .

- 1: **for**  $k \in Y_{1:t-1}$  **do**
  - 2:     Sample a set  $\mathcal{X}^k$  of  $m$  new samples from  $X_t$  which are closest to  $P_{t-1}^k$  based on L2 distance.
  - 3:     **for**  $i$  iterations **do**
  - 4:          $L(f_{t-1}, \mathcal{X}^k, P_{t-1}^k) = \frac{1}{|\mathcal{X}^k|} \sum_{x \in \mathcal{X}^k} \|f_{t-1}(x) - P_{t-1}^k\|_2^2$ .
  - 5:          $x_{adv} = x - \alpha \frac{\nabla_x L(f_{t-1}, x, P_{t-1}^k)}{\|\nabla_x L(f_{t-1}, x, P_{t-1}^k)\|_2} \forall x \in \mathcal{X}^k$
  - 6:          $x \leftarrow x_{adv}$
  - 7:     **end for**
  - 8:     add  $x$  in  $\mathcal{X}_{adv}^k$  if  $\operatorname{argmin}_{y \in Y_{1:t-1}} \|f_{t-1}(x) - P_{t-1}^y\|_2 = k$ .  
        $\triangleright$  Store only those samples in  $\mathcal{X}_{adv}^k$  which are successfully misclassified as  $k$
  - 9:      $\Delta_{t-1 \rightarrow t}^k = \frac{1}{|\mathcal{X}_{adv}^k|} \sum_{x_{adv} \in \mathcal{X}_{adv}^k} f_t(x_{adv}) - f_{t-1}(x_{adv})$
  - 10:      $P_t^k = P_{t-1}^k + \Delta_{t-1 \rightarrow t}^k$       $\triangleright$  Prototype resurrection
  - 11: **end for**
-

# Experiments

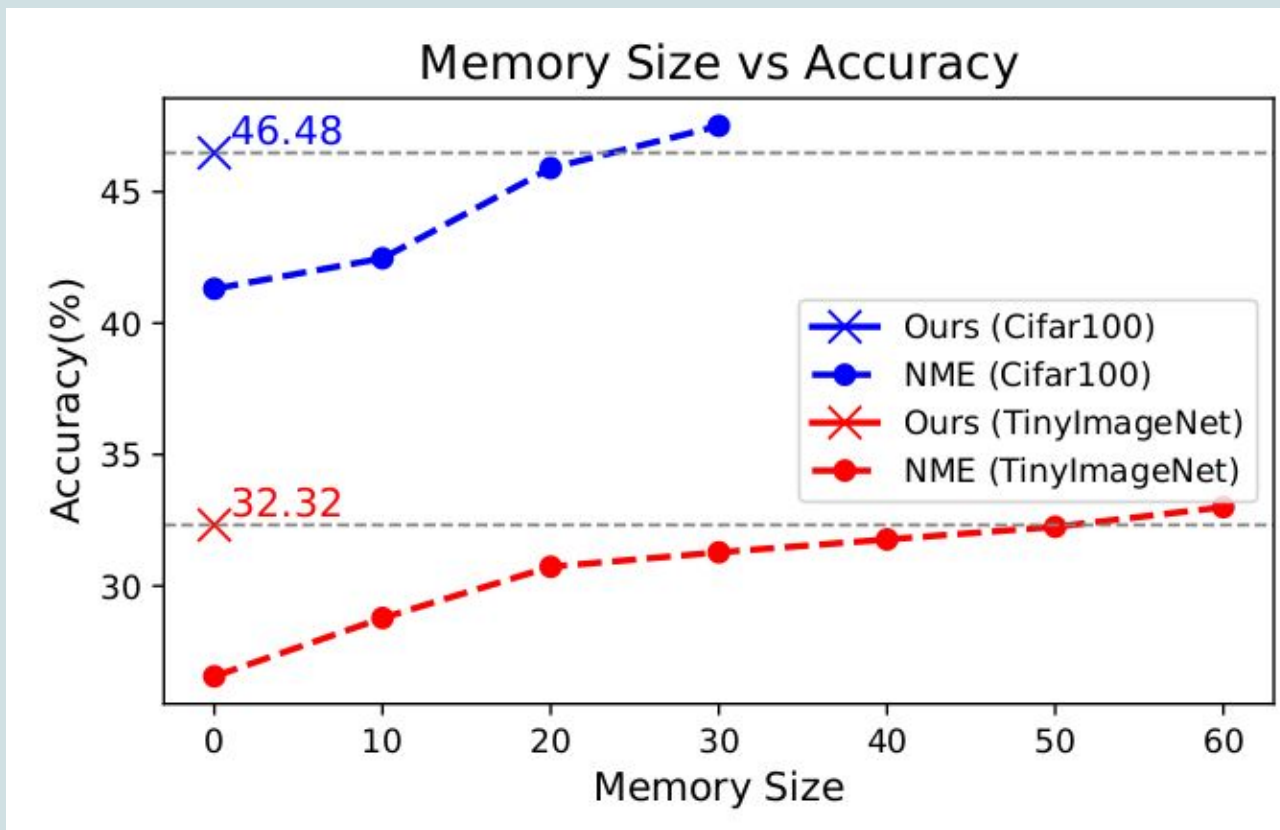
Method	CIFAR-100				TinyImageNet				ImageNet-Subset			
	T = 5		T=10		T = 5		T =10		T = 5		T = 10	
	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$
LwF [24]	45.35	61.94	26.14	46.14	38.81	49.70	<u>27.42</u>	38.77	50.88	69.11	37.90	61.60
NCM	53.53	<u>66.35</u>	41.31	57.85	38.69	50.45	26.56	<u>41.04</u>	57.74	71.99	<u>45.86</u>	65.04
SDC [52]	<u>54.94</u>	64.82	<u>41.36</u>	<u>58.02</u>	<u>40.05</u>	<u>50.82</u>	27.15	40.46	<u>59.82</u>	<u>74.10</u>	43.72	<u>65.83</u>
PASS [55]	<u>49.75</u>	63.39	<u>37.78</u>	52.18	36.44	48.64	26.58	38.65	50.96	66.15	38.90	<u>54.74</u>
SSRE [56]	42.39	56.57	29.44	44.38	30.13	43.20	22.48	34.93	40.30	57.57	28.12	45.87
FeTrIL [35]	45.11	60.42	36.69	52.11	29.91	43.99	23.88	36.35	49.18	63.83	40.26	55.12
FeCAM [11]	47.28	61.37	33.82	48.58	25.62	39.85	23.21	35.32	54.18	67.21	42.68	57.45
ADC (Ours)	<b>59.14</b>	<b>69.62</b>	<b>46.48</b>	<b>61.35</b>	<b>41.0</b>	<b>50.94</b>	<b>32.32</b>	<b>43.04</b>	<b>62.40</b>	<b>74.84</b>	<b>47.58</b>	<b>67.07</b>

Table 1. Evaluation of EFCIL methods on small-start settings. Best results in **bold** and second best results are underlined.

# Experiments

Method	CUB-200				Stanford Cars			
	T = 5		T=10		T = 7		T =14	
	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$	$A_{last}$	$A_{inc}$
LwF [24]	<u>58.68</u>	<u>71.31</u>	41.96	60.15	<u>45.18</u>	61.14	30.33	49.93
NCM	52.74	67.13	38.47	57.83	42.22	59.06	31.60	51.34
SDC [52]	55.20	68.64	41.63	60.43	45.03	<u>61.75</u>	32.15	<u>53.18</u>
PASS [55]	34.04	49.00	26.37	41.08	20.71	37.13	12.30	25.46
FeTrIL [35]	54.66	67.45	49.09	62.42	36.92	54.09	34.29	50.41
FeCAM [11]	53.47	66.39	<u>51.78</u>	<u>64.97</u>	40.64	56.24	<u>37.50</u>	52.78
ADC (Ours)	<b>64.46</b>	<b>73.49</b>	<b>57.97</b>	<b>68.91</b>	<b>54.86</b>	<b>67.07</b>	<b>45.07</b>	<b>61.39</b>

# Experiments



## Drift Estimation Quality

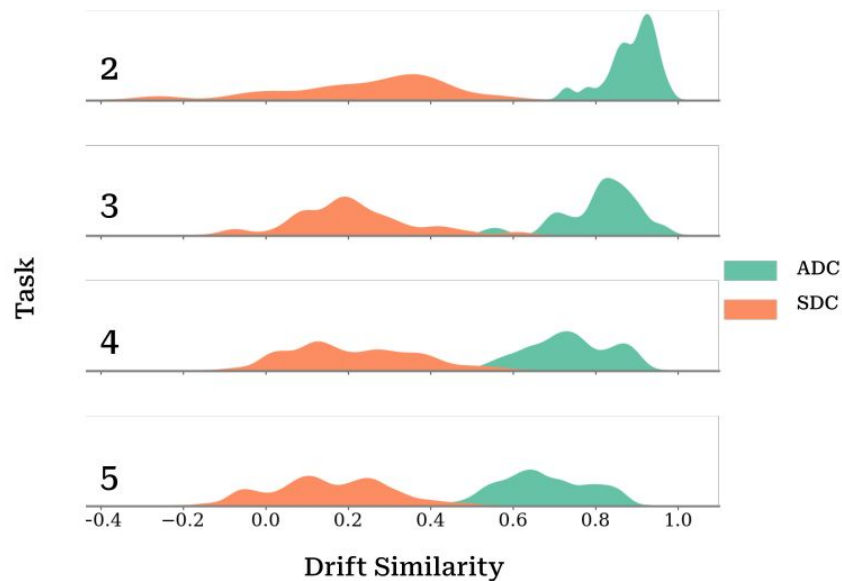
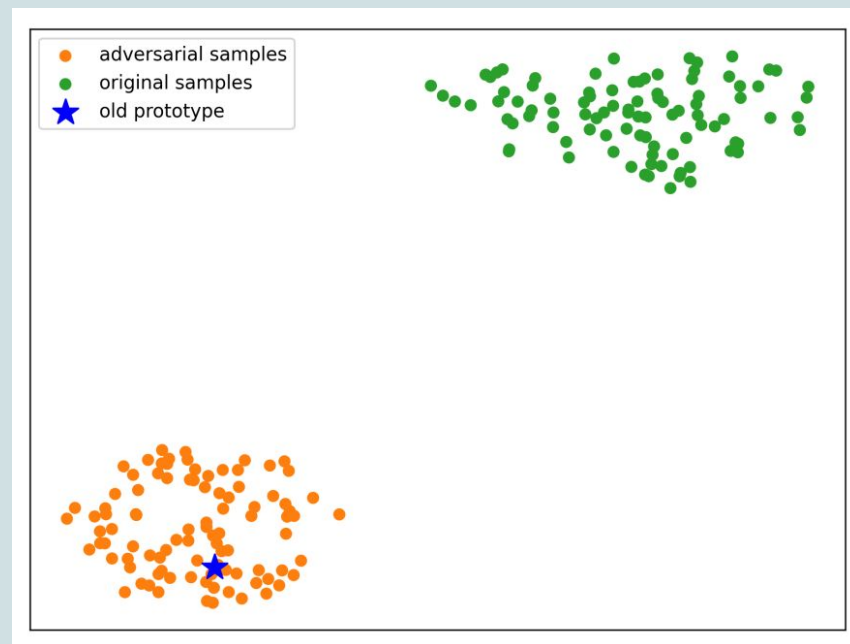


Figure 6. Comparison between the two drift estimation methods SDC [52], and the proposed ADC, on CIFAR-100 (5 tasks). We compute the drift for each class with the two methods and report the distribution of drift estimation quality, measured by computing the cosine similarity between the estimated drift vector and the true drift (obtained using old data), for all previous class prototypes.

## t-SNE Visualization



Thanks