

# Dataset Distillation

**Yufang Liu**

*East China Normal University*



# Background

- Knowledge Distillation
  - large model -> small model
- Dataset Distillation
  - large dataset -> small dataset
  - coreset & **synthetic**

# DATASET DISTILLATION

**Tongzhou Wang**

Facebook AI Research, MIT CSAIL

**Jun-Yan Zhu**

MIT CSAIL

**Antonio Torralba**

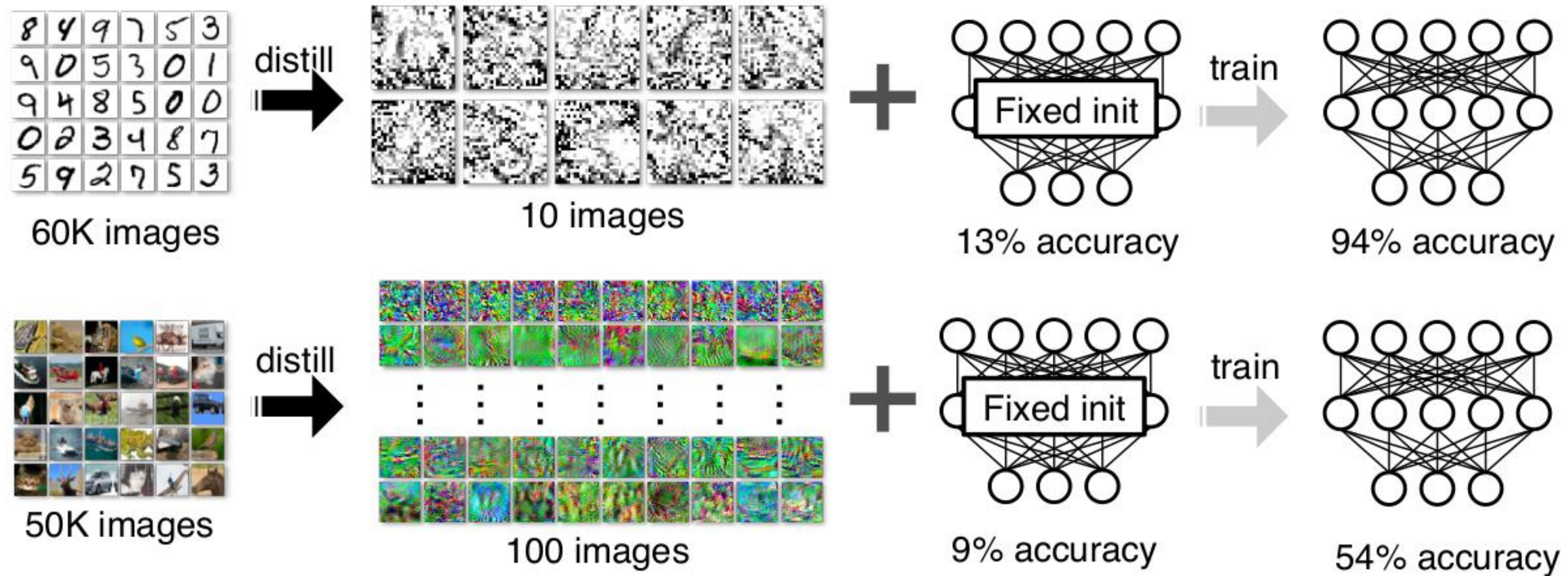
MIT CSAIL

**Alexei A. Efros**

UC Berkeley

*2018*

# Motivation



(a) Dataset distillation on MNIST and CIFAR10



# Methods

training loss  $\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell(x_i, \theta) \triangleq \arg \min_{\theta} \ell(\mathbf{x}, \theta), \quad \mathbf{x} = \{x_i\}_{i=1}^N$

SGD update  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \ell(\mathbf{x}_t, \theta_t),$

SGD update  
on syntactic  
data  $\theta_1 = \theta_0 - \tilde{\eta} \nabla_{\theta_0} \ell(\tilde{\mathbf{x}}, \theta_0) \quad \tilde{\mathbf{x}} = \{\tilde{x}_i\}_{i=1}^M \text{ with } M \ll N$

meta-learning  $\tilde{\mathbf{x}}^*, \tilde{\eta}^* = \arg \min_{\tilde{\mathbf{x}}, \tilde{\eta}} \mathcal{L}(\tilde{\mathbf{x}}, \tilde{\eta}; \theta_0) = \arg \min_{\tilde{\mathbf{x}}, \tilde{\eta}} \ell(\mathbf{x}, \theta_1) = \arg \min_{\tilde{\mathbf{x}}, \tilde{\eta}} \ell(\mathbf{x}, \theta_0 - \tilde{\eta} \nabla_{\theta_0} \ell(\tilde{\mathbf{x}}, \theta_0)),$

dependent on  $\theta_0$

# Methods

- Distillation for random initializations

$$\tilde{\mathbf{x}}^*, \tilde{\eta}^* = \arg \min_{\tilde{\mathbf{x}}, \tilde{\eta}} \mathbb{E}_{\theta_0 \sim p(\theta_0)} \mathcal{L}(\tilde{\mathbf{x}}, \tilde{\eta}; \theta_0)$$

---

**Algorithm 1** Dataset Distillation

---

**Input:**  $p(\theta_0)$ : distribution of initial weights;  $M$ : the number of distilled data

**Input:**  $\alpha$ : step size;  $n$ : batch size;  $T$ : the number of optimization iterations;  $\tilde{\eta}_0$ : initial value for  $\tilde{\eta}$

```
1: Initialize  $\tilde{\mathbf{x}} = \{\tilde{x}_i\}_{i=1}^M$  randomly,  $\tilde{\eta} \leftarrow \tilde{\eta}_0$ 
2: for each training step  $t = 1$  to  $T$  do
3:   Get a minibatch of real training data  $\mathbf{x}_t = \{x_{t,j}\}_{j=1}^n$ 
4:   Sample a batch of initial weights  $\theta_0^{(j)} \sim p(\theta_0)$ 
5:   for each sampled  $\theta_0^{(j)}$  do
6:     Compute updated parameter with GD:  $\theta_1^{(j)} = \theta_0^{(j)} - \tilde{\eta} \nabla_{\theta_0^{(j)}} \ell(\tilde{\mathbf{x}}, \theta_0^{(j)})$ 
7:     Evaluate the objective function on real training data:  $\mathcal{L}^{(j)} = \ell(\mathbf{x}_t, \theta_1^{(j)})$ 
8:   end for
9:   Update  $\tilde{\mathbf{x}} \leftarrow \tilde{\mathbf{x}} - \alpha \nabla_{\tilde{\mathbf{x}}} \sum_j \mathcal{L}^{(j)}$ , and  $\tilde{\eta} \leftarrow \tilde{\eta} - \alpha \nabla_{\tilde{\eta}} \sum_j \mathcal{L}^{(j)}$ 
10: end for
```

**Output:** distilled data  $\tilde{\mathbf{x}}$  and optimized learning rate  $\tilde{\eta}$

---

# Analysis

- origin loss

$$\ell(\mathbf{x}, \theta) = \ell((\mathbf{d}, \mathbf{t}), \theta) = \frac{1}{2N} \|\mathbf{d}\theta - \mathbf{t}\|^2. \quad \mathbf{d} \rightarrow N \times D$$

- DD loss

$$\ell(\mathbf{x}, \theta_0 - \tilde{\eta} \nabla_{\theta_0} \ell(\tilde{\mathbf{x}}, \theta_0)) \quad \tilde{\mathbf{d}} \rightarrow M \times D$$

$$\theta_1 = \theta_0 - \tilde{\eta} \nabla_{\theta_0} \ell(\tilde{\mathbf{x}}, \theta_0) = \theta_0 - \frac{\tilde{\eta}}{M} \tilde{\mathbf{d}}^T (\tilde{\mathbf{d}}\theta_0 - \tilde{\mathbf{t}}) = (\mathbf{I} - \frac{\tilde{\eta}}{M} \tilde{\mathbf{d}}^T \tilde{\mathbf{d}}) \theta_0 + \frac{\tilde{\eta}}{M} \tilde{\mathbf{d}}^T \tilde{\mathbf{t}}.$$

- $\mathbf{d}^T \mathbf{d} \theta^* = \mathbf{d}^T \mathbf{t}$  for global minimum

$$\mathbf{d}^T \mathbf{d} (\mathbf{I} - \frac{\tilde{\eta}}{M} \tilde{\mathbf{d}}^T \tilde{\mathbf{d}}) \theta_0 + \frac{\tilde{\eta}}{M} \mathbf{d}^T \tilde{\mathbf{d}} \tilde{\mathbf{d}}^T \tilde{\mathbf{t}} = \mathbf{d}^T \mathbf{t}$$

- Mild assumption: columns of  $\mathbf{d}$  are independent (i.e.  $\mathbf{d}^T \mathbf{d}$  has full rank)

$$\text{Satisfy for any } \theta_0 \longrightarrow \mathbf{I} - \frac{\tilde{\eta}}{M} \tilde{\mathbf{d}}^T \tilde{\mathbf{d}} = \mathbf{0}, \longrightarrow \tilde{\mathbf{d}}^T \tilde{\mathbf{d}} \text{ has full rank and } M \geq D.$$

any small number of distilled data fail to generalize to arbitrary initial  $\theta_0$



# Methods

- Multiple GD steps and epochs

---

**Algorithm 1** Dataset Distillation

---

**Input:**  $p(\theta_0)$ : distribution of initial weights;  $M$ : the number of distilled data

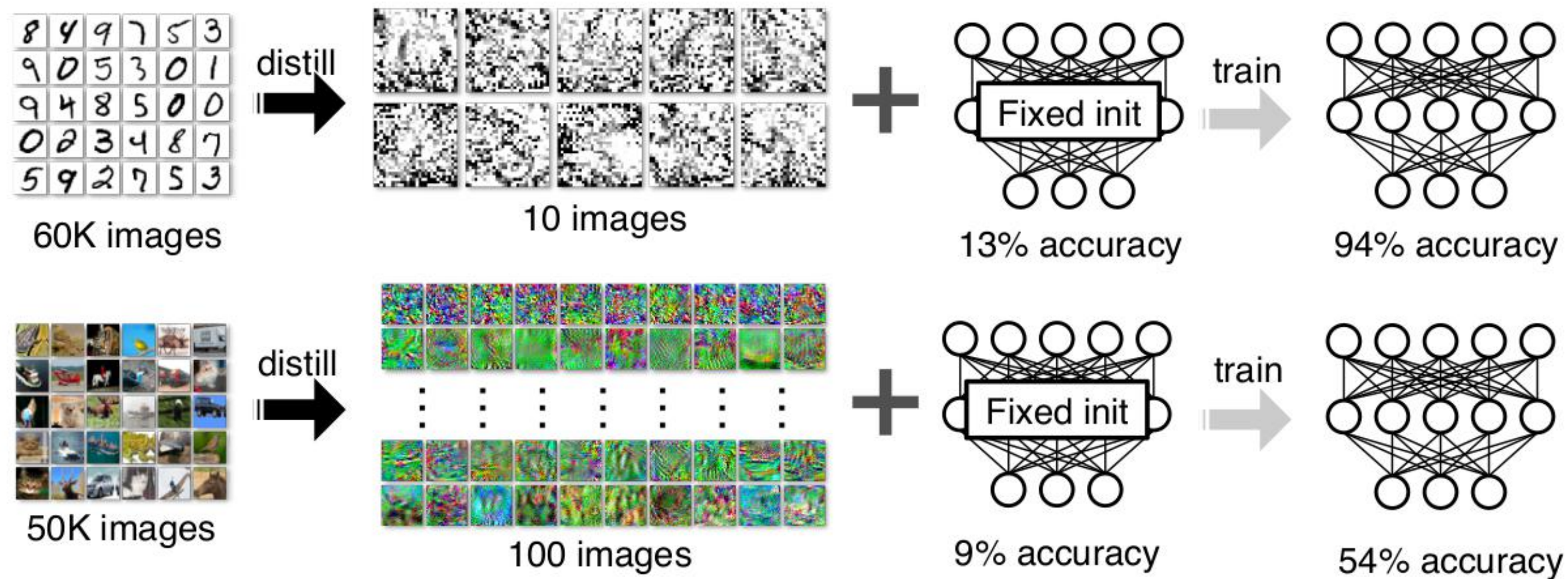
**Input:**  $\alpha$ : step size;  $n$ : batch size;  $T$ : the number of optimization iterations;  $\tilde{\eta}_0$ : initial value for  $\tilde{\eta}$

- 1: Initialize  $\tilde{\mathbf{x}} = \{\tilde{x}_i\}_{i=1}^M$  randomly,  $\tilde{\eta} \leftarrow \tilde{\eta}_0$
  - 2: **for each** training step  $t = 1$  to  $T$  **do** → multiple epochs over the same sequence of distilled data
  - 3:     Get a minibatch of real training data  $\mathbf{x}_t = \{x_{t,j}\}_{j=1}^n$
  - 4:     Sample a batch of initial weights  $\theta_0^{(j)} \sim p(\theta_0)$
  - 5:     **for each** sampled  $\theta_0^{(j)}$  **do**
  - 6:         Compute updated parameter with GD:  $\theta_1^{(j)} = \theta_0^{(j)} - \tilde{\eta} \nabla_{\theta_0^{(j)}} \ell(\tilde{\mathbf{x}}, \theta_0^{(j)})$  →  $\theta_{i+1} = \theta_i - \tilde{\eta}_i \nabla_{\theta_i} \ell(\tilde{\mathbf{x}}_i, \theta_i),$
  - 7:         Evaluate the objective function on real training data:  $\mathcal{L}^{(j)} = \ell(\mathbf{x}_t, \theta_1^{(j)})$
  - 8:     **end for**
  - 9:     Update  $\tilde{\mathbf{x}} \leftarrow \tilde{\mathbf{x}} - \alpha \nabla_{\tilde{\mathbf{x}}} \sum_j \mathcal{L}^{(j)}$ , and  $\tilde{\eta} \leftarrow \tilde{\eta} - \alpha \nabla_{\tilde{\eta}} \sum_j \mathcal{L}^{(j)}$  → back propagate through all steps, using HVP to accelerate
  - 10: **end for**
- Output:** distilled data  $\tilde{\mathbf{x}}$  and optimized learning rate  $\tilde{\eta}$
-



# Experiment

- Fixed initialization



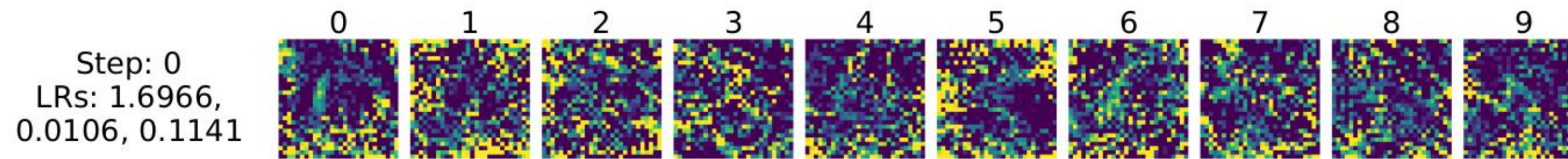
(a) Dataset distillation on MNIST and CIFAR10



# Experiment

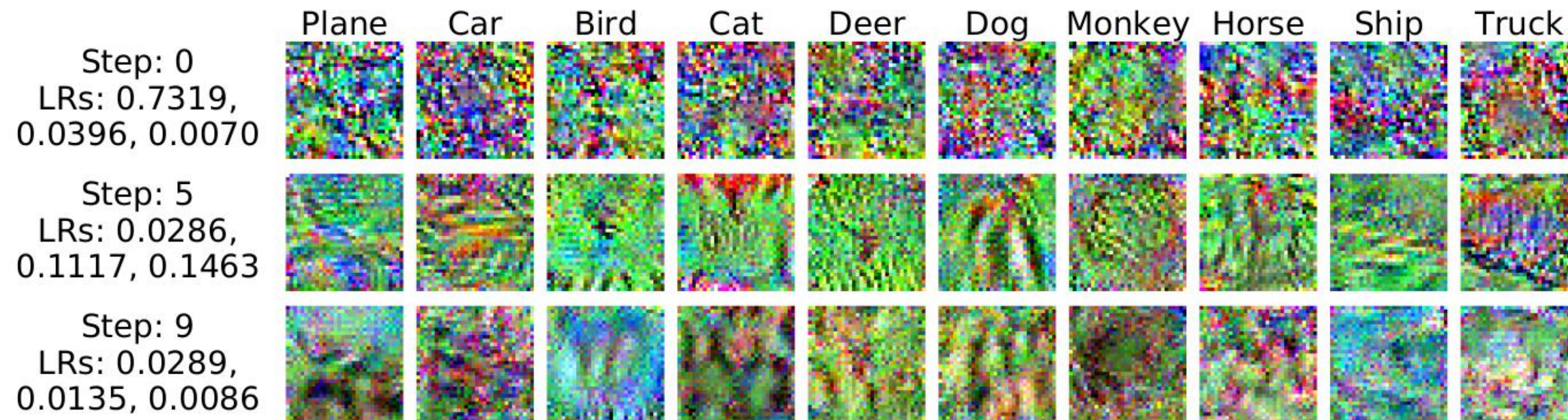
- Fixed initialization

1 GD step and  
3 epochs  
(10 imgs)



(a) MNIST. These distilled images train a fixed initialization from 12.90% test accuracy to 93.76%.

10 GD step  
and 3 epochs  
(100 imgs)



(b) CIFAR10. These distilled images train a fixed initialization from 8.82% test accuracy to 54.03%.

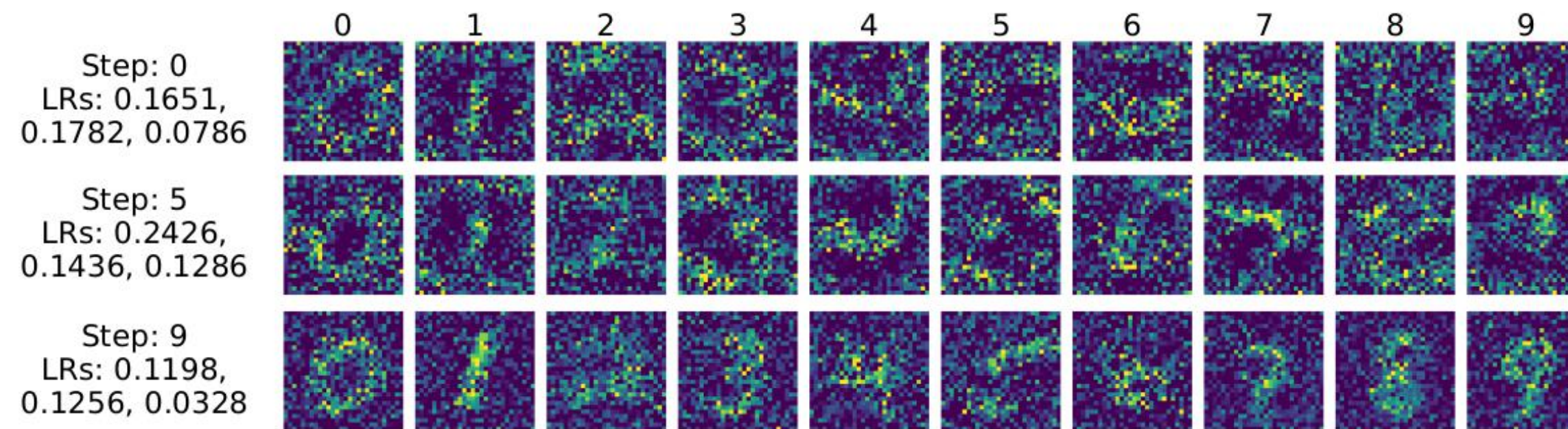
look like random noise as it encodes  $x$  and  $\theta_0$



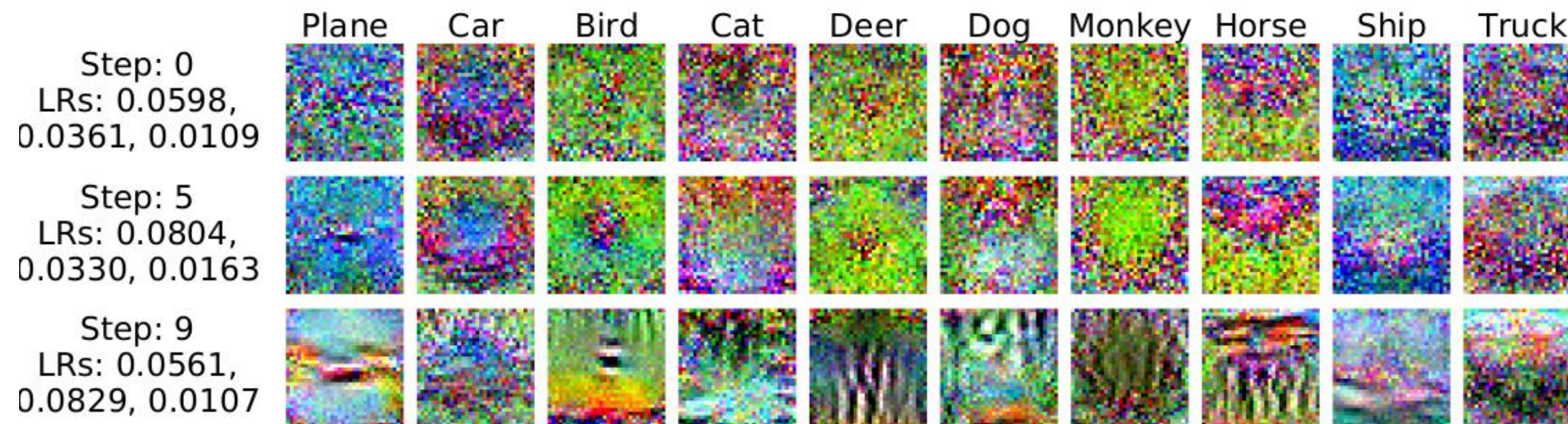
# Experiment

- Random initialization (Xavier Initialization)

10 GD step  
and 3 epochs  
(100 imgs)



(a) MNIST. These distilled images unknown random initializations to  $79.50\% \pm 8.08\%$  test accuracy.



(b) CIFAR10. These distilled images unknown random initializations to  $36.79\% \pm 1.18\%$  test accuracy.

each step containing  
one image per category.

Images used in early  
steps tend to look  
noisier



# Experiment

- Compare with baseline (Lenet)

99%  
80%

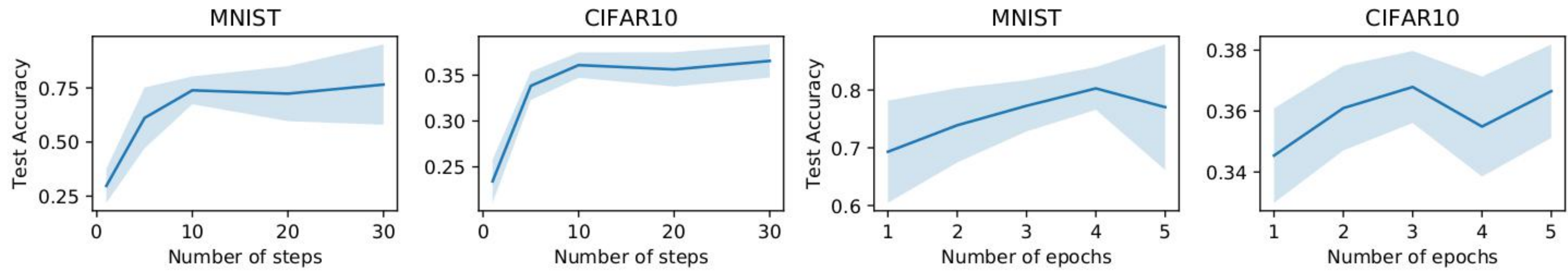
	Ours		Baselines					
	Fixed init.	Random init.	Used as training data in same number of GD steps				Used in K-NN classification	
			Random real	Optimized real	<i>k</i> -means	Average real	Random real	<i>k</i> -means
MNIST	<b>96.6</b>	79.5 ± 8.1	68.6 ± 9.8	73.0 ± 7.6	76.4 ± 9.5	77.1 ± 2.7	71.5 ± 2.1	<b>92.2 ± 0.1</b>
CIFAR10	<b>54.0</b>	<b>36.8 ± 1.2</b>	21.3 ± 1.5	23.4 ± 1.3	22.5 ± 3.1	22.3 ± 0.7	18.8 ± 1.3	29.4 ± 0.3

- Random real images
- Optimized real images -> top 20% from random sets
- *k*-means -> use cluster centroids for each category
- Average real images -> average image
- All methods use ten images per category (100 in total) except average real



# Experiment

- Hyper-parameter sensitivity studies



(a) Test accuracy w.r.t. the number of steps

(b) Test accuracy w.r.t. the number of epochs

Figure 4: Hyper-parameter sensitivity studies: we evaluate models in *random initialization* settings. (a) Average test accuracy w.r.t. the number of GD steps. We use two epochs. (b) Average test accuracy w.r.t. the number of epochs. We use 10 steps, with each step containing 10 images.

# Experiment

- Few-shot domain adaptation

	Ours w/ fixed pre-trained	Ours w/ random pre-trained	Random real	Optimized real	$k$ -means	Average real	Adaptation Motiian et al. (2017)	No adaptation	Train on <b>full</b> target dataset
$\mathcal{M} \rightarrow \mathcal{U}$	<b>97.9</b>	<b>95.4 <math>\pm</math> 1.8</b>	94.9 $\pm$ 0.8	95.2 $\pm$ 0.7	92.2 $\pm$ 1.6	93.9 $\pm$ 0.8	<b>96.7 <math>\pm</math> 0.5</b>	90.4 $\pm$ 3.0	97.3 $\pm$ 0.3
$\mathcal{U} \rightarrow \mathcal{M}$	<b>93.2</b>	<b>92.7 <math>\pm</math> 1.4</b>	87.1 $\pm$ 2.9	87.6 $\pm$ 2.1	85.6 $\pm$ 3.1	78.4 $\pm$ 5.0	89.2 $\pm$ 2.4	67.5 $\pm$ 3.9	98.6 $\pm$ 0.5
$\mathcal{S} \rightarrow \mathcal{M}$	<b>96.2</b>	85.2 $\pm$ 4.7	84.6 $\pm$ 2.1	85.2 $\pm$ 1.2	<b>85.8 <math>\pm</math> 1.2</b>	74.9 $\pm$ 2.6	74.0 $\pm$ 1.5	51.6 $\pm$ 2.8	98.6 $\pm$ 0.5

Table 2: Performance of our method and baselines in adapting models among MNIST ( $\mathcal{M}$ ), USPS ( $\mathcal{U}$ ), and SVHN ( $\mathcal{S}$ ). 100 distilled images are trained for ten GD steps and three epochs. Our method outperforms few-shot domain adaptation (Motiian et al., 2017) and other baselines in most settings.

# DATASET CONDENSATION WITH GRADIENT MATCHING

**Bo Zhao, Konda Reddy Mopuri, Hakan Bilen**  
School of Informatics, The University of Edinburgh  
`{bo.zhao, kmopuri, hbilen}@ed.ac.uk`

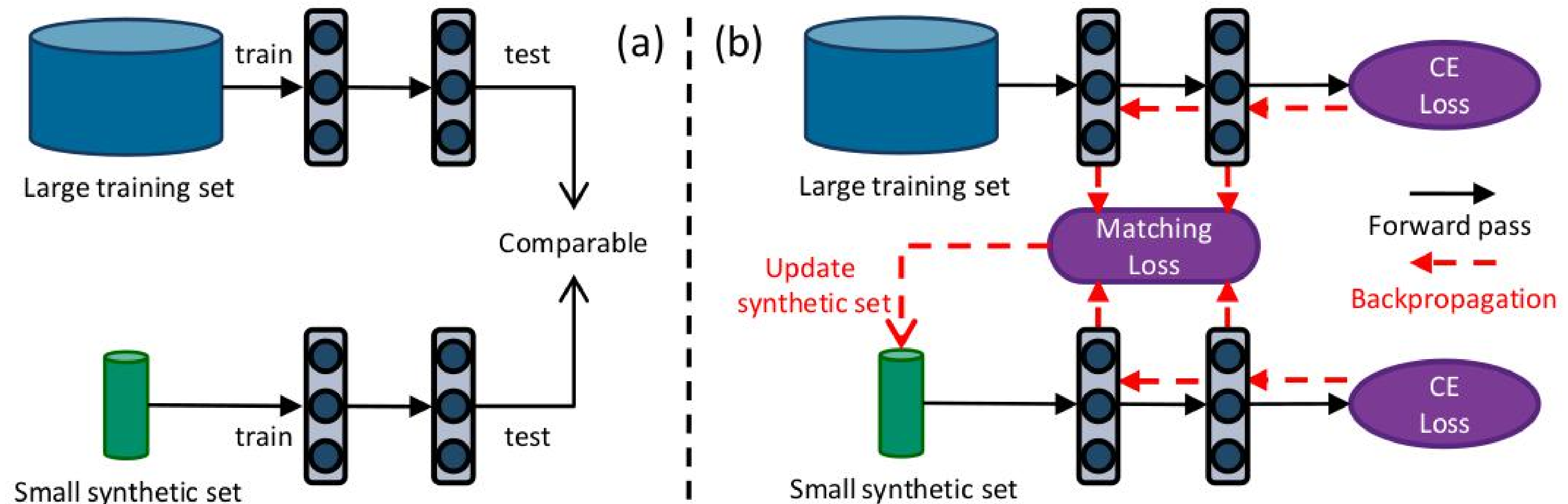
*ICLR 2021*

# Motivation

- Coreset construction
  - continual learning, active learning ...
  - compactness, diversity, forgetfulness ...
  - not guarantee for downstream tasks and presence of representative samples
- Dataset Distillation
  - synthesize informative samples by optimization



# Methods



learning a synthetic set such that a deep network trained on it and the large set produces similar gradients w.r.t. its weights.

# Methods

$$\begin{array}{ccc} \boldsymbol{\theta}^{\mathcal{T}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}) & \xleftrightarrow[|S| \ll |\mathcal{T}|]{} & \boldsymbol{\theta}^{\mathcal{S}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}) \\ \mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}, y) \in \mathcal{T}} \ell(\phi_{\boldsymbol{\theta}}(\mathbf{x}), y) & & \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{s}, y) \in \mathcal{S}} \ell(\phi_{\boldsymbol{\theta}}(\mathbf{s}), y) \end{array}$$

- Discussion

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathcal{L}^{\mathcal{T}}(\boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S})) \quad \text{subject to} \quad \boldsymbol{\theta}^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\boldsymbol{\theta}} \mathcal{L}^{\mathcal{S}}(\boldsymbol{\theta}).$$

- nested loop optimization
- unrolling the recursive computation graph over multiple steps

# Methods

- Optimization

- comparable performance && similar in parameter space  $\theta^S \approx \theta^T$
- assume similar weights imply similar performance

$$\min_S D(\theta^S, \theta^T) \quad \text{subject to} \quad \theta^S(S) = \arg \min_{\theta} \mathcal{L}^S(\theta)$$

- $\theta^T$  depends on  $\theta_0$

$$\min_S \mathbb{E}_{\theta_0 \sim P_{\theta_0}} [D(\theta^S(\theta_0), \theta^T(\theta_0))] \quad \text{subject to} \quad \theta^S(S) = \arg \min_{\theta} \mathcal{L}^S(\theta(\theta_0))$$

- back-optimization approach (fixed steps ( $\varsigma$ ))

$$\theta^S(S) = \text{opt-}\arg_{\theta}(\mathcal{L}^S(\theta), \varsigma)$$

distance can be large and fixed steps may be not enough

# Methods

- Curriculum based approach
  - follow a similar path throughout the optimization

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[ \sum_{t=0}^{T-1} D(\theta_t^{\mathcal{S}}, \theta_t^{\mathcal{T}}) \right] \quad \text{subject to}$$

$$\theta_{t+1}^{\mathcal{S}}(\mathcal{S}) = \text{opt-alg}_{\theta}(\mathcal{L}^{\mathcal{S}}(\theta_t^{\mathcal{S}}), \varsigma^{\mathcal{S}}) \quad \text{and} \quad \theta_{t+1}^{\mathcal{T}} = \text{opt-alg}_{\theta}(\mathcal{L}^{\mathcal{T}}(\theta_t^{\mathcal{T}}), \varsigma^{\mathcal{T}})$$

- In preliminary experiments  $D(\theta_t^{\mathcal{S}}, \theta_t^{\mathcal{T}}) \rightarrow 0$

$$\theta_{t+1}^{\mathcal{S}} \leftarrow \theta_t^{\mathcal{S}} - \eta_{\theta} \nabla_{\theta} \mathcal{L}^{\mathcal{S}}(\theta_t^{\mathcal{S}}) \quad \text{and} \quad \theta_{t+1}^{\mathcal{T}} \leftarrow \theta_t^{\mathcal{T}} - \eta_{\theta} \nabla_{\theta} \mathcal{L}^{\mathcal{T}}(\theta_t^{\mathcal{T}}),$$

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[ \sum_{t=0}^{T-1} D(\nabla_{\theta} \mathcal{L}^{\mathcal{S}}(\theta_t), \nabla_{\theta} \mathcal{L}^{\mathcal{T}}(\theta_t)) \right].$$

does not require the expensive unrolling



# Methods

---

**Algorithm 1:** Dataset condensation with gradient matching

---

**Input:** Training set  $\mathcal{T}$

1 **Required:** Randomly initialized set of synthetic samples  $\mathcal{S}$  for  $C$  classes, probability distribution over randomly initialized weights  $P_{\theta_0}$ , deep neural network  $\phi_{\theta}$ , number of outer-loop steps  $K$ , number of inner-loop steps  $T$ , number of steps for updating weights  $\varsigma_{\theta}$  and synthetic samples  $\varsigma_{\mathcal{S}}$  in each inner-loop step respectively, learning rates for updating weights  $\eta_{\theta}$  and synthetic samples  $\eta_{\mathcal{S}}$ .

2 **for**  $k = 0, \dots, K - 1$  **do**     Loop for initialization

3     Initialize  $\theta_0 \sim P_{\theta_0}$

4     **for**  $t = 0, \dots, T - 1$  **do**     steps

5         **for**  $c = 0, \dots, C - 1$  **do**     class

6             Sample a minibatch pair  $B_c^{\mathcal{T}} \sim \mathcal{T}$  and  $B_c^{\mathcal{S}} \sim \mathcal{S}$       $\triangleright B_c^{\mathcal{T}}$  and  $B_c^{\mathcal{S}}$  are of the same class  $c$ .

7             Compute  $\mathcal{L}_c^{\mathcal{T}} = \frac{1}{|B_c^{\mathcal{T}}|} \sum_{(\mathbf{x}, y) \in B_c^{\mathcal{T}}} \ell(\phi_{\theta_t}(\mathbf{x}), y)$  and  $\mathcal{L}_c^{\mathcal{S}} = \frac{1}{|B_c^{\mathcal{S}}|} \sum_{(\mathbf{s}, y) \in B_c^{\mathcal{S}}} \ell(\phi_{\theta_t}(\mathbf{s}), y)$

8             Update  $\mathcal{S}_c \leftarrow \text{opt-alg}_{\mathcal{S}}(D(\nabla_{\theta} \mathcal{L}_c^{\mathcal{S}}(\theta_t), \nabla_{\theta} \mathcal{L}_c^{\mathcal{T}}(\theta_t)), \varsigma_{\mathcal{S}}, \eta_{\mathcal{S}})$

9         Update  $\theta_{t+1} \leftarrow \text{opt-alg}_{\theta}(\mathcal{L}^{\mathcal{S}}(\theta_t), \varsigma_{\theta}, \eta_{\theta})$       $\triangleright$  Use the whole  $\mathcal{S}$

**Output:**  $\mathcal{S}$

---

# Methods

- Matching loss  $D(\cdot, \cdot)$ 
  - different layers -> different parameter shape
  - MLP: 2D (out×in), CNN: 4D (out×in×h×w)

$$D(\nabla_{\theta} \mathcal{L}^S, \nabla_{\theta} \mathcal{L}^T) = \sum_{l=1}^L d(\nabla_{\theta^{(l)}} \mathcal{L}^S, \nabla_{\theta^{(l)}} \mathcal{L}^T)$$

$$d(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^{\text{out}} \left( 1 - \frac{\mathbf{A}_{i\cdot} \cdot \mathbf{B}_{i\cdot}}{\|\mathbf{A}_{i\cdot}\| \|\mathbf{B}_{i\cdot}\|} \right)$$

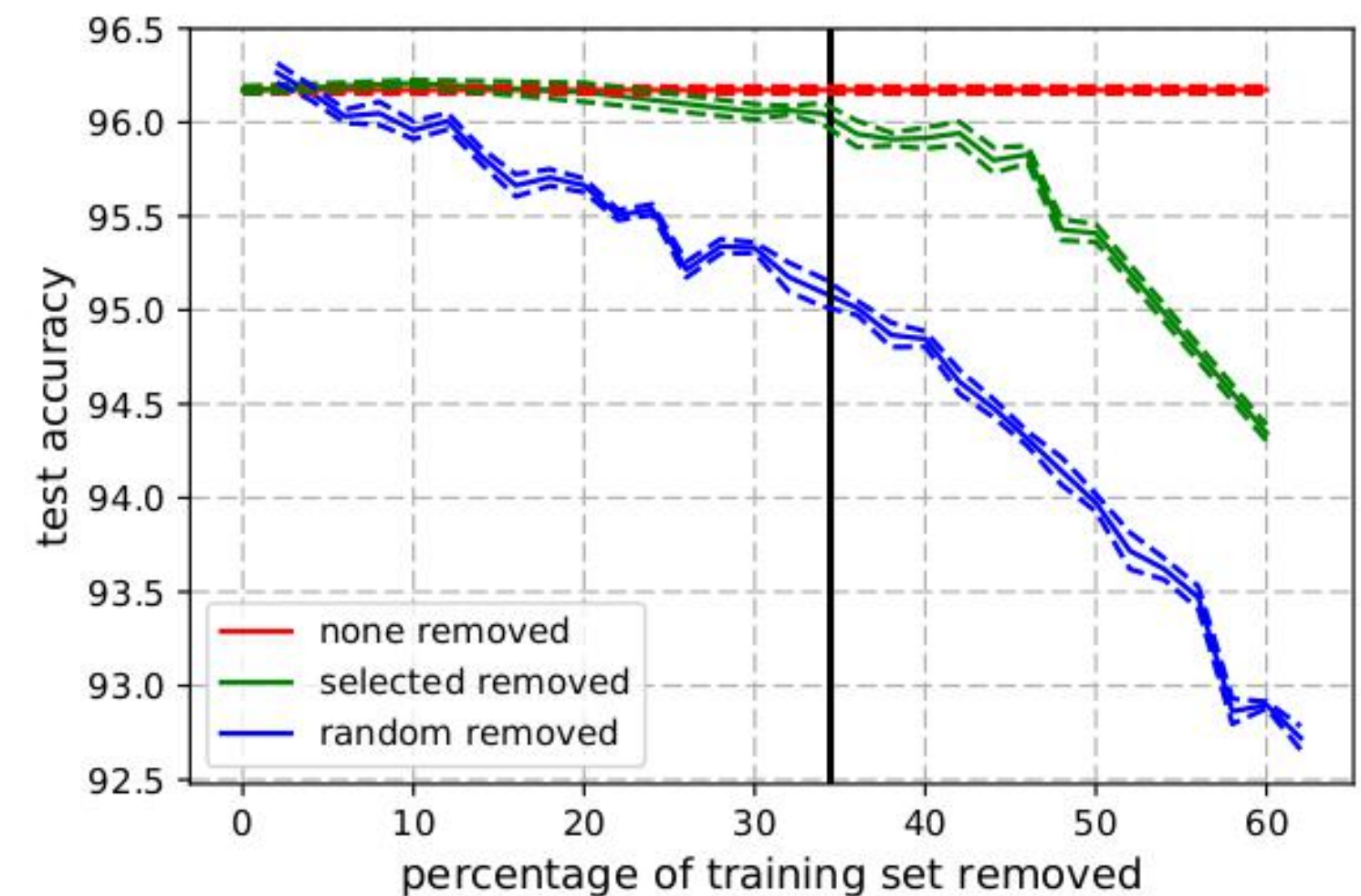
cosine similarity

$\mathbf{A}_{i\cdot}$  and  $\mathbf{B}_{i\cdot}$  are flattened gradients corresponding to each output node  $i$

# Experiment

- Baseline
  - random, DD
  - Herding baseline-> nearest to cluster center
  - K-Center-> multiple center points
  - Forgetting method <sup>[1]</sup>

the unforgettable examples contain less information



[1] An Empirical Study of Example Forgetting during Deep Neural Network Learning, Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, Geoffrey J. Gordon, ICLR 2019.



# Experiment

	Img/Cls	Ratio %	Coreset Selection			Forgetting	Ours	Whole Dataset
			Random	Herding	K-Center			
MNIST	1	0.017	64.9±3.5	89.2±1.6	89.3±1.5	35.5±5.6	<b>91.7±0.5</b>	99.6±0.0
	10	0.17	95.1±0.9	93.7±0.3	84.4±1.7	68.1±3.3	<b>97.4±0.2</b>	
	50	0.83	97.9±0.2	94.9±0.2	97.4±0.3	88.2±1.2	<b>98.8±0.2</b>	
FashionMNIST	1	0.017	51.4±3.8	67.0±1.9	66.9±1.8	42.0±5.5	<b>70.5±0.6</b>	93.5±0.1
	10	0.17	73.8±0.7	71.1±0.7	54.7±1.5	53.9±2.0	<b>82.3±0.4</b>	
	50	0.83	82.5±0.7	71.9±0.8	68.3±0.8	55.0±1.1	<b>83.6±0.4</b>	
SVHN	1	0.014	14.6±1.6	20.9±1.3	21.0±1.5	12.1±1.7	<b>31.2±1.4</b>	95.4±0.1
	10	0.14	35.1±4.1	50.5±3.3	14.0±1.3	16.8±1.2	<b>76.1±0.6</b>	
	50	0.7	70.9±0.9	72.6±0.8	20.1±1.4	27.2±1.5	<b>82.3±0.3</b>	
CIFAR10	1	0.02	14.4±2.0	21.5±1.2	21.5±1.3	13.5±1.2	<b>28.3±0.5</b>	84.8±0.1
	10	0.2	26.0±1.2	31.6±0.7	14.7±0.9	23.3±1.0	<b>44.9±0.5</b>	
	50	1	43.4±1.0	40.4±0.6	27.0±1.4	23.3±1.1	<b>53.9±0.5</b>	

ConvNet is used for training and testing, commonly used in fewshot learning



# Experiment

- Comapre with DD

Dataset	Img/Cls	DD	Ours	Whole Dataset
MNIST	1	-	$85.0 \pm 1.6$	$99.5 \pm 0.0$
	10	$79.5 \pm 8.1$	<b><math>93.9 \pm 0.6</math></b>	
CIFAR10	1	-	$24.2 \pm 0.9$	$83.1 \pm 0.2$
	10	$36.8 \pm 1.2$	<b><math>39.1 \pm 1.2</math></b>	

Table 3: Comparison to DD (Wang et al., 2018) in terms of testing accuracy (%).

2 times faster and more consistent results over multiple runs

# Experiment

- Visualization



Figure 2: Visualization of condensed 1 image/class with ConvNet for MNIST, Fashion-MNIST, SVHN and CIFAR10.

- Cross-architecture performance

C\T	MLP	ConvNet	LeNet	AlexNet	VGG	ResNet
MLP	$70.5 \pm 1.2$	$63.9 \pm 6.5$	$77.3 \pm 5.8$	$70.9 \pm 11.6$	$53.2 \pm 7.0$	$80.9 \pm 3.6$
ConvNet	$69.6 \pm 1.6$	<b><math>91.7 \pm 0.5</math></b>	$85.3 \pm 1.8$	$85.1 \pm 3.0$	<b><math>83.4 \pm 1.8</math></b>	<b><math>90.0 \pm 0.8</math></b>
LeNet	$71.0 \pm 1.6$	$90.3 \pm 1.2$	$85.0 \pm 1.7$	$84.7 \pm 2.4$	$80.3 \pm 2.7$	$89.0 \pm 0.8$
AlexNet	$72.1 \pm 1.7$	$87.5 \pm 1.6$	$84.0 \pm 2.8$	$82.7 \pm 2.9$	$81.2 \pm 3.0$	$88.9 \pm 1.1$
VGG	$70.3 \pm 1.6$	$90.1 \pm 0.7$	$83.9 \pm 2.7$	$83.4 \pm 3.7$	$81.7 \pm 2.6$	$89.1 \pm 0.9$
ResNet	<b><math>73.6 \pm 1.2</math></b>	$91.6 \pm 0.5$	<b><math>86.4 \pm 1.5</math></b>	<b><math>85.4 \pm 1.9</math></b>	<b><math>83.4 \pm 2.4</math></b>	$89.4 \pm 0.9$

Table 2: Cross-architecture performance in testing accuracy (%) for condensed 1 image/class in MNIST.

mlp generated pics performs badly for cnn, while  
cnn generated pics perform well for other mlp



# Experiment

- Continual Learning

testing accuracies are computed by averaging performance after each stage (SVHN→MNIST→USPS)

10 images/class

KD: regularize output vs previous

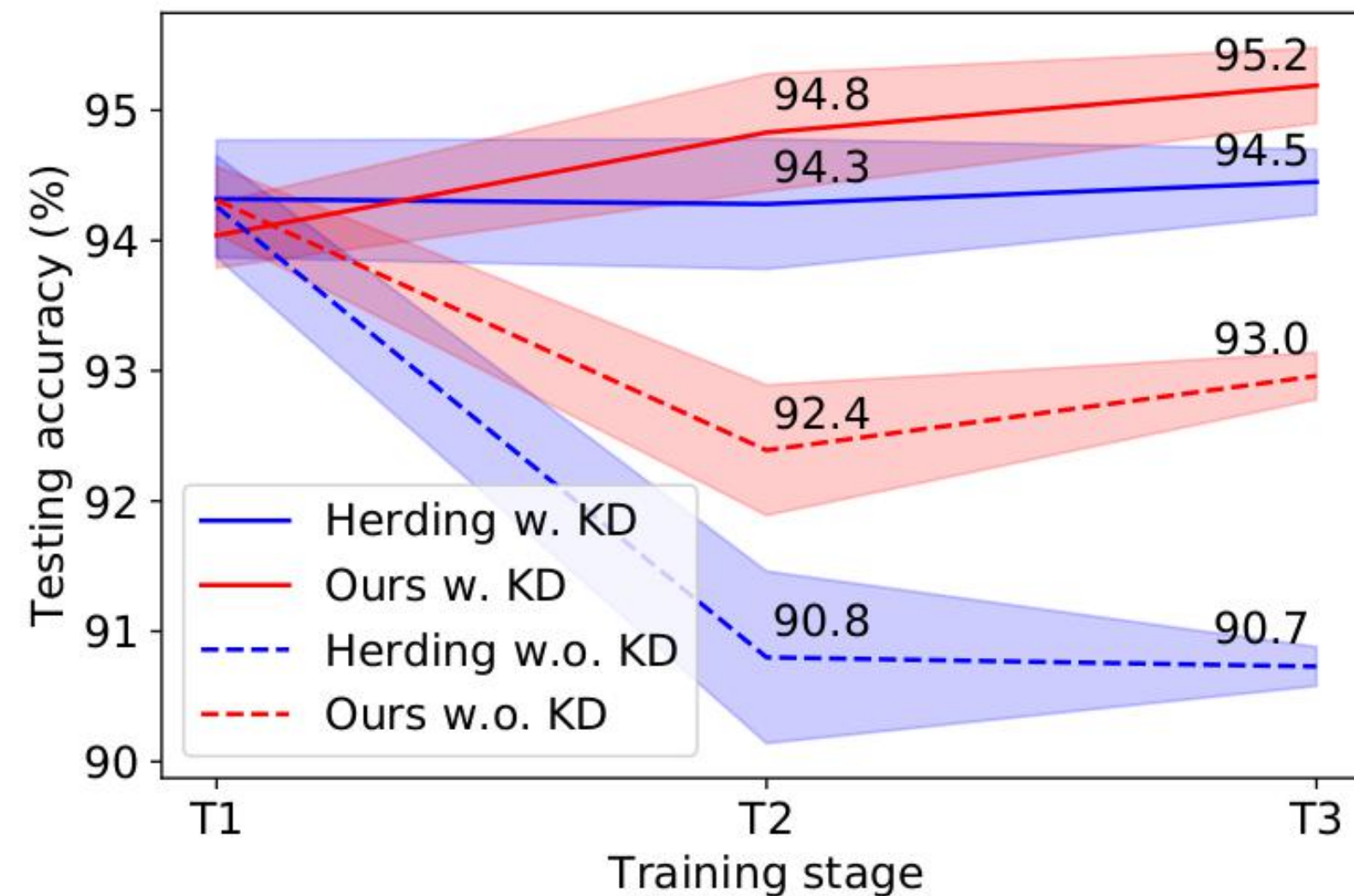


Figure 4: Continual learning performance in accuracy (%). Herding denotes the original E2E (Castro et al. 2018). T1, T2, T3 are three learning stages. The performance at each stage is the mean testing accuracy on all learned tasks.