

M3D: Dataset Condensation by Minimizing Maximum Mean Discrepancy

Hansong Zhang, Shikun Li, Pengju Wang, Dan Zeng, Shiming Ge

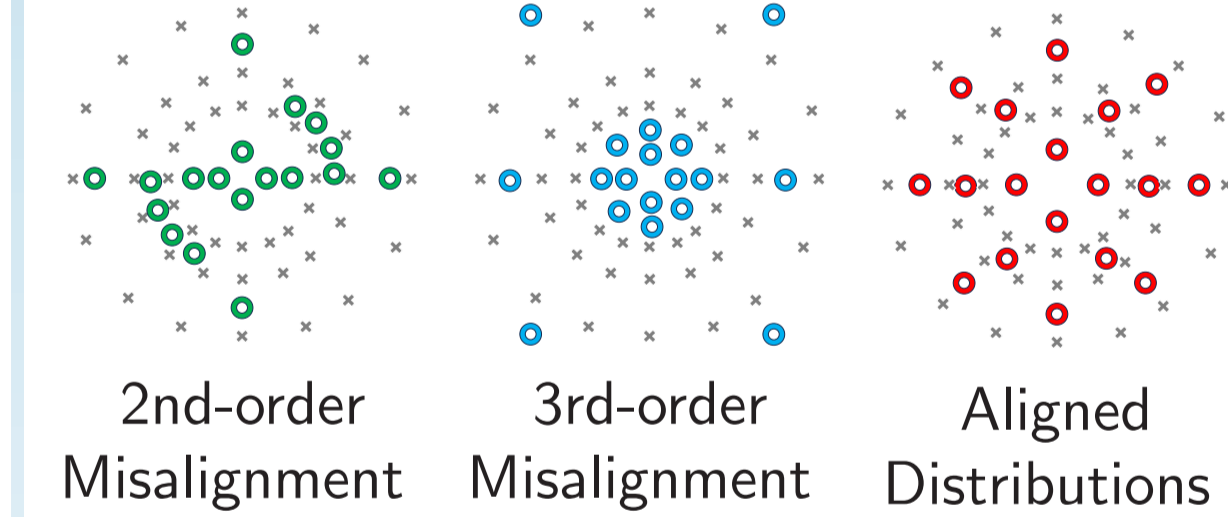
Summary

Background:

Dataset condensation (also known as dataset distillation) is a process aimed at addressing the challenges associated with the extensive data requirements of training state-of-the-art deep models. It involves creating a small synthetic dataset that retains the essential information of the original large-scale dataset.

Motivation:

- Traditional bi-level optimization used in dataset distillation, while effective, are often impractical for larger datasets due to their computational complexity and inefficiency.
- Distribution-Matching (DM) methods focus on aligning the feature distributions of synthetic and real data, offering greater efficiency but producing less informative examples compared to Optimization-Oriented methods.
- Previous DM-based methods only aligns the first-order moment of the synthetic and real data, which is not sufficient for matching two distributions.

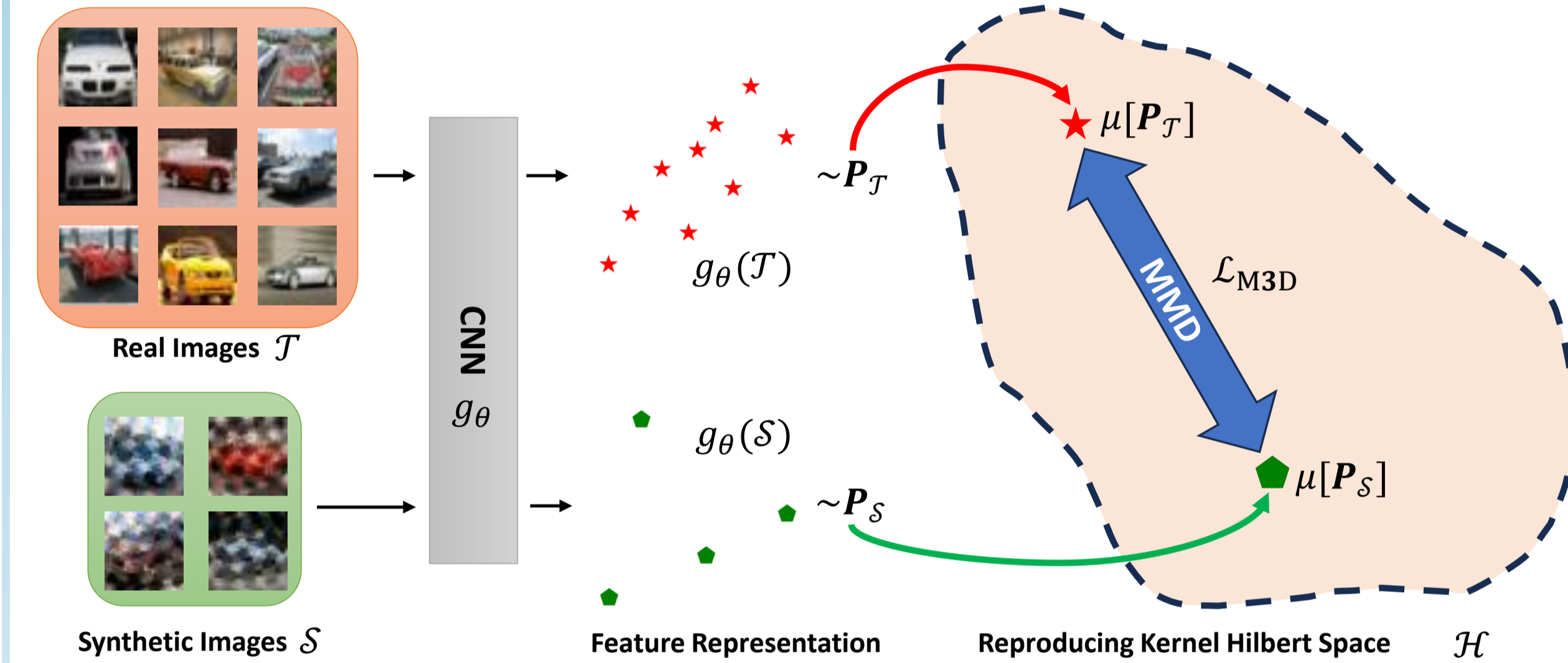


Contribution:

- We introduce a novel DM-based method named M3D for dataset condensation. This method focuses on minimizing the Maximum Mean Discrepancy between feature representations of synthetic and real images, which is achieved by embedding their distributions in a reproducing kernel Hilbert space.
- Unlike previous DM methods that primarily align only the first moment of distributions, M3D aligns all orders of moments of the distributions of real and synthetic images.
- Extensive experiments are conducted on both low-resolution and high-resolution datasets, where the results indicate the superiority of M3D.

Methodology

Main idea: The core idea of M3D is (1) extract the feature representations of both synthetic and real data; (2) based on classical kernel method, further embed the feature representations to a reproducing kernel Hilbert space, where we can map the infinite order of moments into a form of kernel function.



Reproducing Kernel Hilbert Space (RKHS): Given a kernel \mathcal{K} , \mathcal{H} is a Hilbert space of functions $\mathcal{X} \rightarrow \mathbb{R}$ with dot product $\langle \cdot, \cdot \rangle$, if $\forall \phi$, satisfying the following properties:

Reproducing : $\langle \phi(\cdot), \mathcal{K}(x, \cdot) \rangle = \phi(x)$.

Symmetry : $\mathcal{K}(x, x') = \mathcal{K}(x', x)$

Positive : $\mathcal{K}(\cdot, \cdot) \geq 0$

Previous Distribution-Matching: Previous DM methods only align the first-order moment of synthetic and real data, which can be formulated as:

$$S^* = \arg \min_S E_{\theta \sim P_\theta} \left\| \frac{1}{|\mathcal{T}|} \sum_{i=1}^{|\mathcal{T}|} g_\theta(x_i) - \frac{1}{|S|} \sum_{j=1}^{|S|} g_\theta(s_j) \right\|^2$$

Embedding feature representations into RKHS: Our M3D further embeds the feature representations into a Reproducing Kernel Hilbert Space, where we can effectively align all order of moments between synthetic and real data:

$$\mathcal{L}_{M3D} = \text{MMD}^2(P_{\mathcal{T}}, P_S) = \hat{\mathcal{K}}_{\mathcal{T}, \mathcal{T}} + \hat{\mathcal{K}}_{S, S} - 2\hat{\mathcal{K}}_{\mathcal{T}, S}$$

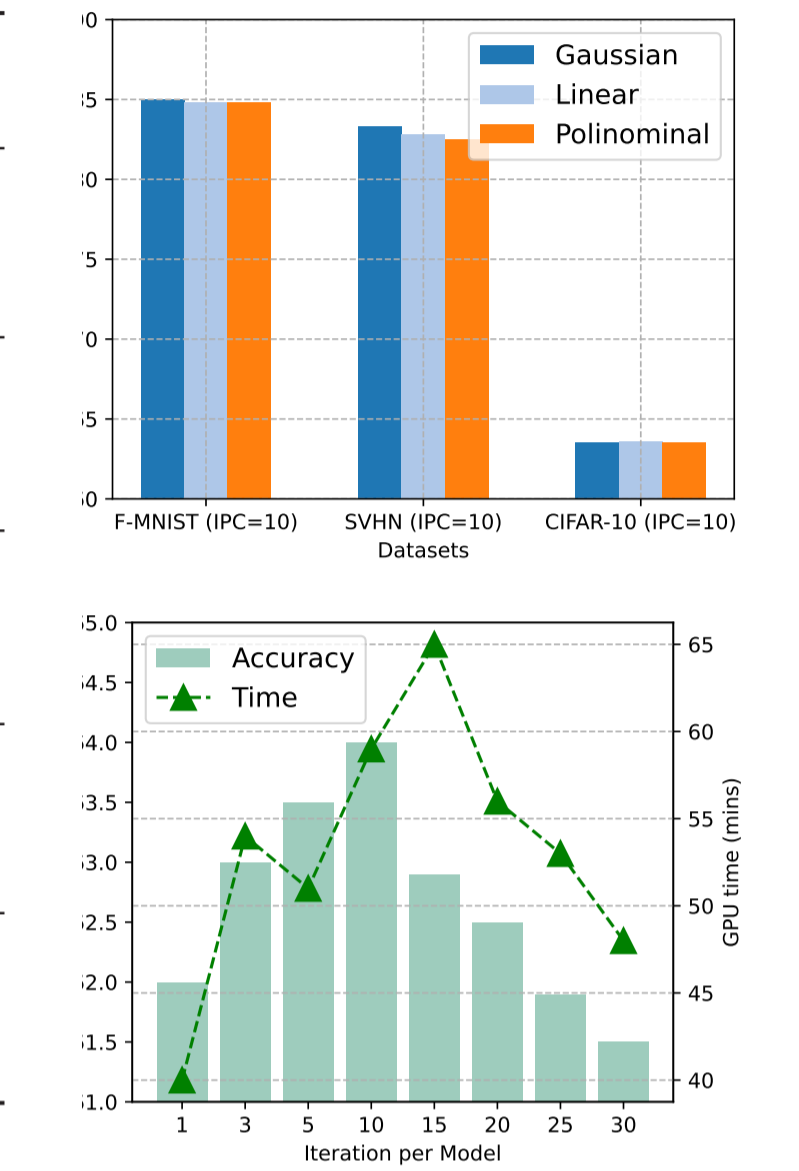
where $\hat{\mathcal{K}}_{X, Y} = \frac{1}{|X| \cdot |Y|} \sum_{i=1}^{|X|} \sum_{j=1}^{|Y|} \mathcal{K}(g_\theta(x_i), g_\theta(y_j))$ with $\{x_i\}_{i=1}^{|X|} \sim X, \{y_j\}_{j=1}^{|Y|} \sim Y$.

Experiments

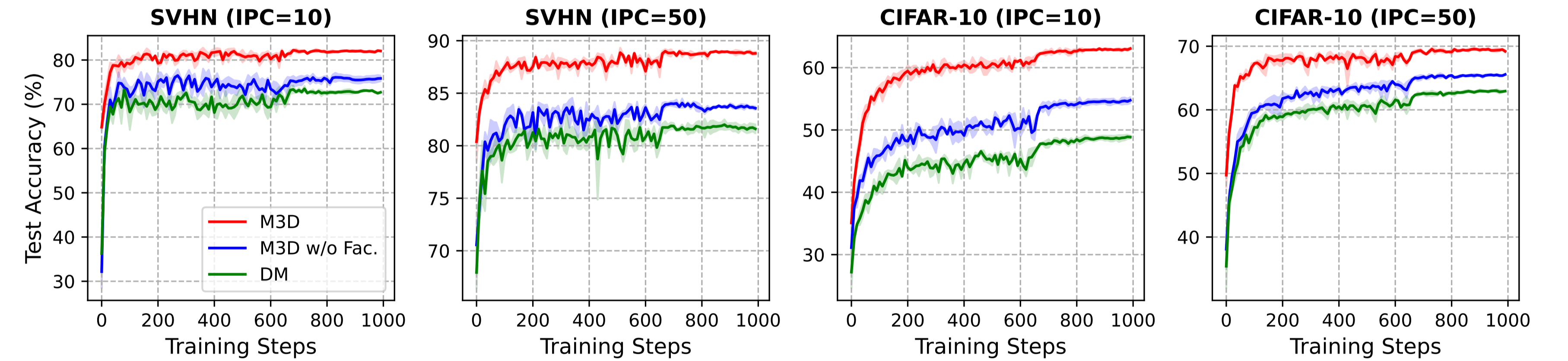
Experiments on low-resolution datasets:

Dataset	IPC	Ratio (%)	Coreset Selection			Dataset Condensation						Whole	
			Random	Herding	K-Center	DC	DSA	CAFE	CAFE+DSA	DM	IDM		M3D
MNIST	1	0.017	64.9 \pm 3.5	89.2 \pm 1.6	89.3 \pm 1.5	91.7 \pm 0.5	88.7 \pm 0.6	93.1 \pm 0.3	90.8 \pm 0.5	89.7 \pm 0.6	-	94.4 \pm 0.2	99.6 \pm 0.0
	10	0.17	95.1 \pm 0.9	93.7 \pm 0.3	84.4 \pm 1.7	97.4 \pm 0.2	97.8 \pm 0.1	97.2 \pm 0.2	97.5 \pm 0.1	97.5 \pm 0.1	-	97.6 \pm 0.1	
	50	0.83	97.9 \pm 0.2	94.8 \pm 0.2	97.4 \pm 0.3	98.8 \pm 0.2	99.2 \pm 0.1	98.6 \pm 0.2	98.9 \pm 0.2	98.6 \pm 0.1	-	98.2 \pm 0.2	
F-MNIST	1	0.017	51.4 \pm 3.8	67.0 \pm 1.9	66.9 \pm 1.8	70.5 \pm 0.6	70.6 \pm 0.6	77.1 \pm 0.9	73.7 \pm 0.7	70.7 \pm 0.6 †	-	80.7 \pm 0.3	93.5 \pm 0.1
	10	0.17	73.8 \pm 0.7	71.1 \pm 0.7	54.7 \pm 1.5	82.3 \pm 0.4	84.6 \pm 0.3	83.0 \pm 0.4	83.0 \pm 0.3	83.5 \pm 0.3 †	-	85.0 \pm 0.1	
	50	0.83	82.5 \pm 0.7	71.9 \pm 0.8	68.3 \pm 0.8	83.6 \pm 0.4	88.7 \pm 0.2	84.8 \pm 0.4	88.2 \pm 0.3	88.1 \pm 0.6 †	-	86.2 \pm 0.3	
SVHN	1	0.014	14.6 \pm 1.6	20.9 \pm 1.3	21.0 \pm 1.5	31.2 \pm 1.4	27.5 \pm 1.4	42.6 \pm 3.3	42.9 \pm 3.0	30.3 \pm 0.1 †	-	62.8 \pm 0.5	95.4 \pm 0.1
	10	0.14	35.1 \pm 4.1	50.5 \pm 3.3	14.0 \pm 1.3	76.1 \pm 0.6	79.2 \pm 0.5	75.9 \pm 0.6	77.9 \pm 0.6	73.5 \pm 0.5 †	-	83.3 \pm 0.7	
	50	0.7	70.9 \pm 0.9	72.6 \pm 0.8	20.1 \pm 1.4	82.3 \pm 0.3	84.4 \pm 0.4	81.3 \pm 0.3	82.3 \pm 0.4	82.0 \pm 0.2 †	-	89.0 \pm 0.2	
CIFAR-10	1	0.02	14.4 \pm 2.0	21.5 \pm 1.2	21.5 \pm 1.3	28.3 \pm 0.5	28.8 \pm 0.7	30.3 \pm 1.1	31.6 \pm 0.8	26.0 \pm 0.8	45.6 \pm 0.7	45.3 \pm 0.3	84.8 \pm 0.1
	10	0.2	26.0 \pm 1.2	31.6 \pm 0.7	14.7 \pm 0.9	44.9 \pm 0.5	52.1 \pm 0.5	46.3 \pm 0.6	50.9 \pm 0.5	48.9 \pm 0.6	58.6 \pm 0.1	63.5 \pm 0.2	
	50	1	43.4 \pm 1.0	40.4 \pm 0.6	27.0 \pm 1.4	53.9 \pm 0.5	60.6 \pm 0.5	55.5 \pm 0.6	62.3 \pm 0.4	63.0 \pm 0.4	67.5 \pm 0.1	69.9 \pm 0.5	
CIFAR-100	1	0.2	4.2 \pm 0.3	8.4 \pm 0.3	8.3 \pm 0.3	12.8 \pm 0.3	13.9 \pm 0.3	12.9 \pm 0.3	14.0 \pm 0.3	11.4 \pm 0.3	20.1 \pm 0.3	26.2 \pm 0.3	56.2 \pm 0.3
	10	2	14.6 \pm 0.5	17.3 \pm 0.3	7.1 \pm 0.2	25.2 \pm 0.3	32.3 \pm 0.3	27.8 \pm 0.3	31.5 \pm 0.2	29.7 \pm 0.3	45.1 \pm 0.1	42.4 \pm 0.2	
	50	10	30.0 \pm 0.4	33.7 \pm 0.5	30.5 \pm 0.3	-	42.8 \pm 0.4	37.9 \pm 0.3	42.9 \pm 0.2	43.6 \pm 0.4	50.0 \pm 0.2	50.9 \pm 0.7	

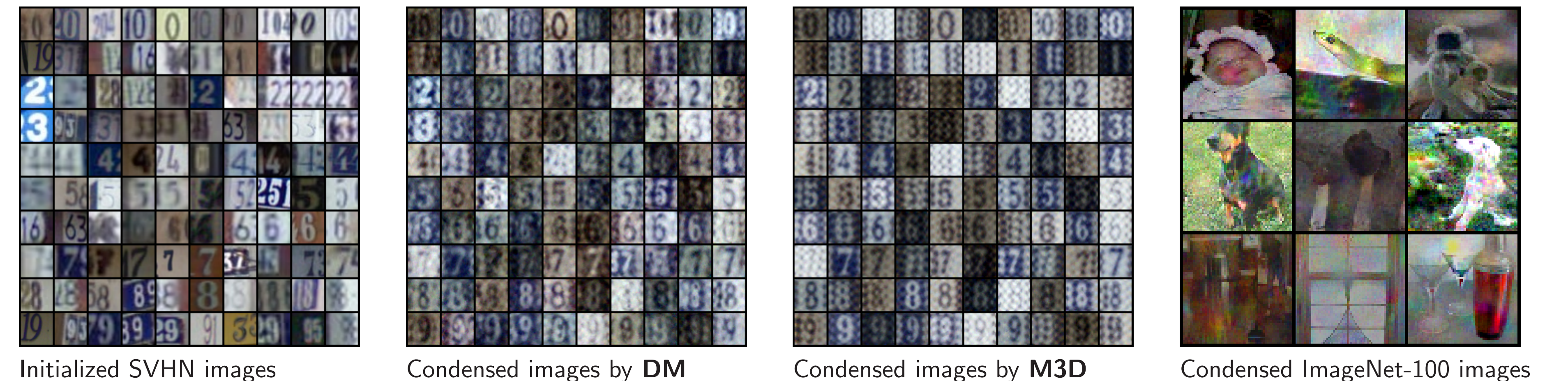
Ablation studies:



Experiments across varying training steps:



Visualization Results:



Initialized SVHN images

Condensed images by DM

Condensed images by M3D

Condensed ImageNet-100 images