
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

Abstract

Efficient training of deep neural networks is an increasingly important problem in the era of sophisticated architectures and large-scale datasets. This paper proposes a training set synthesis technique, called *Dataset Condensation*, that learns to produce a small set of informative samples for training deep neural networks from scratch in a small fraction of the required computational cost on the original data while achieving comparable results. We rigorously evaluate its performance in several computer vision benchmarks and show that it significantly outperforms the state-of-the-art methods. Finally we show promising applications of our method in continual learning and domain adaptation.

1 Introduction

Training deep neural networks fast, accurately and efficiently is a long-standing goal in machine learning with many practical benefits. As the state-of-the-art in computer vision [1, 2], natural language processing [3] and speech recognition [4, 5] rely on more sophisticated deep networks and larger datasets, their training becomes more computationally expensive, memory and storage intensive and even demands for specialized equipment and infrastructure. Thus a widespread adoption of the current state-of-the-art to different data and problems, and development of the future state-of-the-art which typically require training multiple models for validating various decisions (*e.g.* architecture designs, hyperparameters, loss formulations) call for efficient training strategies.

A well studied problem that aims at reducing train time by decreasing the number of models to be trained is hyperparameter optimization. The standard solutions include discretizing the parameter search space via a grid search, random sampling [6], evolutionary optimization [7], modeling the validation loss as a function of hyperparameters which can be fit to the results of previously explored ones to estimate the optimum [8, 9] and utilizing reinforcement learning to maximize the expected accuracy of the selected hyperparameters [10]. However, these strategies are not directly applicable to evaluating more complex design decisions which cannot be considered as hyperparameters.

This paper focuses on an orthogonal direction that aims at *reducing training time for each individual model* by using a significantly smaller training set such that the model learned on this set performs as closely as possible to the model trained on the entire dataset (illustrated in Figure 1(a)). A classical technique to reducing a large set of high-dimensional points to a small set such that the reduced set approximately preserves a target property (*e.g.* width, diameter) of the original data is coresets selection [11, 12, 13]. Recent examples of this approach can be found in active learning [14] and continual learning [15, 16, 17, 18] where there is typically a fixed budget in labeling and storing training samples respectively. Sener and Savarese [14] propose an active learning method that picks data points to label from a large set by using a greedy K-center technique. Rebuffi *et al.* [15] and Castro *et al.* [17] use the herding selection strategy in [19] for each seen category by picking the closest samples to the category mean in a feature space to construct the rehearsal memory in continual learning. Aljundi *et al.* [18] cast sample selection as a constraint reduction problem and propose to select a diverse set of samples based on the similarity between their gradients. Toneva *et al.* [16] study

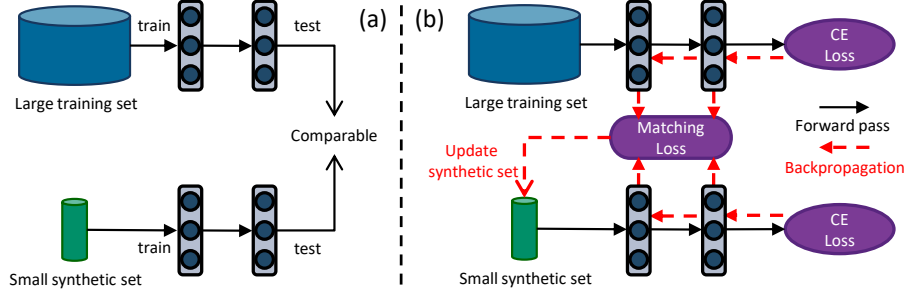


Figure 1: Dataset Condensation (left) aims to generate a small set of synthetic images that can match the performance of a network trained on a large image dataset. Our method (right) realizes this goal by learning a synthetic set such that a deep network trained on it and the large set produces similar gradients w.r.t. the parameters. The synthetic data can later be used to train a network from scratch in a fraction of the original computational load. CE denotes Cross-Entropy.

the relation between the sample selection and rate of misclassification during training, and show that the least forgotten (*i.e.* misclassified) samples, a significant fraction of training set, can be omitted while still maintaining good performance. However, these methods have two shortcomings: they rely on i) heuristics (*e.g.* picking cluster centers, diversity) that do not guarantee any optimum solution for the downstream task, ii) presence of representative samples. Motivated by these shortcomings, we focus on learning to *synthesize new samples* that are optimized to train neural networks for downstream tasks and not limited to individual samples in the data. We show in section 3 that our method produces a more compact set and achieves better generalization performance than the coresets when a deep neural network is trained on them.

Our method is inspired by a recent work, Dataset Distillation (DD) [20] that goes beyond finding coresets by learning to synthesize a small set of informative images from large training data. In particular, the authors model the network parameters as a function of the synthetic training data and learn them by minimizing the training loss over the original training data w.r.t. synthetic data. Like DD [20], our goal is also to achieve comparable generalization performance with a model trained on the synthesized images to a model trained on the original images (see Figure 1(a)). To this end, we propose a *Dataset Condensation* method to learn a small set of “condensed” synthetic samples such that a deep network trained on them obtains not only similar performance but also a close solution to a network trained on the large training data in the network parameter space. We formulate this goal as a minimization problem between two sets of gradients for the network parameters that are computed for a training loss over a large fixed training set and a learnable condensed set (see Figure 1(b)). Our method enables a more effective learning of synthetic images and we show that neural networks trained on the condensed images outperforms [20] with a wide margin in multiple computer vision benchmarks.

Our method is also related to knowledge distillation (KD) [21] techniques [22, 23, 24] that transfer the knowledge of an ensemble of models to a single one. Though our method can also be seen as KD from a teacher to a student, we distill knowledge of a large training set into a small synthetic set rather than the knowledge between models. Recent zero-shot KD methods [25, 26] aim to perform KD from a trained model in the absence of training data by generating synthetic data as the intermediate production to further use. Unlike them, our method does not require pretrained teacher models to provide the knowledge, *i.e.* to obtain the features and labels. Our method is also related to Generative Adversarial Networks [27, 28, 29] and Variational AutoEncoders [30] that synthesize high fidelity samples by capturing the data distribution. In contrast, our goal is to generate informative samples for training deep neural networks rather than to produce “real-looking” samples. Finally our method is loosely related to the methods [31, 32, 33] that recover images by projecting the feature activations back to the input pixel space [31], reconstruct the input image by matching the feature activations [32], recover private training images for given training gradients [33]. Our goal is however to synthesize a set of condensed training images not to recover the original training images.

In the remainder of this paper, we first review the problem of dataset condensation and introduce our method in section 2, present and analyze our results in several image recognition benchmarks in section 3.1, showcase applications in continual learning and few-shot domain adaptation in section 3.2, and conclude the paper with future directions in section 4.

2 Method

2.1 Dataset condensation

Suppose we are given a large dataset consisting of m pairs of a training image and its class label $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$ where $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^d$, $y \in \{1, \dots, C\}$, \mathcal{X} is a d -dimensional input space and C is the number of classes. We wish to learn a differentiable function ϕ (*i.e.* deep neural network) with parameters θ that correctly predicts labels of previously unseen images, *i.e.* $y = \phi_{\theta}(\mathbf{x})$. One can learn the parameters of this function by minimizing an empirical loss term over the training set:

$$\theta^{\mathcal{T}} = \arg \min_{\theta} \mathcal{L}^{\mathcal{T}}(\theta) \quad (1)$$

where $\mathcal{L}^{\mathcal{T}}(\theta) = \sum_{(\mathbf{x}, y) \in \mathcal{T}} \ell(\phi_{\theta}(\mathbf{x}), y)$, $\ell(\cdot, \cdot)$ is a task specific loss (*i.e.* cross-entropy) and $\theta^{\mathcal{T}}$ is the minimizer of $\mathcal{L}^{\mathcal{T}}$. The generalization performance of the obtained model $\phi_{\theta^{\mathcal{T}}}$ can be written as $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}}[\ell(\phi_{\theta^{\mathcal{T}}}(\mathbf{x}), y)]$ where $P_{\mathcal{D}}$ is the data distribution. Our goal is to generate a small set of condensed synthetic samples with their labels, $\mathcal{S} = \{(\mathbf{s}_i, y_i)\}_{i=1}^n$ where $\mathbf{s} \in \mathbb{R}^d$ and $y \in \mathcal{Y}$, $n \ll m$. Similar to eq. (1), once the condensed set is learned, one can train ϕ on them as follows

$$\theta^{\mathcal{S}} = \arg \min_{\theta} \mathcal{L}^{\mathcal{S}}(\theta) \quad (2)$$

where $\mathcal{L}^{\mathcal{S}}(\theta) = \sum_{(\mathbf{s}, y) \in \mathcal{S}} \ell(\phi_{\theta}(\mathbf{s}), y)$ and $\theta^{\mathcal{S}}$ is the minimizer of $\mathcal{L}^{\mathcal{S}}$. As the synthetic set \mathcal{S} is significantly smaller (*e.g.* 2-3 orders of magnitude), we expect the optimization in eq. (2) to be significantly faster than that in eq. (1). We also wish that the generalization performance of $\phi_{\theta^{\mathcal{S}}}$ is comparable to $\phi_{\theta^{\mathcal{T}}}$, *i.e.* $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}}[\ell(\phi_{\theta^{\mathcal{T}}}(\mathbf{x}), y)] \simeq \mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}}[\ell(\phi_{\theta^{\mathcal{S}}}(\mathbf{x}), y)]$.

Discussion. The problem of achieving comparable generalization performance by training on the condensed data can be formulated in different ways. One approach, which is proposed in [20], is to pose the parameters $\theta^{\mathcal{S}}$ as a function of the synthetic data \mathcal{S} :

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathcal{L}^{\mathcal{T}}(\theta^{\mathcal{S}}(\mathcal{S})) \quad \text{subject to} \quad \theta^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}^{\mathcal{S}}(\theta). \quad (3)$$

The method aims to find the optimum set of synthetic images \mathcal{S}^* such that the model $\phi_{\theta^{\mathcal{S}}}$ trained on them minimizes the empirical loss over the original training set. Optimizing eq. (3) involves a nested loop optimization and solving the inner loop for $\theta^{\mathcal{S}}(\mathcal{S})$ at each iteration to recover the gradients for \mathcal{S} which requires a computationally expensive procedure – unrolling the recursive computation graph for \mathcal{S} over multiple optimization steps for θ (see [34, 35]). Hence, it does not scale to large models and/or accurate inner-loop optimizers with many steps. Next we propose an alternative formulation for dataset condensation.

2.2 Dataset condensation with parameter matching

Here we aim to learn \mathcal{S} such that the model $\phi_{\theta^{\mathcal{S}}}$ trained on them achieves not only comparable generalization performance to $\phi_{\theta^{\mathcal{T}}}$ but also converges to a similar solution in the parameter space (*i.e.* $\theta^{\mathcal{S}} \approx \theta^{\mathcal{T}}$). Let ϕ_{θ} be a locally smooth function¹, similar weights ($\theta^{\mathcal{S}} \approx \theta^{\mathcal{T}}$) imply similar mappings in a local neighborhood and thus generalization performance, *i.e.* $\mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}}[\ell(\phi_{\theta^{\mathcal{T}}}(\mathbf{x}), y)] \simeq \mathbb{E}_{\mathbf{x} \sim P_{\mathcal{D}}}[\ell(\phi_{\theta^{\mathcal{S}}}(\mathbf{x}), y)]$. Now we can formulate this goal as

$$\min_{\mathcal{S}} D(\theta^{\mathcal{S}}, \theta^{\mathcal{T}}) \quad \text{subject to} \quad \theta^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}^{\mathcal{S}}(\theta) \quad (4)$$

where $\theta^{\mathcal{T}} = \arg \min_{\theta} \mathcal{L}^{\mathcal{T}}(\theta)$ and $D(\cdot, \cdot)$ is a distance function. In a deep neural network, $\theta^{\mathcal{T}}$ typically depends on its initial values θ_0 . However, the optimization in eq. (4) aims to obtain an optimum set of synthetic images only for one model $\phi_{\theta^{\mathcal{T}}}$ with the initialization θ_0 , while our actual goal is to generate samples that can work with a distribution of random initializations P_{θ_0} . Thus we modify eq. (4) as follows:

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}}[D(\theta^{\mathcal{S}}(\theta_0), \theta^{\mathcal{T}}(\theta_0))] \quad \text{subject to} \quad \theta^{\mathcal{S}}(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}^{\mathcal{S}}(\theta(\theta_0)) \quad (5)$$

¹Local smoothness is frequently used to obtain explicit first-order local approximations in deep networks (*e.g.* see [36, 37, 38]).

where $\theta^T = \arg \min_{\theta} \mathcal{L}^T(\theta(\theta_0))$. For brevity, we use only θ^S and θ^T to indicate $\theta^S(\theta_0)$ and $\theta^T(\theta_0)$ respectively in the next sections. The standard approach to solving eq. (5) employs implicit differentiation (see [35] for details), which involves solving an inner loop optimization for θ^S . As the inner loop optimization $\theta^S(\mathcal{S}) = \arg \min_{\theta} \mathcal{L}^S(\theta)$ can be computationally expensive in case of large-scale models, one can adopt the back-optimization approach in [35] which re-defines θ^S as the output of an incomplete optimization:

$$\theta^S(\mathcal{S}) = \text{opt-}\text{alg}_{\theta}(\mathcal{L}^S(\theta), \varsigma) \quad (6)$$

where $\text{opt-}\text{alg}$ is a specific optimization procedure with a fixed number of steps (ς).

In practice, θ^T for different initializations can be trained first in an offline stage and then used as the target parameter vector in eq. (5). However, there are two potential issues by learning to regress θ^T as the target vector. First the distance between θ^T and intermediate values of θ^S can be too big in the parameter space with multiple local minima traps along the path and thus it can be too challenging to reach. Second $\text{opt-}\text{alg}$ involves a limited number of optimization steps as a trade-off between speed and accuracy which may not be sufficient to take enough steps for reaching the optimal solution. These problems are similar to those of [20], as they both involve parameterizing θ^S with \mathcal{S} and θ_0 .

2.3 Dataset condensation with curriculum gradient matching

As such we propose a curriculum based approach to address the above mentioned challenges. The key idea is that we wish θ^S to be close to not only the final θ^T but also to follow a similar path to θ^T throughout the optimization. While this can restrict the optimization dynamics for θ , we argue that it also enables a more guided optimization and effective use of the incomplete optimizer. We can now rewrite eq. (5) as a sum of subproblems:

$$\begin{aligned} \min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=1}^T D(\theta_t^S, \theta_t^T) \right] \quad \text{subject to} \\ \theta_t^S(\mathcal{S}) = \text{opt-}\text{alg}_{\theta}(\mathcal{L}^S(\theta_{t-1}^S), \varsigma^S) \quad \text{and} \quad \theta_t^T = \text{opt-}\text{alg}_{\theta}(\mathcal{L}^T(\theta_{t-1}^T), \varsigma^T) \end{aligned} \quad (7)$$

where T is the number of iterations, ς^S and ς^T are the numbers of optimization steps for θ^S and θ^T respectively. In words, we wish to generate a set of condensed samples \mathcal{S} such that the network parameters trained on them (θ_t^S) are similar to the ones trained on the original training set (θ_t^T) at each iteration t . We find in our preliminary experiments that θ_t^S , which is parameterized with \mathcal{S} , can successfully track θ_t^T by updating \mathcal{S} and minimize $D(\theta_{t-1}^S, \theta_{t-1}^T)$ close to zero.

In the case of one step gradient descent optimization for $\text{opt-}\text{alg}$, we have the following update rule

$$\theta_t^S \leftarrow \theta_{t-1}^S - \eta_{\theta} \nabla_{\theta} \mathcal{L}^S(\theta_{t-1}^S) \quad \text{and} \quad \theta_t^T \leftarrow \theta_{t-1}^T - \eta_{\theta} \nabla_{\theta} \mathcal{L}^T(\theta_{t-1}^T), \quad (8)$$

where η_{θ} is the learning rate. Based on our observation ($D(\theta_{t-1}^S, \theta_{t-1}^T) \approx 0$), we simplify the formulation in eq. (7) by replacing θ_{t-1}^T with θ_{t-1}^S and use a single symbol θ to denote θ^S in the rest of the paper:

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[\sum_{t=1}^T D(\nabla_{\theta} \mathcal{L}^S(\theta_{t-1}), \nabla_{\theta} \mathcal{L}^T(\theta_{t-1})) \right]. \quad (9)$$

We now have a single deep network with parameters θ trained on the synthetic set \mathcal{S} which is optimized such that the distance between the gradients for the loss over the training samples \mathcal{L}^T w.r.t. θ and the gradients for the loss over the condensed samples \mathcal{L}^S w.r.t. θ is minimized. In words, our goal reduces to matching the gradients for the real and synthetic training loss w.r.t. θ via updating the condensed samples. This approximation has the key advantage over the one in eq. (5) and also [20] that it does not require the expensive unrolling of the recursive computation graph over the previous parameters $\{\theta_0, \dots, \theta_{t-1}\}$. The important consequence is that the optimization is significantly faster, memory efficient and thus scales up to the state-of-the-art deep neural networks (e.g. ResNet [2]).

Generating labeled synthetic samples. The set of synthetic data contains pairs of a sample and its label (s, y) that can be jointly learned by optimizing eq. (9) in theory. However, their joint optimization is challenging, as the content of the samples depend on their label and vice-versa. Thus in our experiments we learn to synthesize images for fixed labels, e.g. one synthetic image per class.

Algorithm. We depict the optimization details in Alg. 1. At the outer level, the optimization contains a loop over random weight initializations, as we want the condensed samples to train previously unseen models. Once θ , which denotes θ^S for brevity, is randomly initialized, we use ϕ_θ to first compute the loss over the training samples (\mathcal{L}^T) and the synthetic samples (\mathcal{L}^S) and their gradients w.r.t. θ , then optimize the synthetic samples to match these gradients $\nabla_\theta \mathcal{L}^S$ to $\nabla_\theta \mathcal{L}^T$ by applying ς_S gradient descent steps with learning rate η_S . We use the stochastic gradient descent optimization for both `opt- alg_θ` and `opt- alg_S` . Next we train θ on the synthetic samples by minimizing the loss \mathcal{L}^S with learning rate η_θ for ς_θ steps. Note that each real and synthetic batch sampled from \mathcal{T} and \mathcal{S} contains samples from a single class and the synthetic data for each class are individually updated at each iteration (t) for the following reasons: i) this reduces the memory use at train time, ii) imitating the mean gradients w.r.t. the data from single class is easier compared to that of multiple classes.

Algorithm 1: Dataset condensation with gradient matching

Input: Training set \mathcal{T}

- 1 **Required:** randomly initialized set of synthetic samples \mathcal{S} , probability distribution over randomly initialized weights P_{θ_0} , deep neural network ϕ_θ , number of outer-loop steps K , number of inner-loop steps T , number of steps for updating weights ς_θ and synthetic samples ς_S in each inner-loop step respectively, learning rates for updating weights η_θ and synthetic samples η_S .
- 2 **for** $k = 0, \dots, K - 1$ **do**
- 3 Initialize $\theta_0 \sim P_{\theta_0}$
- 4 **for** $t = 0, \dots, T - 1$ **do**
- 5 Sample a minibatch $B^T \sim \mathcal{T}$ and $B^S \sim \mathcal{S}$
- 6 Compute $\mathcal{L}^T = \sum_{(x,y) \in B^T} \ell(\phi_{\theta_t}(x), y)$ and $\mathcal{L}^S = \sum_{(s,y) \in B^S} \ell(\phi_{\theta_t}(s), y)$
- 7 Update synthetic samples: $\mathcal{S} \leftarrow \text{opt-}\text{alg}_S(D(\nabla_\theta \mathcal{L}^S(\theta_t), \nabla_\theta \mathcal{L}^T(\theta_t)), \varsigma_S, \eta_S)$
- 8 Update model parameters: $\theta_{t+1} \leftarrow \text{opt-}\text{alg}_\theta(\mathcal{L}^S(\theta_t), \varsigma_\theta, \eta_\theta)$

Output: \mathcal{S}

Gradient matching loss. The matching loss $D(\cdot, \cdot)$ in eq. (9) quantifies the distance between the gradients for \mathcal{L}^S and \mathcal{L}^T w.r.t. θ . When ϕ_θ is a multi-layered neural network, the gradients correspond to a set of learnable 2-dimensional ($\text{out} \times \text{in}$) weights at each fully connected (FC) and 4-dimensional ones ($\text{out} \times \text{in} \times \text{h} \times \text{w}$) at each convolutional layer where out , in , h , w are number of output and input channels, kernel height and width respectively. The matching loss can be decomposed into a sum of layerwise losses as $D(\nabla_\theta \mathcal{L}^S, \nabla_\theta \mathcal{L}^T) = \sum_{l=1}^L d(\nabla_{\theta^{(l)}} \mathcal{L}^S, \nabla_{\theta^{(l)}} \mathcal{L}^T)$ where l is the layer index, L is the number of layers with weights and

$$d(\mathbf{A}, \mathbf{B}) = \sum_{i=1}^{\text{out}} \left(1 - \frac{\mathbf{A}_i \cdot \mathbf{B}_i}{\|\mathbf{A}_i\| \|\mathbf{B}_i\|} \right) \quad (10)$$

where \mathbf{A}_i and \mathbf{B}_i are flattened vectors of gradients corresponding to each output node *i.e.* in dimensional for FC weights and $\text{in} \times \text{h} \times \text{w}$ dimensional for convolutional weights. In contrast to the previous works [39, 18, 33] that ignore the layer-wise structure by flattening tensors over all layers to one vector and then computing the distance between two vectors, we group them for each output node. We found that this is a better distance for gradient matching and enables using a single learning rate across all layers in our preliminary experiments.

3 Experiments

3.1 Dataset condensation

First we evaluate classification performance with the condensed images on four standard benchmark datasets: digit recognition on MNIST [40], SVHN [41] and object classification on FashionMNIST [42], CIFAR10 [43]. We test our method using six standard deep network architectures: MLP, ConvNet [44], LeNet [40], AlexNet [1], VGG-11 [45] and ResNet-18 [2]. MLP is a multilayer perceptron with two nonlinear hidden layers, each has 128 units. ConvNet is a commonly used modular architecture in few-shot learning [46, 47, 44] with D duplicate blocks, and each block has a convolutional layer with W (3×3) filters, a normalization layer N , an activation layer A and a pooling layer P , denoted as $[W, N, A, P] \times D$. The default ConvNet (unless specified otherwise)

	Img/Cls	Ratio %	Core-set Selection				Ours	Whole Dataset
			Random	Herdng	K-Center	Forgetting		
MNIST	1	0.017	64.9 \pm 3.5	89.2 \pm 1.6	89.3 \pm 1.5	35.5 \pm 5.6	91.7\pm0.5	99.6 \pm 0.0
	10	0.17	95.1 \pm 0.9	93.7 \pm 0.3	84.4 \pm 1.7	68.1 \pm 3.3	97.4\pm0.2	
FashionMNIST	1	0.017	51.4 \pm 3.8	67.0 \pm 1.9	66.9 \pm 1.8	42.0 \pm 5.5	70.5\pm0.6	93.5 \pm 0.1
	10	0.17	73.8 \pm 0.7	71.1 \pm 0.7	54.7 \pm 1.5	53.9 \pm 2.0	82.3\pm0.4	
SVHN	1	0.014	14.6 \pm 1.6	20.9 \pm 1.3	21.0 \pm 1.5	12.1 \pm 1.7	31.2\pm1.4	95.4 \pm 0.1
	10	0.14	35.1 \pm 4.1	50.5 \pm 3.3	14.0 \pm 1.3	16.8 \pm 1.2	76.1\pm0.6	
CIFAR10	1	0.02	14.4 \pm 2.0	21.5 \pm 1.2	21.5 \pm 1.3	13.5 \pm 1.2	28.3\pm0.5	84.8 \pm 0.1
	10	0.2	26.0 \pm 1.2	31.6 \pm 0.7	14.7 \pm 0.9	23.3 \pm 1.0	44.9\pm0.5	

Table 1: The performance comparison to coreset methods. This table shows the testing accuracies (%) of different methods on four datasets. ConvNet is used for training and testing. Img/Cls: image(s) per class, Ratio (%): the ratio of condensed images to whole training set.

includes 3 blocks, each with 128 filters, followed by InstanceNorm [48], ReLU and AvgPooling modules. The final block is followed by a linear classifier. More details about the datasets and networks can be found in the supplementary.

The pipeline for dataset condensation has two stages: learning the condensed images (denoted as Synth) and training classifiers from scratch on them (denoted as Cls). In the coreset baselines, the coreset is selected in the first stage. Note that the model architectures used in two stages might be different. We investigate two settings: 1 and 10 image/class learning, which means that the synthetic set or core-set contains 1 and 10 images per class respectively. Each method is run for 5 times, and 5 synthetic sets are generated in the first stage. Each generated synthetic set is evaluated on 20 randomly initialized models in the second stage, which equals to evaluating 100 models in the second stage. In all experiments, we report the mean and standard deviation of these 100 testing results.

Baselines. We compare our method to four core-set selection baselines (Random, Herding, K-Center and Forgetting) and also to dataset distillation (DD) [20]. In random, the training samples are randomly selected as the core-set. Herding baseline, which selects closest samples to the cluster center, is based on the method of [19] and used in [15, 17, 49, 50]. K-Center [51, 14] selects multiple center points such that the largest distance between a data point and its nearest center is minimized. For herding and K-Center, we use models trained on the whole dataset to extract features, compute l_2 distance to centers. Forgetting method [16] considers those training samples which are easy to forget during training as the important ones and uses forgetting statistics to measure the importance. The forgetting statistics are obtained during the training process [16].

Comparison to coreset methods. We first compare our method to the coreset baselines on MNIST, FashionMNIST, SVHN and CIFAR10 in Table 1 using the default ConvNet in terms of classification accuracy. Whole dataset indicates a model trained on the whole training set which serves as an approximate upper-bound performance. First we observe that our method outperforms all the baselines significantly and achieves a comparable result (97.4%) in case of 10 images per class to the upper bound (99.6%) in MNIST which uses two orders of magnitude more training images per class (6000). We also obtain promising results in FashionMNIST, however, the gap between our method and upper bound is bigger in SVHN and CIFAR10 which contain more diverse images with varying foregrounds and backgrounds. We also observe that, (i) the random selection baseline is competitive to other coreset methods in case of 10 images per class and (ii) herding method is on average the best coreset technique. We visualize the synthetic images produced by our method under 1 image/class setting in Figure 2. Interestingly they are interpretable and look like “prototypes” of each class.

Comparison to Dataset Distillation [20]. Unlike the setting in Table 1, DD [20] reports results only for 10 images per class on MNIST and CIFAR10 over LeNet and AlexCifarNet (a customized AlexNet). For a fair comparison, we strictly follow the experimental setting in [20], use the same architectures and report our and their original results in Table 3. Our method achieves significantly better performance than DD on both MNIST and CIFAR10. Especially, on MNIST, we obtain 5% better accuracy with only 1 synthetic sample per class than DD with 10 per class. In addition, our method obtains consistent results over multiple runs with a standard deviation of only 0.6% on MNIST, while DD’s performance significantly vary over different runs (8.1%). Finally our method trains 2 times faster than DD and requires 50% less memory on CIFAR10 experiments. More detailed runtime and qualitative comparison can be found in the supplementary.



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

Synth/Cls	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	91.6\pm0.5	85.3 \pm 1.8	85.1 \pm 3.0	83.4\pm1.8	90.0\pm0.8
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	73.6\pm1.2	91.6\pm0.5	86.4\pm1.5	85.4\pm1.9	83.4\pm2.4	89.4 \pm 0.9

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

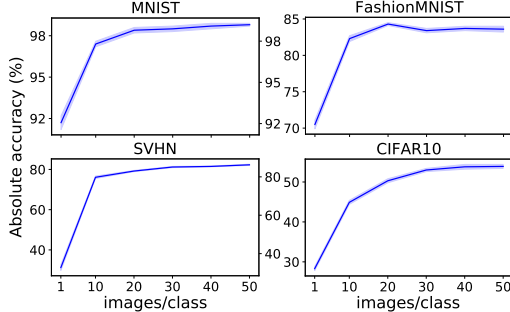


Figure 3: Absolute and relative testing accuracies for varying the number of condensed images/class for MNIST, FashionMNIST, SVHN and CIFAR10. The relative accuracy means the ratio compared to its upper-bound, *i.e.* training with the whole dataset.

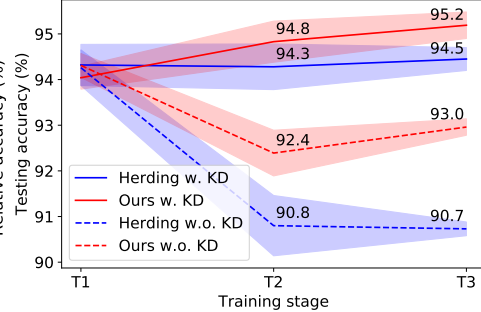


Figure 4: Continual learning performance in accuracy (%). Herding denotes the original E2E [17]. T1, T2, T3 are three learning stages. The performance at each stage is the mean testing accuracy on all learned tasks.

Computational cost. A key advantage of our method is the significant reduction in the size of training set and thus in the training time. Once the condensed images are learned, one can train a network (from scratch) on them only in a fraction of the training time taken by the original data. For instance, the full-dataset training times for learning a ConvNet (from scratch) over MNIST, FashionMNIST, SVHN and CIFAR10 on a NVIDIA GTX1080-Ti GPU are respectively 773, 776, 931 and 644 seconds (all above 10 minutes). However, the training times for 1 and 10 condensed images/class are only 0.6 and 1.3 seconds respectively for all datasets which is a remarkable reduction of 2-3 orders of magnitude. More details can be found in the supplementary.

Cross-architecture generalization. Another key advantage of our method is that the condensed images learned using one architecture can be used to train another unseen one. Here we learn to generate 1 condensed image per class for MNIST over a diverse set of networks including MLP, ConvNet [44], LeNet [40], AlexNet [1], VGG-11 [45] and ResNet-18 [2] (see Table 2). Once the condensed sets are synthesized, we train every network on all the sets separately from scratch and evaluate their cross architecture performance in terms of classification accuracy on the MNIST test set. Table 2 shows that the condensed images, especially the ones that are trained with convolutional networks, perform well and are thus architecture generic. MLP generated images do not work well for training convolutional architectures which is possibly due to the mismatch between translation invariance properties of MLP and convolutional networks. Interestingly, MLP achieves better performance with convolutional network generated images than the MLP generated ones. The best results are obtained in most cases with ResNet generated images and ConvNet or ResNet as classifiers which is inline with the performances of the architectures when trained on the original dataset.

Number of condensed images. We also study the relation between the number of condensed images per class and test performance of a ConvNet trained on them for MNIST, FashionMNIST, SVHN and CIFAR10 in Figure 3 in absolute and relative terms – normalized by its upper-bound. While increasing the number of condensed images improves the accuracies in all benchmarks and further closes the gap with the upper-bound performance especially in MNIST and FashionMNIST, the performance saturates around 50 images/class. We plan to investigate effective regularization strategies on the condensed images to further improve the performance in the future.

Dataset	Img/Cls	DD [20]	Ours	Whole Dataset		I/C	Random	Herding	K-Center	Forgetting	Ours	No Adapt
MNIST	1	-	85.0 \pm 1.6	99.5 \pm 0.0	S \rightarrow M	1	75.7 \pm 4.5	93.5 \pm 0.8	93.1 \pm 1.1	41.5 \pm 6.4	94.4\pm0.4	76.5 \pm 1.7
	10	79.5 \pm 8.1	93.9\pm0.6			10	97.0 \pm 0.5	95.8 \pm 0.4	90.8 \pm 0.7	74.8 \pm 1.8	97.7\pm0.2	
CIFAR10	1	-	24.2 \pm 0.9	83.1 \pm 0.2	M \rightarrow S	1	17.2 \pm 2.2	25.0 \pm 3.2	25.0 \pm 3.2	12.8 \pm 1.4	35.7\pm2.1	31.0 \pm 1.2
	10	36.8 \pm 1.2	39.1\pm1.2			10	56.1 \pm 1.8	61.7 \pm 2.1	19.0 \pm 1.6	32.3 \pm 2.2	74.8\pm0.8	

Table 3: Comparison to DD [20] in terms of classification accuracy (%).

Table 4: Few-shot domain adaption in terms of classification accuracy (%). I/C: image(s) per class. M: MNIST. S: SVHN.

Activation, normalization & pooling. We also study the effect of various activation (sigmoid, ReLU [52, 53], leaky ReLU [54]), pooling (max, average) and normalization functions (batch [55], group [56], layer [57], instance norm [48]) and have the following observations: i) leaky ReLU over ReLU and average pooling over max pooling enable learning better condensed images, as they allow for denser gradient flow; ii) instance normalization obtains better classification performance than its alternatives when used in the networks that are trained on a small set of condensed images. We refer to the supplementary for detailed results and discussion.

3.2 Applications

Continual Learning First we apply our method to a continual-learning scenario [15, 17] where new tasks are learned incrementally and the goal is to preserve the performance on the old tasks while learning the new ones. In particular, we build our model on E2E method in [17] that uses a limited budget rehearsal memory (we consider 10 images/class here) to keep representative samples from the old tasks and knowledge distillation (KD) to regularize the network’s output w.r.t. to previous predictions. We replace its sample selection mechanism (herding) with ours such that a set of condensed images are generated and stored in the memory, keep the rest of the model same and evaluate this model on the task-incremental learning problem on the digit recognition datasets, SVHN [41], MNIST [40] and USPS [58] in the same order. MNIST and USPS images are reshaped to 32×32 RGB images. We compare our method to E2E [17], depicted as herding in Figure 4, with and without KD regularization. The experiment contains 3 incremental training stages (SVHN \rightarrow MNIST \rightarrow USPS) and testing accuracies are computed by averaging over the test sets of the previous and current tasks after each stage. The desired outcome is to obtain high mean classification accuracy at T3. We show that our method achieves better performance than E2E for both settings and the condensed images are more informative than the ones sampled by herding. In particular, our method outperforms E2E by a large margin (2.3% at T3) when KD is not employed.

Few-shot domain adaptation. Next we apply our method to a supervised domain adaptation scenario [59] where the goal is to efficiently transfer a model that is trained on a large labeled source domain to a small target domain with few selected or synthesized samples. To this end, we use a simple but effective transfer learning technique, *finetuning*, that trains the pretrained network on the few given samples. We report results on two domain adaptation scenarios: from MNIST (RGB) to SVHN, and from SVHN to MNIST (RGB). Table 4 depicts the adaptation results using ConvNet architecture for all methods. We compare our method to the selected coreset methods and also to the pretrained network without finetuning (No adaptation). We show that our method synthesizes more informative images and thus a better learning of the target task. We also run an additional experiment with LeNet architecture for comparing against DD [20] which reports results only for SVHN to MNIST for 10 images per class. Our method obtains $93.7 \pm 0.8\%$ and outperforms DD ($85.2 \pm 4.7\%$) in this setting though DD requires many pretrained models to synthesize images while ours employ randomly initialized ones.

4 Conclusion

In this paper, we have proposed a dataset condensation method that learns to synthesize a small set of informative images. These images can be used to efficiently train neural networks 2-3 orders of magnitude faster with a modest drop in performance. Our method outperforms the state-of-the-art significantly in multiple benchmarks and also obtains promising results in continual learning and few-shot adaptation. As future work, we plan to explore the use of condensed images in more diverse and thus challenging datasets such as ImageNet [60] that contain higher resolution images with larger variations in pose, appearance of objects, background.

Acknowledgment. We thank Iain Murray and Oisín Mac Aodha for their valuable feedback.

A Implementation details

In this part, we explain the implementation details for the dataset condensation, continual learning and few-shot domain adaptation experiments.

Dataset condensation. The presented experiments involve tuning of six hyperparameters – the number of outer-loop K and inner-loop steps T , learning rates η_S and number of optimization steps ς_S for the condensed samples, learning rates η_θ and number of optimization steps ς_θ for the model weights. In all experiments, we set $K = 1000$, $\eta_S = 0.1$, $\eta_\theta = 0.01$, $\varsigma_S = 1$ and employ Stochastic Gradient Descent (SGD) as the optimizer. The only exception is that we set η_S to 0.01 for synthesizing data with MLP in cross-architecture experiments (Table 2), as MLP requires a slightly different treatment. Note that while K is the maximum number of outer-loop steps, the optimization early-stops automatically if it converges before K steps. For the remaining hyperparameters, we use two different sets of for each 1 and 10 image/class setting. We set $T = 1$, $\varsigma_\theta = 1$ for 1 image/class and $T = 10$, $\varsigma_\theta = 50$ for 10 images/class. Note that when $T = 1$, it is not required to update the model parameters (see the Step 8 in Algorithm 1), as this model is not further used. For those experiments where more than 10 images/class are synthesized, we set T to be the same number as the synthetic images per class and $\varsigma_\theta = 500/T$, *e.g.* $T = 20$, $\varsigma_\theta = 25$ for 20 images/class learning.

The minibatches (Step 5 in Algorithm 1) that are sampled at each inner iteration contain samples from a class at a time. Specifically, we randomly sample 256 real images of a class (c) as a batch to calculate the mean gradient $\nabla_{\theta} \mathcal{L}^{\mathcal{T}_c}(\theta_t)$ and match it with the mean gradient $\nabla_{\theta} \mathcal{L}^{\mathcal{S}_c}(\theta_t)$ that is averaged over the condensed samples with the corresponding class. We train the condensed samples class-by-class in every inner-loop.

In all experiments, we use the standard train/test splits of the datasets – the train/test statistics are shown in Table T1. We apply data augmentation (crop, scale and rotate) only for experiments (coreset methods and ours) on the MNIST dataset. The only exception is that we also use data augmentation when compared to DD [20] on CIFAR10 with AlexCifarNet, as data augmentation is also used in [20].

	MNIST	FashionMNIST	SVHN	CIFAR10	USPS
Train	60,000	60,000	73,257	50,000	7,291
Test	10,000	10,000	26,032	10,000	2,007

Table T1: Train/test statistics for MNIST, FashionMNIST, SVHN, CIFAR10 and USPS datasets.

In the first stage – while training the condensed images –, we use Batch Normalization in the VGG and ResNet networks. For reliable estimation of the running mean and variance, we sample many real training data to estimate the running mean and variance and then freeze them ahead of Step 7. In the second stage – while training a deep network on the condensed set –, we replace Batch Normalization layers with Instance Normalization in VGG and ResNet, due to the fact that the batch statistics are not reliable when training networks with few condensed images. Another minor modification that we apply to the standard network ResNet architecture in the first stage is replacing the strided convolutions where $stride = 2$ with convolutional layers where $stride = 1$ coupled with an average pooling layer. We observe that this change enables more detailed (per pixel) gradients w.r.t. the condensed images and leads to better condensed images.

Continual learning. In this experiment, we focus on a task-incremental learning on SVHN, MNIST and USPS with the given order. The three tasks share the same label space, however have significantly different image statistics. The images of the three datasets are reshaped to 32×32 RGB size for standardization. We use the standard splits for training sets and randomly sample 2,000 test images for each datasets to obtain a balanced evaluation over three datasets. Thus each model is tested on a growing test set with 2,000, 4,000 and 6,000 images at the three stages respectively. We use the default ConvNet in this experiment and set the weight of distillation loss to 1.0 and the temperature to 2. We run 5,000 and 500 iterations for training and balanced finetuning as in [17] with the learning rates 0.01 and 0.001 respectively. We run 5 experiments and report the mean and standard variance in Figure 4.

Synth\Cls	Sigmoid	ReLu	LeakyReLu
Sigmoid	86.7 \pm 0.7	91.2 \pm 0.6	91.2 \pm 0.6
ReLu	86.1 \pm 0.9	91.7 \pm 0.5	91.7 \pm 0.5
LeakyReLu	86.3 \pm 0.9	91.7 \pm 0.5	91.7 \pm 0.4

Table T3: Cross-activation experiments in accuracy (%) for 1 condensed image/class in MNIST.

Synth\Cls	None	MaxPooling	AvgPooling
None	78.7 \pm 3.0	80.8 \pm 3.5	88.3 \pm 1.0
MaxPooling	81.2 \pm 2.8	89.5 \pm 1.1	91.1 \pm 0.6
Avgpooling	81.8 \pm 2.9	90.2 \pm 0.8	91.7 \pm 0.5

Table T4: Cross-pooling experiments in accuracy (%) for 1 condensed image/class in MNIST.

Few-shot domain adaptation. Here we reshape MNIST data to 32×32 RGB images in this experiment so that both MNIST and SVHN images have the same dimensionality and the same architecture (ConvNet) can be used on both images. The learning rate for finetuning is initialized with 0.01 and then lowered to 0.001 after 300 iterations. We run each experiment for 5 times with different random initializations and report their mean and standard variance.

B Further analysis

Next we provide additional results on the computational cost of training deep networks on various condensed sets of images, ablative studies over various deep network layers including activation, pooling and normalization functions and also over depth and width of deep network architecture, an additional qualitative analysis on the learned condensed images.

Computational cost. In Table T2 we report the number of training images and training times for training a ConvNet from scratch for the standard training on the whole dataset, and for our method on 1 and 10 condensed image(s)/class in MNIST, FashionMNIST, SVHN and CIFAR10. All the experiments are run for 100 epochs – which is just enough for convergence in all the settings – on a NVIDIA GTX 1080-Ti GPU and the training times are averaged over 10 experiments. No data augmentation is applied. We see in Table T2 that training on the condensed sets are about 2 to 3 orders of magnitude faster than the standard training on the whole data.

		MNIST	FashionMNIST	SVHN	CIFAR10
Images	Whose Dataset	60000	60000	73257	50000
	Ours 1 img/clc	10	10	10	10
	Ours 10 img/clc	100	100	100	100
Time (s)	Whose Dataset	773	776	931	644
	Ours 1 img/clc	0.6	0.6	0.6	0.6
	Ours 10 img/clc	1.3	1.3	1.3	1.3

Table T2: Training time for a randomly initialized ConvNet on whole dataset and our condensed set.

Ablation study on activation functions. Here we study the use of three activation functions – Sigmoid, ReLU, LeakyReLU (negative slope is set to 0.01) – in two stages, when training condensed images (denoted as Synth) and when training a ConvNet from scratch on the learned condensed images (denoted as Cls). The experiments are conducted in MNIST dataset for 1 condensed image/class setting. Table T3 shows that all three activation functions are good for the first stage while generating good condensed images, however, Sigmoid performs poor in the second stage while learning a classifier on the condensed images – its testing accuracies are lower than ReLu and LeakyReLU by around 5%. This suggests that ReLU can provide sufficiently informative gradients for learning condensed images, though the gradient of ReLU w.r.t. to its input is typically sparse.

Ablation study on pooling functions. Next we investigate the performance of two pooling functions – average pooling and max pooling – also no pooling for 1 image/class dataset condensation with ConvNet architecture in MNIST in terms of classification accuracy. Table T4 shows that max and average pooling both perform significantly better than no pooling (None) when they are used in the second stage. When the condensed samples are trained and tested on models with average pooling, the best testing accuracy ($91.7 \pm 0.5\%$) is obtained, possibly, because average pooling provides more informative and smooth gradients for the whole image rather than only for its discriminative parts.

Ablation study on normalization functions. Next we study the performance of four normalization options – No normalization, Batch [55], Layer [57], Instance [48] and Group Normalization [56] (number of groups is set to be four) – for 1 image/class dataset condensation with ConvNet architecture in MNIST classification accuracy. Table T5 shows that the normalization layer has little influence for

Synth \ Cls	Cls					
	None	BatchNorm	LayerNorm	InstanceNorm	GroupNorm	
None	79.0 \pm 2.2	80.8 \pm 2.0	85.8 \pm 1.7	90.7 \pm 0.7	85.9 \pm 1.7	
BatchNorm	78.6 \pm 2.1	80.7 \pm 1.8	85.7 \pm 1.6	90.9 \pm 0.6	85.9 \pm 1.5	
LayerNorm	81.2 \pm 1.8	78.6 \pm 3.0	87.4 \pm 1.3	90.7 \pm 0.7	87.3 \pm 1.4	
InstanceNorm	72.9 \pm 7.1	56.7 \pm 6.5	82.7 \pm 5.3	91.7 \pm 0.5	84.3 \pm 4.2	
GroupNorm	79.5 \pm 2.1	81.8 \pm 2.3	87.3 \pm 1.2	91.6 \pm 0.5	87.2 \pm 1.2	

Table T5: Cross-normalization experiments in accuracy (%) for 1 condensed image/class in MNIST.

Synth\Cls	1	2	3	4
1	61.3 \pm 3.5	78.2 \pm 3.0	77.1 \pm 4.0	76.4 \pm 3.5
2	78.3 \pm 2.3	89.0 \pm 0.8	91.0 \pm 0.6	89.4 \pm 0.8
3	81.6 \pm 1.5	89.8 \pm 0.8	91.7 \pm 0.5	90.4 \pm 0.6
4	82.5 \pm 1.3	89.9 \pm 0.8	91.9 \pm 0.5	90.6 \pm 0.4

Table T6: Cross-depth performance in accuracy (%) for 1 condensed image/class in MNIST.

Synth\Cls	32	64	128	256
32	90.6 \pm 0.8	91.4 \pm 0.5	91.5 \pm 0.5	91.3 \pm 0.6
64	91.0 \pm 0.8	91.6 \pm 0.6	91.8 \pm 0.5	91.4 \pm 0.6
128	90.8 \pm 0.7	91.5 \pm 0.6	91.7 \pm 0.5	91.2 \pm 0.7
256	91.0 \pm 0.7	91.6 \pm 0.6	91.7 \pm 0.5	91.4 \pm 0.5

Table T7: Cross-width performance in accuracy (%) for 1 condensed image/class in MNIST.

learning the condensed set, while the choice of normalization layer is important for training networks on the condensed set. LayerNorm and GroupNorm have similar performance, and InstanceNorm is the best choice for training a model on condensed images. BatchNorm obtains lower performance which is similar to None (no normalization), as it is known to perform poorly when training models on few condensed samples as also observed in [56]. Note that Batch Normalization does not allow for a stable training in the first stage (Synth); thus we replace its running mean and variance for each batch with those of randomly sampled real training images.

Ablation study on network depth and width. Here we study the effect of network depth and width for 1 image/class dataset condensation with ConvNet architecture in MNIST in terms of classification accuracy. To this end we conduct multiple experiments by varying the depth and width of the networks that are trained to synthesize condensed images and that are trained to classify testing data in ConvNet architecture and report the results in Table T6 and Table T7. In Table T6, we observe that deeper ConvNets with more blocks generate better condensed images that results in better classification performance when a network is trained on them, while ConvNet with 3 blocks performs best as classifier. Interestingly, Table T7 shows that the best results are obtained with the classifier that has 128 filters at each block, while network width (number of filters at each block) in generation has little overall impact on the final classification performance.

Further qualitative analysis We first depict the condensed images that are learned on MNIST, FashionMNIST, SVHN and CIFAR10 datasets in one experiment by using the default ConvNet in 10 images/class setting in Figure F5. It is interesting that the 10 images/class results in Figure F5 are diverse which cover the main variations, while the condensed images for 1 image/class setting (see Figure 2) look like the “prototype” of each class. For example, in Figure F5 (a), the ten images of “four” indicate ten different styles. The ten “bag” images in Figure F5 (b) are significantly different from each other, similarly “wallet” (1st row), “shopping bag” (3rd row), “handbag” (8th row) and “schoolbag” (10th row). Figure F5 (c) also shows the diverse house numbers with different shapes, colors and shadows. Besides, different poses of a “horse” have been learned in Figure F5 (d).

C Further comparison to DD [20]

Next we compare our method to DD [20] first quantitatively in terms of cross-architecture generalization, then qualitatively in terms of synthetic image quality, and finally in terms of computational load for training synthetic images. Note that we use the original source code to obtain the results for DD that is provided by the authors of DD in the experiments.

Generalization ability comparison. Here we compare the generalization ability across different deep network architectures to DD. To this end, we use the synthesized 10 images/class data learned with LeNet on MNIST to train MLP, ConvNet, LeNet, AlexNet, VGG11 and ResNet18 and report the results in Table T8. We see that that the condensed set produced by our method achieves good classification performances with all architectures, while the synthetic set by DD perform poorly when used to trained some architectures, *e.g.* AlexNet, VGG and ResNet. Note that DD generates learning rates to be used in every training step in addition to the synthetic data. This is in contrast to our

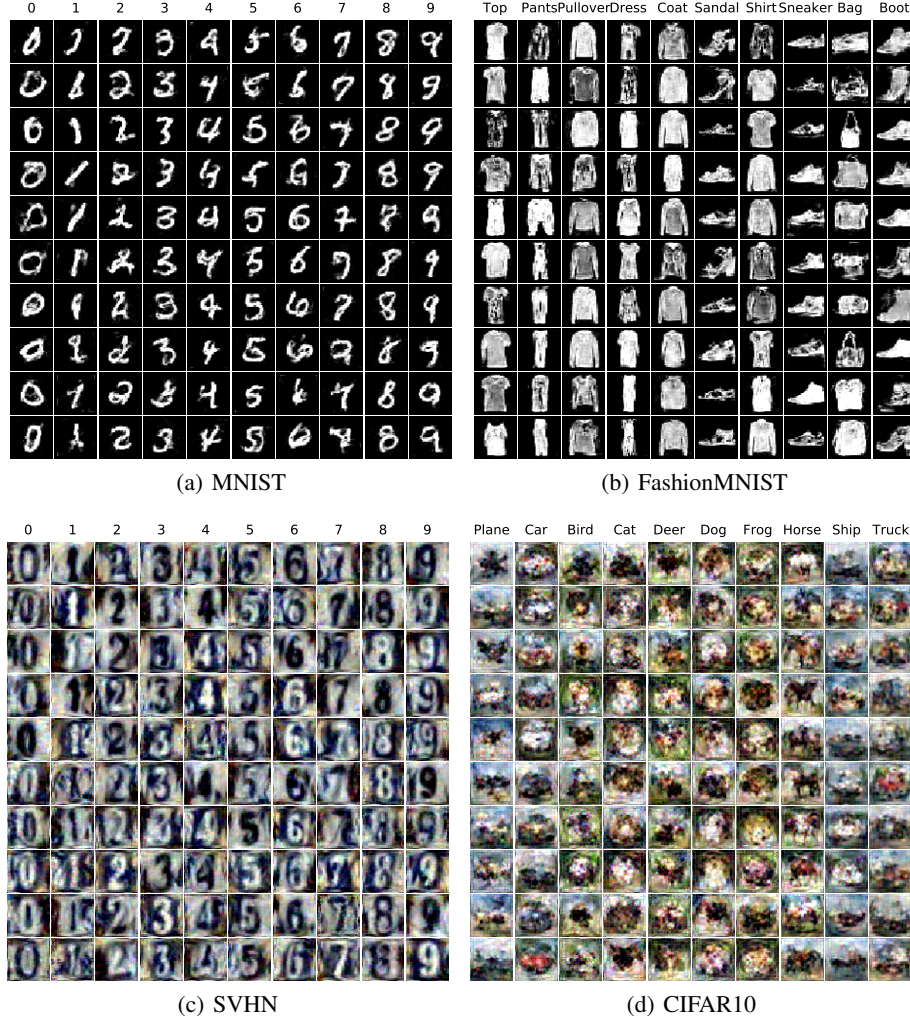


Figure F5: The synthetic images for MNIST, FashionMNIST, SVHN and CIFAR10 produced by our method with ConvNet under 10 images/class setting.

method which does not learn learning rates for specific training steps. Although the tied learning rates improve the performance of DD while training a network on the synthetic data, our method outperforms it in this experiment.

Qualitative comparison. We also provide a qualitative comparison to DD in terms of image quality in Figure F6. Note that both of the synthetic sets are trained with LeNet on MNIST and AlexCifarNet on CIFAR10. Our method produces more interpretable and realistic images than DD, although it is not our goal. The MNIST images produced by DD are noisy, and the CIFAR10 images produced by DD do not show any clear structure of the corresponding class. In contrast, the MNIST and CIFAR10 images produced by our method are both visually meaningful and diverse.

Training memory and time. One advantage of our method is that we decouple the model weights from its previous states in training, while DD requires to maintain the recursive computation graph which is not scalable to large models and inner-loop optimizers with many steps. Hence, our method requires less training time and memory cost. We compare the training time and memory cost required by DD and our method with one NVIDIA GTX1080-Ti GPU. Table T9 shows that our method requires significantly less memory and training time than DD and provides an approximation reduction of 20% and 45% in memory and 400% and 200% in train time to learn MNIST and CIFAR10 datasets respectively.

Method	MLP	ConvNet	LeNet	AlexNet	VGG	ResNet
DD	72.7 \pm 2.8	77.6 \pm 2.9	79.5 \pm 8.1	51.3 \pm 19.9	11.4 \pm 2.6	63.6 \pm 12.7
Ours	83.0\pm2.5	92.9\pm0.5	93.9\pm0.6	90.6\pm1.9	92.9\pm0.5	94.5\pm0.4

Table T8: Generalization ability comparison to DD. The 10 condensed images per class are trained with LeNet, and tested on various architectures. It shows that condensed images generated by our method have better generalization ability.

Method	Dataset	Architecture	Memory (MB)	Time (min)	Test Acc.
DD	MNIST	LeNet	785	160	79.5 \pm 8.1
Ours	MNIST	LeNet	653	46	93.9 \pm 0.6
DD	CIFAR10	AlexCifarNet	3211	214	36.8 \pm 1.2
Ours	CIFAR10	AlexCifarNet	1445	105	39.1 \pm 1.2

Table T9: Time and memory use for training DD and our method in 10 images/class setting.

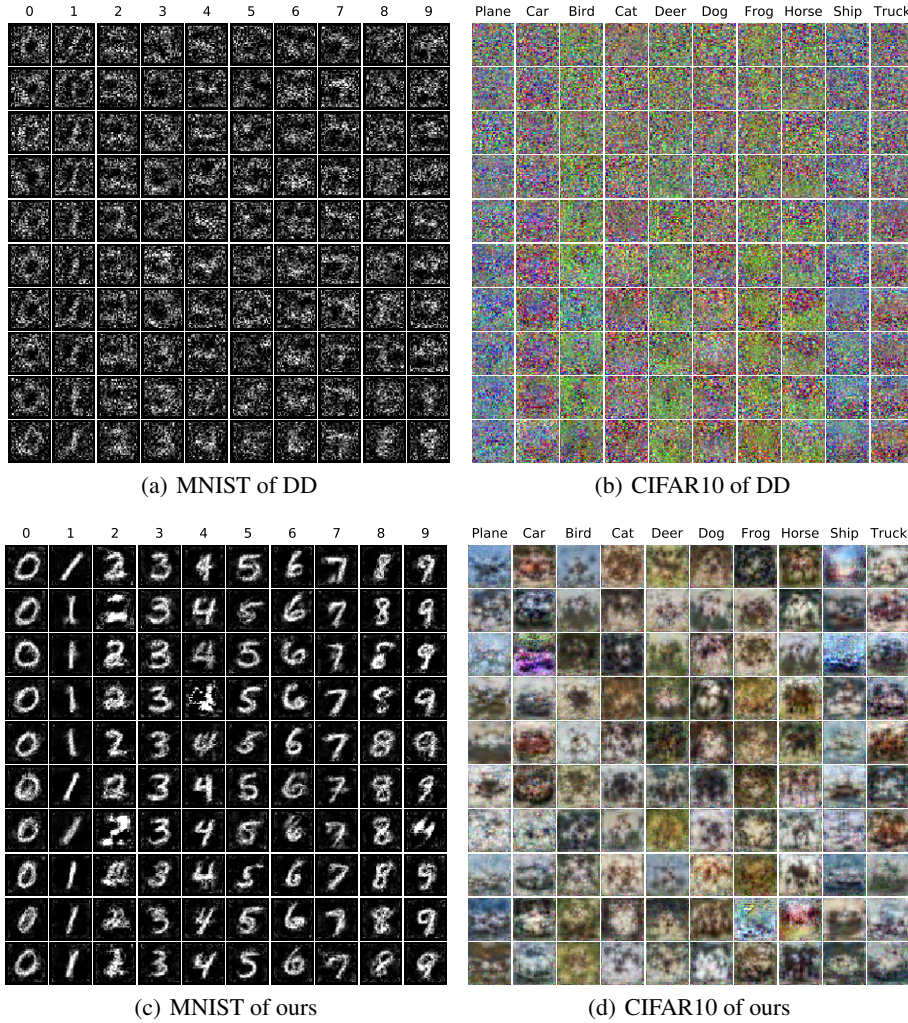


Figure F6: Qualitative comparison between the condensed images produced by DD and ours under 10 images/class setting. LeNet and AlexCifarNet are utilized for MNIST and CIFAR10 respectively.

References

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [4] Dario Amodei, Sundaram Ananthanarayanan, Rishita Anubhai, Jingliang Bai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro, Qiang Cheng, Guoliang Chen, et al. Deep speech 2: End-to-end speech recognition in english and mandarin. In *International conference on machine learning*, pages 173–182, 2016.
- [5] George Saon, Gakuto Kurata, Tom Sercu, Kartik Audhkhasi, Samuel Thomas, Dimitrios Dimitriadis, Xiaodong Cui, Bhuvana Ramabhadran, Michael Picheny, Lynn-Li Lim, et al. English conversational telephone speech recognition by humans and machines. *arXiv preprint arXiv:1703.02136*, 2017.
- [6] James Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *Journal of machine learning research*, 13(Feb):281–305, 2012.
- [7] Nikolaus Hansen, Sibylle D Müller, and Petros Koumoutsakos. Reducing the time complexity of the derandomized evolution strategy with covariance matrix adaptation (cma-es). *Evolutionary computation*, 11(1):1–18, 2003.
- [8] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *International conference on learning and intelligent optimization*, pages 507–523. Springer, 2011.
- [9] Jasper Snoek, Hugo Larochelle, and Ryan P Adams. Practical bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems*, pages 2951–2959, 2012.
- [10] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. *arXiv preprint arXiv:1611.01578*, 2016.
- [11] Pankaj K Agarwal, Sarel Har-Peled, and Kasturi R Varadarajan. Approximating extent measures of points. *Journal of the ACM (JACM)*, 51(4):606–635, 2004.
- [12] Sarel Har-Peled and Soham Mazumdar. On coresets for k-means and k-median clustering. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 291–300, 2004.
- [13] Dan Feldman, Melanie Schmidt, and Christian Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca and projective clustering. In *Proceedings of the twenty-fourth annual ACM-SIAM symposium on Discrete algorithms*, pages 1434–1453. SIAM, 2013.
- [14] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. *ICLR*, 2018.
- [15] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- [16] Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *ICLR*, 2019.
- [17] Francisco M Castro, Manuel J Marín-Jiménez, Nicolás Guil, Cordelia Schmid, and Karteek Alahari. End-to-end incremental learning. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 233–248, 2018.
- [18] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In *Advances in Neural Information Processing Systems*, pages 11816–11825, 2019.

- [19] Max Welling. Herding dynamical weights to learn. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 1121–1128. ACM, 2009.
- [20] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation. *arXiv preprint arXiv:1811.10959*, 2018.
- [21] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [22] Cristian Buciluă, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 535–541, 2006.
- [23] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In *Advances in neural information processing systems*, pages 2654–2662, 2014.
- [24] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [25] Raphael Gontijo Lopes, Stefano Fenu, and Thad Starner. Data-free knowledge distillation for deep neural networks. In *LLD Workshop at Neural Information Processing Systems (NIPS)*, 2017.
- [26] Gaurav Kumar Nayak, Konda Reddy Mopuri, Vaisakh Shaj, Venkatesh Babu Radhakrishnan, and Anirban Chakraborty. Zero-shot knowledge distillation in deep networks. In *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- [27] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- [28] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.
- [29] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- [30] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [31] Matthew D Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *European conference on computer vision*, pages 818–833. Springer, 2014.
- [32] Aravindh Mahendran and Andrea Vedaldi. Understanding deep image representations by inverting them. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5188–5196, 2015.
- [33] Ligeng Zhu, Zhijian Liu, and Song Han. Deep leakage from gradients. In *Advances in Neural Information Processing Systems*, pages 14747–14756, 2019.
- [34] Kegan GG Samuel and Marshall F Tappen. Learning optimized map estimates in continuously-valued mrf models. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 477–484. IEEE, 2009.
- [35] Justin Domke. Generic methods for optimization-based modeling. In *Artificial Intelligence and Statistics*, pages 318–326, 2012.
- [36] Salah Rifai, Yoshua Bengio, Yann Dauphin, and Pascal Vincent. A generative process for sampling contractive auto-encoders. *arXiv preprint arXiv:1206.6434*, 2012.
- [37] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- [38] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1885–1894. JMLR. org, 2017.
- [39] David Lopez-Paz et al. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems*, pages 6467–6476, 2017.
- [40] Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- [41] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [42] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [43] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- [44] Spyros Gidaris and Nikos Komodakis. Dynamic few-shot visual learning without forgetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4367–4375, 2018.
- [45] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [46] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in neural information processing systems*, pages 4077–4087, 2017.
- [47] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. In *Advances in neural information processing systems*, pages 3630–3638, 2016.
- [48] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [49] Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 374–382, 2019.
- [50] Eden Belouadah and Adrian Popescu. Scail: Classifier weights scaling for class incremental learning. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1266–1275, 2020.
- [51] G W Wolf. Facility location: concepts, models, algorithms and case studies. 2011.
- [52] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [53] Matthew D. Zeiler, Marc’ Aurelio Ranzato, Rajat Monga, Mark Z. Mao, Kyeongcheol Yang, Quoc V. Le, Patrick Nguyen, Andrew W. Senior, Vincent Vanhoucke, Jeffrey Dean, and Geoffrey E. Hinton. On rectified linear units for speech processing. *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 3517–3521, 2013.
- [54] Andrew L Maas, Awni Y Hannun, and Andrew Y Ng. Rectifier nonlinearities improve neural network acoustic models. In *International conference on machine learning (ICML)*, 2013.
- [55] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ArXiv*, abs/1502.03167, 2015.
- [56] Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19, 2018.
- [57] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [58] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
- [59] Saeid Motiian, Quinn Jones, Seyed Iranmanesh, and Gianfranco Doretto. Few-shot adversarial domain adaptation. In *Advances in Neural Information Processing Systems*, pages 6670–6680, 2017.
- [60] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pages 248–255. Ieee, 2009.