

---

# Remember the Past: Distilling Datasets into Addressable Memories for Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose an algorithm that compresses the critical information of a large dataset  
2 into compact addressable memories. These memories can then be recalled to  
3 quickly re-train a neural network and recover the performance (instead of storing  
4 and re-training on the full original dataset).

5 Building upon the dataset distillation framework, we make a key observation that a  
6 *shared common representation* allows for more efficient and effective distillation.  
7 Concretely, we learn a set of bases (aka “memories”) which are shared between  
8 classes and combined through learned flexible addressing functions to generate  
9 a diverse set of training examples. This leads to several benefits: 1) the size of  
10 compressed data does not necessarily grow linearly with the number of classes; 2)  
11 an overall higher compression rate with more effective distillation is achieved; and  
12 3) more generalized queries are allowed beyond recalling the original classes.

13 We demonstrate state-of-the-art results on the dataset distillation task across five  
14 benchmarks, including up to 16.5% and 9.7% accuracy improvement when distilling  
15 CIFAR10 and CIFAR100 respectively. We then leverage our framework to  
16 perform continual learning, achieving state-of-the-art results on four benchmarks,  
17 with 23.2% accuracy improvement on MANY.

## 18 1 Introduction

19 Compressing a large amount of information into a small memory storage space is one of the key  
20 components in human intelligence [1, 2, 3] – a person can retrieve the memories from the past and  
21 fastly recover the corresponding skills. Deep learning methods have made large strides in building  
22 task-specific models, but are shown to easily forget the past knowledge when learning new tasks [4, 5].

23 To equip neural network learners with the memorizing ability, dataset distillation [6] is proposed as a  
24 potential solution, which performs pragmatic compression and prioritize information that affects the  
25 model training for compression. Concretely, a compressed set of examples (memories) is learned to  
26 summarize the key information in a dataset that affects model training; these examples can then be  
27 used to quickly retrain models and recover the corresponding skills. This differs from the standard  
28 reconstruction-based compression algorithms [7, 8, 9] and shows strong performance [10, 11, 12, 13].

29 A critical question in building powerful compressed memories is: what structures and representations  
30 should we use to build the memories? An effective structure and organization of memories can lead  
31 to different fundamental assumptions about data and affect the compression and learning behaviours.  
32 Existing works [10, 11, 14, 12, 15, 13, 6] follow a simple representation, where a set of learnable  
33 examples is assigned for each class. However, under this assumption, the size of the memories can  
34 linearly grow with the number of classes, making the distillation of dataset with large number of  
35 classes challenging. Naturally, this can potentially lead to redundancies in the learned memories,

36 due to the separation of data among classes. Furthermore, this representation is less generalizable to  
37 continuous label space, where infinite number of label values exists.

38 In our paper, we make the observation that there are information sharing among classes, and hypoth-  
39 esize that a common and compact representation exists for all classes. Following the hypothesis,  
40 we propose to formulate the problem as a memory addressing process, where the memories store a  
41 common set of basis shared by all labels, and the recombination of basis is performed through the  
42 addressing function. This decomposition between memories and addressing functions enables the  
43 possibility that all common information is stored in one part of the representation, and the accessing of  
44 the common information depends on the specific labels and is handled through an extra function. We  
45 find that this formulation can significantly improve both the compression rate and the performance.

46 We adopt the back-propagation through time learning framework to train the memories and addressing  
47 functions, and identify several critical factors that can improve the performance. Specifically, we find  
48 that adopting the momentum term, and performing long unrolls in the inner optimization loop, can  
49 drastically improve the performance. This differs from the common usage of bi-level optimization  
50 algorithm on this task [6, 16], and leads to strong performances outperforming single-step gradient  
51 matching methods [10, 11] even with the simple data representation.

52 In the experiments, we extensively evaluate our algorithm on five benchmarks of the Dataset Distilla-  
53 tion task, and show that it consistently outperforms previous state-of-the-art by a significant margin.  
54 For example, we achieve 66.4% accuracy on CIFAR10 with the storage space of 1 image per class,  
55 improving over the previous state-of-the-art KIP method [12, 13] by 16.5%. We further demonstrate  
56 our method on the continual learning tasks, and show that a simple “compress-then-recall” method  
57 using our framework leads to state-of-the-art results on four datasets. For example, we outperform all  
58 prior methods by 23.2% in retained accuracy on the challenging MANY [17] benchmark. Finally, we  
59 demonstrate the generality of our approach by extending to image-based (rather than label-based)  
60 memory recall, and synthesizing new classifiers (unseen during training) from our distilled memories.

## 61 2 Related works

62 **Dataset Distillation.** The task of dataset distillation is fundamentally a compression problem, with a  
63 different prioritization on the information contained in data. There have been mainly several lines of  
64 methods, developed with different criteria to prioritize information. *Generalization loss* with bi-level  
65 optimization framework [18, 19, 20] has been widely adopted in various tasks and is adopted in  
66 the early works of dataset distillation [6, 16]. It emphasizes the loss at the final optimization state.  
67 *Gradient-matching or score-matching* methods [10, 11, 15] are adopted to directly match the induced  
68 gradients from synthetic data. If ideally matched over the gradient field, the compressed dataset can  
69 naturally lead to the same model parameters with gradient descent. *Kernel method* [13, 12] shows  
70 that with the connection to Gaussian processes, a kernel inducing points method can be adopted to  
71 achieve strong performance, but large computation costs. These are also connected with the recent  
72 progress on pragmatic compression methods [21, 22], which compress or match distributions based  
73 on a decision process (in dataset distillation’s case, the gradient descent search process).

74 **Continual learning.** Broadening the learning paradigms, continual learning problem aims to build  
75 agents that learn through a stream of tasks and accrue knowledge along the process. “Catastrophic  
76 forgetting” [5, 23, 4] is a well-known phenomenon in this setting, where the neural network forgets  
77 previous skills when learning new ones. Various methods [24, 25, 26, 27, 17, 28, 29, 30, 31, 32, 33,  
78 34, 35] on regularization, replay, or dynamic model, have been proposed to alleviate the issue and  
79 address the “stability-plasticity dilemma” [36, 37]. Memory buffer has been a critical component in the  
80 past methods [24, 25, 26, 27, 17, 28, 38, 39], but mainly relies on random selection of real samples  
81 with different strategies. Recently, several works extend the usage of memory to storing random  
82 basis [39] (online setting) or SVD bases [38] (offline setting).

## 83 3 Background: dataset distillation

84 The task of Dataset Distillation [6] is proposed to compress the key information of a large-scale  
85 training dataset into a small amount of learned data, which can be stored using limited memory space  
86 and retrieved through label indices or task information to recover the performance of a model.

87 **Problem Setting.** Formally, given a large dataset  $\mathcal{D}_{tr} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$  containing  $N$  pairs of training  
88 data  $(\mathbf{x}_i, \mathbf{y}_i)$ , where  $\mathbf{x}_i$  is an image and  $\mathbf{y}_i$  is the corresponding label in  $C$  classes, a small dataset

89  $\mathcal{D}_s = \{(\mathbf{x}'_i, \mathbf{y}_i)\}_{i=1}^{N'}, N' \ll N$ , can be synthesized or distilled, such that a model trained on  $\mathcal{D}_s$  can  
 90 have the same generalization ability with ones trained on  $\mathcal{D}_{tr}$ :

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{te}} [\mathbf{m}(f(\mathbf{x}; \boldsymbol{\theta}^{(*)}), \mathbf{y})] \simeq \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{te}} [\mathbf{m}(f(\mathbf{x}; \boldsymbol{\theta}'^{(*)}), \mathbf{y})] \quad (1)$$

91 where  $\boldsymbol{\theta}^{(*)}$  and  $\boldsymbol{\theta}'^{(*)}$  are the optimized parameters using  $\mathcal{D}_{tr}$  and  $\mathcal{D}_s$  respectively,  $\mathbf{m}$  is a metric,  
 92 e.g., accuracy, and  $\mathcal{D}_{te}$  is the test dataset. The model is often a neural network classifier  $f(\cdot; \boldsymbol{\theta})$   
 93 parameterized by  $\boldsymbol{\theta}$  and can be trained with a loss function  $\ell(f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y})$ .

94 **Synthetic dataset representations.** The synthetic dataset contains the core parameters that need to  
 95 be learned in dataset distillation. The representation of the synthetic data affects the compactness  
 96 and effectiveness of the distillation process. In the existing methods [6, 10, 11, 14, 12, 15], the  
 97 dataset  $\mathcal{D}_s$  is defined as a collection of learnable data samples  $(\mathbf{x}', \mathbf{y})$ , and the number of samples  
 98 is separately and equally distributed across classes. This representation has several disadvantages:  
 99 firstly, the number of synthetic data samples needed for a dataset linearly grows with the number  
 100 of classes, leading to limited applicability when the number of classes is large or undefined (e.g.,  
 101 language or other continuous labels); secondly, the potentially shared and common information across  
 102 classes is ignored - this results in a less compact representation of the distilled information and lower  
 103 compression rate; lastly, the representation is not able to generalize to new classes or tasks, due to the  
 104 lack of common representation learned across classes.

105 **The underperforming back-propagation through time (BPTT) algorithm.** Bi-level optimization  
 106 using generalization loss to back-propagate gradients through inner loops is widely adopted across  
 107 different tasks [20, 18, 40]. However, albeit used in the original dataset distillation framework [6], it  
 108 has been shown to consistently underperform other competing algorithms, such as gradient matching  
 109 methods [10, 11]. In our paper, we point out several critical aspects that affect the performance and  
 110 show that BPTT can outperform single-step gradient matching methods by a large margin, which  
 111 contrasts with the current common observations in this community.

## 112 4 Model

113 **Overview.** In this section, firstly, we present a new perspective of the problem, where the Dataset  
 114 Distillation problem is formulated as a *memory addressing process* - instead of learning synthetic  
 115 images separately for each class - we construct and learn a common memory representation that can  
 116 be addressed through an addressing matrix to construct synthetic datasets. Under this formulation,  
 117 the number of synthetic images doesn't have to grow linearly with the number of classes, the shared  
 118 information among classes can be exploited to reduce redundancies and improve compression rate,  
 119 and datasets can be distilled with respect to more generic queries. Secondly, we further show  
 120 several critical empirical facets of back-propagation through time framework, which lead to drastic  
 121 improvements on the performance and outperform the single-step gradient matching methods. This is  
 122 in contrast with the current common observation that gradient matching outperforms back-propagation  
 123 through time framework on dataset distillation tasks.

124 In the following, we present the two core components of our method, (1) the new formulation of  
 125 dataset distillation, with memories and addressing matrices in Sec. 4.1, and (2) the learning framework  
 126 under back-propagation through time in Sec. 4.2.

### 127 4.1 Dataset distillation as memory addressing

128 **Problem formulation.** Given a task-specific dataset  $\mathcal{D}_{tr} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$ , we aim to learn a single  
 129 compact representation  $\mathcal{M}$ , referred to as memories, that can be accessed through a learned addressing  
 130 function  $\mathcal{A}(\cdot)$ , which takes all possible values of  $\mathbf{y}$  as input and recall the corresponding synthetic  
 131 data. With a set of  $\{\mathbf{y}_i\}$ , a synthetic dataset recalled using the above process can train a model  
 132  $f_{\boldsymbol{\theta}} : \mathcal{X} \rightarrow \mathcal{Y}$  from scratch and obtain the same generalization ability as trained on  $\mathcal{D}_{tr}$ .

133 Ideally, the memories  $\mathcal{M}$  and the addressing function  $\mathcal{A}$  can jointly capture the critical information  
 134 that defines the task mapping from  $\mathbf{x} \in \mathcal{X}$  to  $\mathbf{y} \in \mathcal{Y}$ , such that, when training a model using the  
 135 synthetic dataset collected through accessing memories with a set of  $\{\mathbf{y}_i\}$ , the trained model can  
 136 obtain similar or same generalization ability with ones trained on  $\mathcal{D}_{tr}$ . For example, in the standard  
 137 classification tasks, we can enumerate all possible values in the label space (discrete) and address the  
 138 memories, to construct the synthetic dataset.

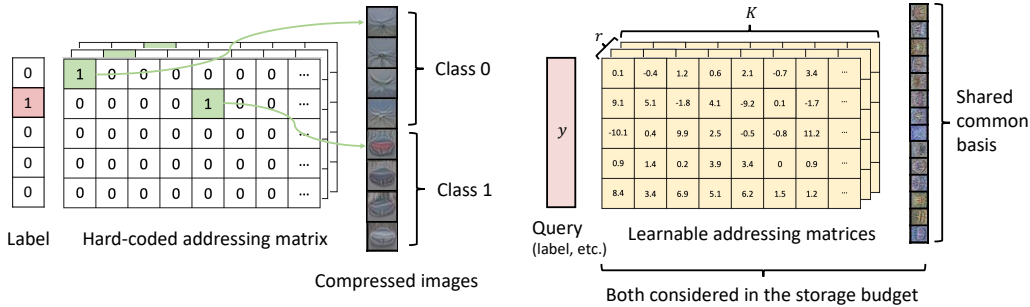


Figure 1: Distill a large-scale dataset into compressed memories. Left: the standard dataset distillation task under the formulation of memory addressing. Right: learnable addressing matrix with shared common bases. The queries can be generalized to any vector representation, besides one-hot labels, i.e. for a general dataset from  $\mathcal{X}$  to  $\mathcal{Y}$ , it can be distilled into memories for recall and model re-training.

139 Under this definition, since we are learning a single, shared and accessible representation for all  $y$ s,  
 140 the size of memories can be defined flexibly regardless of the number of classes, removing the linear  
 141 growth limitation in the standard distillation settings. There is also no constraint on the form of  $y$ s,  
 142 which can be either discrete or continuous. The *storage budget* or *compression rate* is calculated  
 143 by considering the storage space of parameters in both memories and addressing functions, which  
 144 should be as compact as possible.

145 **Memory representation.** We use a set of basis to store in the memories  $\mathcal{M} = \{b_0, \dots, b_{K-1}\}$ , where  
 146 each vector  $b_k \in \mathbb{R}^d$  has the same dimension as  $x \in \mathbb{R}^d$ , and all vectors collectively define the  
 147 intrinsic components in a dataset that characterize the task mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ . Through re-using  
 148 the basis, we can produce a desired synthetic dataset for model re-training. *Spatial redundancies in*  
 149 *images:* note that, as a special case, images can contain redundancies spatially and be stored in a  
 150 downsampled version to improve the compression rate. The downsampled image basis can be passed  
 151 via a deterministic upsampling process (e.g., bilinear interpolation) to recover the original resolution.

152 **Memory addressing.** For each query  $y$ , we use a parameterized function  $\mathcal{A}(y)$  to re-combine the  
 153 basis in the memories  $\mathcal{M}$ . Similar to previous methods on accessing memories, we use  $\mathcal{A}$  to produce  
 154 a set of coefficients, and linearly combine the basis to produce synthetic data. Formally, to retrieve  $r$   
 155 synthetic data for each  $y$ , we define a set of matrices  $\{A_0, \dots, A_{r-1}\}$ ,  $A_i \in \mathbb{R}^{d_y \times K}$ , where  $d_y$  is the  
 156 dimension size of  $y$  and  $r$  is the number of data samples that can be retrieved. With the memories  
 157  $\mathcal{M} = \{b_0, \dots, b_{K-1}\}$ , we define:

$$x'_i = y^T A_i [b_0; \dots; b_{K-1}]^T, x' \in \mathbb{R}^d \quad (2)$$

158 where  $y \in \mathbb{R}^{d_y \times 1}$  is in a vectorized form, such one-hot encoding of categorical labels, and  $v = y^T A_i$   
 159 corresponds to a coefficient vector  $v$  that combines the basis. The produced synthetic data  $x'$  is  
 160 paired with  $y$  as the corresponding label. Our model is shown on the right of figure 1.

161 **Constructing a dataset.** To construct a synthetic dataset  $\mathcal{D}_s$  for model re-training, we are often given  
 162 a set of samples  $\{y_i\}$ , or can enumerate all possible values of  $y$  (if discrete). With the set  $\{y_i\}$ , we  
 163 use eqn.2 to address and retrieve synthetic dataset  $\mathcal{D}_s^{y=y_i} = \{(x'_j, y_i)\}_{j=1}^r$  for each  $y_i$ . The final  
 164 dataset is the union of all  $\mathcal{D}_s = \bigcup \mathcal{D}_s^{y=y_i}$ . The dataset  $\mathcal{D}_s$  can be used in either minibatch forms for  
 165 stochastic gradient descent, or as a whole for batch gradient descent.

166 **Generalized possibilities of queries.** Another advantage of our formulation, besides a compact and  
 167 shared representation, is the various possibilities of queries  $y$ . In principle, under this formulation, the  
 168 label  $y$  can be flexibly defined as other forms, such as language, audio, of which the representation  
 169 resides in continuous space or follow a distribution  $p(h(y))$  defined by a feature extractor  $h(\cdot)$ . The  
 170 set of  $\{y_i\}$  then can be sampled from the distribution  $p$  instead of having to enumerate all possible  
 171 values. This provides a general way of compressing or distilling a large dataset without constraint on  
 172 the forms of labels.

173 **Connection with standard Dataset Distillation.** In the standard setting, the task is defined for  
 174 classification tasks with discrete labels and each label owns its unique set of synthetic data. We can

---

**Algorithm 1**

---

1: **hyperparameters:** Momentum rate  $\beta_\theta, \beta_\phi$ , learning rate  $\alpha_\theta, \alpha_\phi$  for  $\theta$  and  $\phi$  respectively.  
2: **input:** Dataset  $\mathcal{D}_{tr}$ , memories  $\mathcal{M}$ , addressing function  $\mathcal{A}$ , loss function  $\ell(\cdot, \cdot)$   
3: **repeat**  
4:   Sample a subset of labels  $\mathcal{Y}'$   
5:   Address the memories  $\mathcal{M}$  and obtain synthetic dataset  $\mathcal{D}_s^{\mathcal{Y}'}$  with eqn. 2  
6:   Randomly initialize model parameters  $\theta_0$   
7:   Initialize momentum  $m_0 = 0$   
8:   **for**  $t = 1$  **to**  $T$  **do**  
9:     Sample a minibatch  $B_s = \{(x'_i, y_i)\}$  from  $\mathcal{D}_s^{\mathcal{Y}'}$   
10:     Compute  $\mathcal{L} = \frac{1}{|B_s|} \sum_{i=1}^{|B_s|} \ell(f_{\theta_{t-1}}(x'_i), y_i)$   
11:     **Update momentum**  $m_t = \beta_\theta m_{t-1} + \frac{d\mathcal{L}}{d\theta_{t-1}}$   
12:     Update  $\theta_t = \theta_{t-1} - \alpha_\theta m_t$   
13:   **end for**  
14:   Sample a minibatch  $B = \{(x_i, y_i)\}$  from  $\mathcal{D}_{tr}$  with labels in  $\mathcal{Y}'$   
15:   Compute  $J(\phi) = \frac{1}{|B|} \sum_{i=1}^{|B|} \ell(f_{\theta_t}(x_i), y_i)$  Update  $\phi \leftarrow \text{opt-alg}(J(\phi), \alpha_\phi, \beta_\phi)$   
16: **until** Converge

---

175 show that this is a special case of our formulation: if the basis is constructed as the collection of those  
176 label-specific synthetic data ( $K = N'$ ), and the addressing matrix  $A_i$  is defined in space  $\{0, 1\}^{C \times N'}$   
177 with hardcoded 1s at position  $c + i$  of each row  $c$ , the "retrieval" process can also be defined as eqn. 2.

## 178 4.2 Learning framework: back-propagation through time

179 In this section, we build upon the back-propagation through time algorithm and discuss in detail the  
180 learning framework that performs the distillation process from a dataset to memories and addressing  
181 functions. Specifically, we identify several critical aspects in the algorithm that can boost the  
182 performance of the task.

183 Starting from notations, we define the parameters contained in both memories and addressing  
184 functions as  $\phi$ , which are collectively optimized. A loss function is  $\ell(\cdot, \cdot)$  is defined on the task-  
185 specific dataset  $\mathcal{D}_{tr}$ . We denote an optimization algorithm as  $\text{opt}(\cdot, \cdot; \alpha, \beta)$ , where  $\alpha$  and  $\beta$  are the  
186 learning rate and the momentum rate, respectively.

187 To learn the parameters  $\phi = \{\mathcal{M}, \mathcal{A}\}$ , we follow a standard bi-level optimization framework with  
188 back-propagation through time (BPTT), where the inner-loop uses the synthetic dataset  $\mathcal{D}_s$  to train  
189 a randomly initialized model starting from scratch, and a generalization loss is computed using a  
190 minibatch  $B = \{(x_i, y_i)\}$  sampled from  $\mathcal{D}_{tr}$ :

$$\begin{aligned} \min J(\phi) &= \frac{1}{|B|} \sum_{i=1}^{|B|} \ell(f(x_i; \theta^*), y_i), \\ \text{subject to } \theta^* &= \text{opt}(\theta_0, \mathcal{D}_s; \alpha_\theta, \beta_\theta) \end{aligned} \quad (3)$$

191 where  $\theta_0$  represents the initializing parameters and  $\theta^*$  is the optimized model parameters in the  
192 inner optimization loop. In practice, we use a subsampled version  $\mathcal{Y}'$  of all  $y$ s and retrieve the  
193 corresponding subset  $\mathcal{D}_s^{\mathcal{Y}'}$ . This reduces the computation cost in inner loops. We empirically observe  
194 that equivalent results can be achieved with faster speed. The algorithm is summarized in Alg. 1.

195 **Critical factors in BPTT.** Although a natural choice in performing dataset distillation, the BPTT  
196 framework has shown to underperform other algorithms, such as gradient matching methods, in  
197 various works and benchmarks. In our work, we empirically identify several aspects that can  
198 drastically improve the performance.

199 *Long unrolled trajectories.* In previous data distillation works [6], the number of inner loop optimiza-  
200 tion steps is set to relative small values, potentially due to the common practice in meta-learning and  
201 few-shot learning works [28, 18, 41]. Instead we find that unrolling the trajectories long enough (e.g.  
202 200 steps) can potentially produce  $\theta^*$  that better summarizes the information contained in memories  
203 and produce more effective gradients.

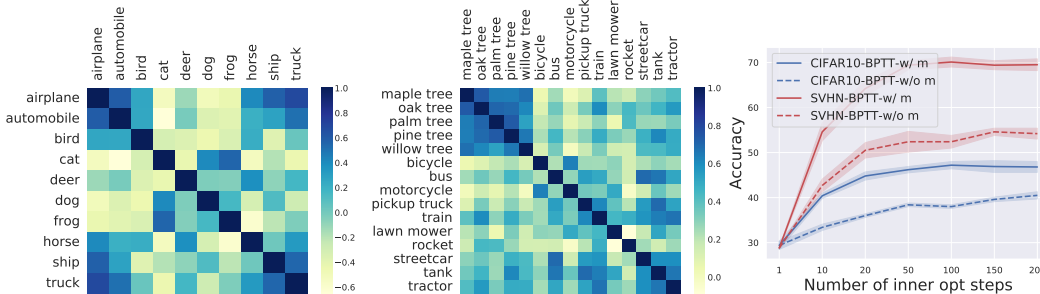


Figure 2: Similarity matrices of learned addressing coefficients for the CIFAR10 dataset (left) and a subset of CIFAR100 classes (right). Figure 3: Analysis on BPTT steps and momentums.

204 *Momentum term.* Another aspect ignored in dataset distillation and not commonly used in meta-  
 205 learning is the momentum term. Indeed, adding momentum term in meta-learning can potentially  
 206 even hurt the performance and lead to less gradient diversity [41]. However, we show that, in dataset  
 207 distillation, momentum term is critical, even in short inner loop optimization settings. We provide  
 208 more detailed analysis in the experiment section.

## 209 5 Experiments

210 We thoroughly evaluate our model under various tasks and benchmarks, and demonstrate the benefits  
 211 of our model over previous methods. First, we show that using a shared representation is core to  
 212 improve the distillation performance and compression rates. Specifically, we observe that there  
 213 is strong evidence that there is information re-using across classes. Moreover, we show that the  
 214 shared representation enables more general queries, which can be continuous. We further show the  
 215 benefits of our model on continual learning and synthesizing new classifiers. For example, we observe  
 216 that a simple compress-then-recall method can achieve performance outperforming state-of-the-art  
 217 continual learning models with complex designs.

### 218 5.1 Dataset Distillation

219 In this section, we follow the standard setting of dataset distillation, and compress a dataset into  
 220 memories that can be recalled with labels.

221 **Datasets.** We test our models on five standard dataset distillation benchmarks: MNIST [42], Fashion-  
 222 MNIST [43], SVHN [44], CIFAR10 [45], CIFAR100 [45]. MNIST contains 10 classes with 60,000  
 223 writing digit images as training and 10,000 as test set. The images are gray-scale with a shape of  
 224  $28 \times 28$ . FashionMNIST is a dataset with clothing and shoe images and consists of a training with  
 225 size 60,000 and a test set with size 10,000. Each image is  $28 \times 28$  in gray scale, and has a label from  
 226 10 classes. SVHN contains street digit images where each image has shape  $32 \times 32 \times 3$ . The dataset  
 227 contains 73257 images for training and 26032 images for testing. CIFAR10 and CIFAR100 are color  
 228 image datasets, with 50,000 training images and 10,000 testing images on each. CIFAR10 has 10  
 229 classes with 5,000 images per class, and CIFAR100 has 100 classes with 500 images per class.

230 **Experiment settings.** We evaluate our distillations model under three different memory budgets for  
 231 each dataset: 1/10/50 images per class. We focus on high compression rate scenarios and consider  
 232 1 and 10 settings for CIFAR100. Following the previous works [14, 11, 12, 15], the main network  
 233 architecture used in experiments is a simple convolutional network (ConvNet) with  $3 \times 3$  filters,  
 234 InstanceNorm, ReLU and  $2 \times 2$  average pooling. For the evaluation protocol, each model is evaluated  
 235 on 20 randomly initialized models, trained for 300 epochs on a synthetic dataset, and tested on a  
 236 held-out testing dataset. We use one GPU per experiment run.

237 **Memory budget calculation.** Since our model uses memories and addressing matrices to store the  
 238 compressed information, we treat the total number of images as a memory storage budget:

$$\text{size}(\text{basis}) + \text{size}(\text{addressing matrices}) \approx \text{size}(\text{images}) \quad (4)$$

239 where images is the total number of images for all classes and  $\text{size}(\cdot)$  is the total size of a tensor. We  
 240 assume all tensors are in the format of float numbers.

	I/C	DC [10]	DSA [11]	KIP [12] (NN)	CAFE* [46]	TM [15]	DM [14]	Ours
MNIST [42]	1	91.7±0.5	88.7±0.6	90.1±0.1	93.1±0.3	-	89.7±0.6	<b>98.7±0.7</b>
	10	97.4±0.2	97.8±0.1	97.5±0.0	97.5±0.1	-	97.5±0.1	<b>99.3±0.5</b>
	50	98.8±0.2	99.2±0.1	98.3±0.1	98.9±0.2	-	98.6±0.1	<b>99.4±0.4</b>
F-MNIST [43]	1	70.5±0.6	70.6±0.6	73.5±0.5	77.1±0.9	-	-	<b>88.5±0.1</b>
	10	82.3±0.4	84.6±0.3	86.8±0.1	83.0±0.3	-	-	<b>90.0±0.7</b>
	50	83.6±0.4	88.7±0.2	88.0±0.1	88.2±0.3	-	-	<b>91.2±0.3</b>
SVHN [44]	1	31.2±1.4	27.5±1.4	57.3±0.1	42.9±3.0	-	-	<b>87.3±0.1</b>
	10	76.1±0.6	79.2±0.5	75.0±0.1	77.9±0.6	-	-	<b>89.1±0.2</b>
	50	82.3±0.3	84.4±0.4	80.5±0.1	82.3±0.4	-	-	<b>89.5±0.2</b>
CIFAR10 [45]	1	28.3±0.5	28.8±0.7	49.9±0.2	31.6±0.8	46.3±0.8	26.0±0.8	<b>66.4±0.4</b>
	10	44.9±0.5	52.1±0.5	62.7±0.3	50.9±0.5	65.3±0.7	48.9±0.6	<b>71.2±0.4</b>
	50	53.9±0.5	60.6±0.5	68.6±0.2	62.3±0.4	71.6±0.2	63.0±0.4	<b>73.8±0.5</b>
CIFAR100 [45]	1	12.8±0.3	13.9±0.3	15.7±0.2	14.0±0.3	24.3±0.3	11.4±0.3	<b>34.0±0.4</b>
	10	25.2±0.3	32.3±0.3	28.3±0.1	31.5±0.2	40.1±0.4	29.7±0.3	<b>42.9±0.7</b>

Table 1: We compare our method with previous works on ConvNet performance recovery. Our algorithm consistently outperforms all previous methods and achieves state-of-the-art. \*Note that we selected the best results from baseline model variants. I/C is images per class (storage budget eqn. 4).

	I/C	Single-step GM	Ours <sup>BPTT</sup>	Ours <sup>BPTT+ds</sup>	Ours <sup>Full w/o Aug.</sup>	Ours <sup>Full</sup>
CIFAR10	1	28.8±0.7	49.1±0.6	55.2±0.5	64.2±0.6	<b>66.4±0.4</b>
	10	52.1±0.5	62.4±0.4	65.9±0.4	70.9±0.4	<b>71.2±0.4</b>
	50	60.6±0.5	70.5±0.4	71.1±0.5	72.1±0.5	<b>73.8±0.4</b>
CIFAR100	1	13.9±0.3	21.3±0.6	25.9±0.4	33.5±0.2	<b>34.0±0.4</b>
	10	32.3±0.3	34.7±0.5	36.5±0.4	40.6±0.3	<b>42.9±0.7</b>

Table 2: Ablation studies of every component and comparison with single-step gradient matching [10]. ds: downsample. Aug.: data augmentation

241 **Model details.** We use basis with size downsampled by a factor of  
242 2 from the standard image size based on datasets. All models are  
243 trained for 50k iterations with SGD optimizer. For the inner loop  
244 optimization, we set the momentum rate as 0.9, and use 150 steps for  
245 small memory budgets 1 and 10, and 200 steps for budget 50. The  
246 selection of number of basis to use is based on a held-out validation  
247 set (10% of training set), shown in figure 4. Number of addressing  
248 matrices is calculated with eqn. 4. More details in appendix.

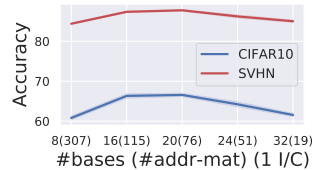


Figure 4: Validation set

249 **Result 1: state-of-the-art accuracy.** We compare our model  
250 with previous methods: Dataset Condensation (DC) [10], Differ-  
251 entiable Siamese Augmentation (DSA) [11], Kernel Inducing Points  
252 (KIP) [12], Distribution Matching (DM) [14], Aligning Features (CAFE) [46] and Trajectory Match-  
253 ing (TM) [15]. The results are summarized in table 1. Following previous methods [11, 12, 15], we  
254 adopt simple data augmentations, flip and rotation, on CIFAR10 and CIFAR100 datasets, and use  
255 ZCA on SVHN, CIFAR10 and CIFAR100. As shown in the table, our model consistently outperforms  
256 previous methods, especially under high compression rate cases where 1 image is allowed per class.

257 **Analysis: information sharing across classes.** The core observation in our method is that a common  
258 representation can enable information sharing across classes and reduce redundancies. To verify this,  
259 we calculate the average coefficients  $\bar{\mathbf{v}} = \frac{1}{r} \sum_{i=0}^{r-1} \mathbf{y}^T \mathbf{A}_i$  for each class  $\mathbf{y}$  and visualize the cosine  
260 similarities of  $\bar{\mathbf{v}}$  from two classes. The visualizations are shown in fig. 2. Higher cosine similarity  
261 scores indicate that two classes are utilizing similar basis components in the memories to produce  
262 synthetic images. For example, classes maple, oak, palm, pine and willow trees have strong sharing,  
263 while lawn mower and rocket are distinct from each other. Similar patterns can be found in CIFAR10.

264 **Result 2: back-propagation through time is a strong baseline.** In figure 3 and table 2, we show  
265 that BPTT is already a strong baseline that can outperform previous single-step gradient methods  
266 by up to 40% (on SVHN 1 image per class). Results shown in the figure are obtained through  
267 standard synthetic dataset organization without downsampling and memory-addressing framework.  
268 This contrasts to previous observation that single-step gradient matching methods always outperform  
269 BPTT with generalization loss [10, 11, 6]. The effects of momentum and long unrolls in BPTT are  
270 shown in figure 3. It is interesting to observe that even for shorter unrolls, momentum term is still  
271 important and can significantly improve the performance.

272 **Ablation studies.** To analyze the effects of each component, we perform ablation studies on CIFAR10  
 273 and CIFAR100, summarized in table 2. We show that the observations on BPTT, downsampling (ds)  
 274 due to spatial redundancies, and shared representations can all improve the results and compression  
 275 rate. We also ablate the effect of data augmentation (Aug.). See more in the appendix.

## 276 5.2 Continual Learning

277 One of the key usage of memories is to prevent forgetting when a model continually learns through  
 278 tasks. In this section, we evaluate our algorithm on the standard continual learning benchmarks and  
 279 show that, due to the strong performance, a simple “compress-then-recall” method with our model  
 280 can already rival with previous state-of-the-arts with complex designs.

281 **Efficient lifelong learning.** Following [47], we work with the problem where all tasks are streamed  
 282 in mini-batches and learned in a *single pass*. A learner is allowed to be equipped with a small memory  
 283 buffer. The data samples after seen will not be available unless stored in the buffer. We use mini-batch  
 284 size 10 to stream the data, following previous works [28, 24].

285 **Evaluation.** The learner’s performance after learning on the task stream is evaluated under two  
 286 metrics: retained accuracy (RA) and backward-transfer and interference (BTI). RA is the average  
 287 accuracy of the final trained model on all tasks, and BTI measures the performance difference between  
 288 after it was learned and after the full training process.

289 **Benchmarks.** We evaluate our method on three tasks widely used in previous Continual Learning  
 290 works. MNIST Rotations [24] contains 20 tasks with 1,000 samples in each. Every task consists  
 291 of images rotated by a fixed angle from 0 to 180 degrees. MNIST Permutations [48] has 20  
 292 tasks, and each task contains 1,000 images generated through shuffling the image pixels by a fixed  
 293 permutation. Many Permutations [17] is a longer variant with 100 tasks in total and 200 samples in  
 294 each. Incremental CIFAR-100 [29, 24] splits the CIFAR100 dataset into 20 5-way classification tasks  
 295 for learning.

296 **Our model.** Based on our distillation method, we adopt a simple framework to perform continual  
 297 learning: “compress then recall”. During the training phase, we do not perform learning on neural  
 298 networks, instead, the dataset of each task is distilled to memories and the paired addressing matrices.  
 299 During test phase, we simply fetch the corresponding memories and addressing matrices for each task,  
 300 and train a new model from scratch to perform classification. *Memory buffer designs.* When a new  
 301 task starts, we use the full remaining memory buffer to store the samples and perform distillation with  
 302 both buffer samples and streamed samples. After a task ends, the distilled memories and addressing  
 303 matrices are stored in the buffer, taking  $1/T$  of the space. Namely, the buffer size keeps shrinking  
 304 when more compressed representation of tasks is stored. Note that we compare our model with  
 305 previous methods under the *exact same memory sizes for fair comparisons*. See appendix.

306 **Results.** We show that this simple method is already a strong baseline that outperforms prior arts  
 307 on four benchmarks, summarized in table 3. Our method is compared with: Online, EWC [48],  
 308 GEM [24], MER [17], C-MAML [28], La-MAML [28] and Sparse-LaMAML [30]. For example, we  
 309 can obtain a 23% boost on MNIST MANY benchmark.

310 We further compare our model with previous works Kernel Continual Learning [39] and Stable  
 311 SGD [37] following their settings. We list the results here due to space limit: our model achieves

	Rotations		Permutations		MANY		CIFAR-100	
	RA $\uparrow$	BTI $\downarrow$	RA $\uparrow$	BTI $\downarrow$	RA $\uparrow$	BTI $\downarrow$	RA $\uparrow$	BTI $\downarrow$
ONLINE	53.38 $\pm$ 1.53	-5.44 $\pm$ 1.70	55.42 $\pm$ 0.65	-13.76 $\pm$ 1.19	32.62 $\pm$ 0.43	-19.06 $\pm$ 0.86	32.62 $\pm$ 0.43	-19.06 $\pm$ 0.86
EWC [48]	57.96 $\pm$ 1.33	-20.42 $\pm$ 1.60	62.32 $\pm$ 1.34	-13.32 $\pm$ 2.24	33.10 $\pm$ 0.14	-18.50 $\pm$ 0.91	-	-
GEM [24]	67.38 $\pm$ 1.75	-18.02 $\pm$ 1.99	55.42 $\pm$ 1.10	-24.42 $\pm$ 1.10	39.50 $\pm$ 0.62	-17.50 $\pm$ 0.41	48.27 $\pm$ 1.10	-13.7 $\pm$ 0.70
MER [17]	77.42 $\pm$ 0.78	-5.60 $\pm$ 0.70	73.46 $\pm$ 0.45	-9.96 $\pm$ 0.45	51.00 $\pm$ 0.54	-13.57 $\pm$ 0.45	51.38 $\pm$ 1.05	-12.83 $\pm$ 1.44
La-M [28]	77.42 $\pm$ 0.65	-8.64 $\pm$ 0.40	74.34 $\pm$ 0.67	-7.60 $\pm$ 0.51	50.43 $\pm$ 0.21	-10.00 $\pm$ 0.36	61.18 $\pm$ 1.44	-9.00 $\pm$ 0.2
sp-La [30]	77.77 $\pm$ 0.58	-8.16 $\pm$ 0.61	76.88 $\pm$ 0.72	-8.39 $\pm$ 0.63	50.81 $\pm$ 0.79	-13.73 $\pm$ 0.73	-	-
Ours	<b>80.32<math>\pm</math>0.28</b>	-	<b>78.48<math>\pm</math>0.76</b>	-	<b>74.07<math>\pm</math>0.51</b>	-	<b>62.58<math>\pm</math>1.1</b>	-

Table 3: We show that “compress-then-recall” is a strong baseline that outperforms previous methods on four continual learning benchmarks. Baseline numbers are from [28] or obtained from public official repos.

312  $87.34^{\pm 0.92}$  and  $88.25^{\pm 0.58}$  on Permuted MNIST and Rotated MNIST with 60,000 samples for each  
 313 task, outperforming both KCL and Stable SGD. Interestingly, our results also are higher than the  
 314 multitask upperbound ( $86.5^{\pm 0.21}$  and  $87.3^{\pm 0.47}$ ), potentially due to that there are task interference in  
 315 joint training, which can be naturally avoided in our method.

### 316 5.3 Synthesizing new classifiers after learning

317 If we want to memorize the past, what is the benefit of storing the compressed representation rather  
 318 than a trained model? In this section, we show that our compressed representation can enable flexible  
 319 synthesis on new classifiers after the learning. Specifically, we demonstrate on extrapolating between  
 320 tasks to train new models, and performing memory recall with images instead of labels, showing the  
 321 generalizability of our framework on other query forms.

#### 322 5.3.1 Extrapolating between tasks

323 In real world, tasks often do not come together and a learner therefore cannot observe all tasks at once.  
 324 In current machine learning paradigms, when models are separately trained for disjoint tasks, it has  
 325 difficulty extrapolating between tasks to build new classifiers. This is different from human learning.  
 326 We show that, if storing our compressed representation, it enables a learner to extrapolate and  
 327 synthesize new classifiers after learning separately on each task. Specifically, we separate CIFAR100  
 328 into 20 disjoint 5-way classification tasks. Independently, we distill the 20 datasets into corresponding  
 329 memories and test the possibility of generating new classifiers on new class combinations unseen  
 330 during training. We find that the compressed data can indeed train classifiers on new combinations,  
 331 for example, we can achieve 72.53% on 2-way classification, and 46.54% on 5-way classification,  
 332 with 1 image per class storage budget. The upperbound with full real dataset is 92.23% and 82.72%.

#### 333 5.3.2 Dataset Distillation extension – recall the past with images

334 We extend the standard setting to recall the past with images: when the label information and task  
 335 scopes are missing, but a few visual observations can be made, we would like to build classifiers  
 336 based on the visual data. For example, when we see a bear image and a deer image, but cannot recall  
 337 the exact word or category, can we recall the memories with images and build a classifier? This is  
 338 possible with our problem formulation, where the forms of queries are not constraint to labels and we  
 339 can *distill a dataset to memories addressable by images*.

340 Formally, after observing a training dataset  $\mathcal{D}_{tr}$  with  $\mathcal{Y} =$   
 341  $\{0, \dots, C - 1\}$ , we would like to flexibly build classifiers for  
 342 a subtask  $\mathcal{Y}_g \subset \mathcal{Y}$  based on visual observations  $\mathcal{X}_g$  from  $\mathcal{Y}_g$ ,  
 343 when the actual information of  $\mathcal{Y}_g$  is unknown. We work with  
 344 1-shot and 5-shot observation cases. As a baseline, we build  
 345 a nearest neighbor classifier, which is pretrained on  $\mathcal{D}_{tr}$  and  
 346 takes features of few-shot data to classify test images. As a  
 347 model variant, we could also “classify-then-recall”, using a  
 348 classifier trained on  $\mathcal{D}_{tr}$  to map the image shots into labels and  
 349 turning into the standard setup. Benefiting from the general design of our formulation, we show that  
 350 a model can directly perform “image addressing”, where a feature network can provide query vectors  
 351  $\mathbf{y}$  in fig.1. The feature network, memories and addressing matrices can be jointly trained on  $\mathcal{D}_{tr}$ .

Methods	1 shot	5 shot
Nearest neighbor	48.55	61.72
classify-then-recall	50.58	58.46
image addressing	<b>55.74</b>	<b>71.20</b>

Table 4: Few-shot perf. recovery.

## 352 6 Conclusion and limitations

353 In this paper, we propose a framework that distills a large dataset into compact addressable memories.  
 354 This framework introduce several benefits, including removing the linear growth constraint on the  
 355 compressed data size, allowing more general queries beside categorical labels, and most importantly,  
 356 high compression rate with strong re-training performance which outperforms previous state-of-  
 357 the-arts. We also demonstrate a “compress-then-recall” method using our framework, leading to  
 358 new state-of-the-arts on four datasets. Our full model has potential limitation on the costly inner  
 359 optimization loop, which might be time consuming on gigantic models or on much larger datasets.  
 360 This limitation might be solved through combining the memory formulation with a different learning  
 361 framework. We do not observe obvious negative societal impacts in our proposed method.

## 362 References

- 363 [1] Timothy F Brady, Talia Konkle, and George A Alvarez. Compression in visual working  
364 memory: using statistical regularities to form more efficient memory representations. *Journal*  
365 *of Experimental Psychology: General*, 138(4):487, 2009.
- 366 [2] Geoffrey R Loftus and Elizabeth F Loftus. *Human memory: The processing of information*.  
367 Psychology Press, 2019.
- 368 [3] John R Anderson and Gordon H Bower. *Human associative memory*. Psychology press, 2014.
- 369 [4] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical  
370 investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint*  
371 *arXiv:1312.6211*, 2013.
- 372 [5] Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks:  
373 The sequential learning problem. In *Psychology of learning and motivation*, volume 24, pages  
374 109–165. Elsevier, 1989.
- 375 [6] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A Efros. Dataset distillation.  
376 *arXiv preprint arXiv:1811.10959*, 2018.
- 377 [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
378 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural*  
379 *information processing systems*, 27, 2014.
- 380 [8] Ian Gemp, Brian McWilliams, Claire Vernade, and Thore Graepel. Eigengame: Pca as a nash  
381 equilibrium. In *International Conference on Learning Representations*, 2020.
- 382 [9] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*  
383 *arXiv:1312.6114*, 2013.
- 384 [10] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching.  
385 In *International Conference on Learning Representations*, 2021.
- 386 [11] Bo Zhao and Hakan Bilen. Dataset condensation with differentiable siamese augmentation. In  
387 *International Conference on Machine Learning*, 2021.
- 388 [12] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with  
389 infinitely wide convolutional networks. *Advances in Neural Information Processing Systems*,  
390 34, 2021.
- 391 [13] Timothy Nguyen, Zhourong Chen, and Jaehoon Lee. Dataset meta-learning from kernel  
392 ridge-regression. In *International Conference on Learning Representations*, 2021.
- 393 [14] Bo Zhao and Hakan Bilen. Dataset condensation with distribution matching. *arXiv preprint*  
394 *arXiv:2110.04181*, 2021.
- 395 [15] George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A Efros, and Jun-Yan Zhu.  
396 Dataset distillation by matching training trajectories. *arXiv preprint arXiv:2203.11932*, 2022.
- 397 [16] Iliia Sucholutsky and Matthias Schonlau. Soft-label dataset distillation and text dataset distilla-  
398 tion. *arXiv preprint arXiv:1910.02551*, 2019.
- 399 [17] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Ger-  
400 ald Tesauero. Learning to learn without forgetting by maximizing transfer and minimizing  
401 interference. *arXiv preprint arXiv:1810.11910*, 2018.
- 402 [18] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adapta-  
403 tion of deep networks. In *International Conference on Machine Learning*, pages 1126–1135.  
404 PMLR, 2017.
- 405 [19] Luca Franceschi, Michele Donini, Paolo Frasconi, and Massimiliano Pontil. Forward and  
406 reverse gradient-based hyperparameter optimization. In *International Conference on Machine*  
407 *Learning*, pages 1165–1173. PMLR, 2017.

- 408 [20] Aniruddh Raghu, Maithra Raghu, Simon Kornblith, David Duvenaud, and Geoffrey Hinton.  
409 Teaching with commentaries. In *International Conference on Learning Representations*, 2020.
- 410 [21] Sid Reddy, Anca Dragan, and Sergey Levine. Pragmatic image compression for human-in-the-  
411 loop decision-making. *Advances in Neural Information Processing Systems*, 34, 2021.
- 412 [22] Shengjia Zhao, Abhishek Sinha, Yutong He, Aidan Perreault, Jiaming Song, and Stefano Ermon.  
413 Comparing distributions by measuring differences that affect decision making. In *International  
414 Conference on Learning Representations*, 2021.
- 415 [23] Roger Ratcliff. Connectionist models of recognition memory: constraints imposed by learning  
416 and forgetting functions. *Psychological review*, 97(2):285, 1990.
- 417 [24] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning.  
418 *Advances in neural information processing systems*, 30:6467–6476, 2017.
- 419 [25] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K  
420 Dokania, Philip HS Torr, and Marc’Aurelio Ranzato. On tiny episodic memories in continual  
421 learning. *arXiv preprint arXiv:1902.10486*, 2019.
- 422 [26] Tom Veniat, Ludovic Denoyer, and Marc’Aurelio Ranzato. Efficient continual learning with  
423 modular networks and task-driven priors. *arXiv preprint arXiv:2012.12631*, 2020.
- 424 [27] Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual  
425 learning. *arXiv preprint arXiv:1710.10628*, 2017.
- 426 [28] Gunshi Gupta, Karmesh Yadav, and Liam Paull. La-maml: Look-ahead meta learning for  
427 continual learning. *arXiv preprint arXiv:2007.13904*, 2020.
- 428 [29] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl:  
429 Incremental classifier and representation learning. In *Proceedings of the IEEE conference on  
430 Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- 431 [30] Johannes Von Oswald, Dominic Zhao, Seijin Kobayashi, Simon Schug, Massimo Caccia,  
432 Nicolas Zucchet, and João Sacramento. Learning where to learn: Gradient sparsity in meta and  
433 continual learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- 434 [31] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that  
435 questions our progress in continual learning. In *European conference on computer vision*, pages  
436 524–540. Springer, 2020.
- 437 [32] Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis,  
438 Gregory Slabaugh, and Tinne Tuytelaars. Continual learning: A comparative study on how to  
439 defy forgetting in classification tasks. *arXiv preprint arXiv:1909.08383*, 2(6), 2019.
- 440 [33] Raia Hadsell, Dushyant Rao, Andrei A Rusu, and Razvan Pascanu. Embracing change: Con-  
441 tinual learning in deep neural networks. *Trends in cognitive sciences*, 24(12):1028–1040,  
442 2020.
- 443 [34] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick,  
444 Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv  
445 preprint arXiv:1606.04671*, 2016.
- 446 [35] Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with  
447 dynamically expandable networks. In *International Conference on Learning Representations*,  
448 2018.
- 449 [36] Martial Mermillod, Aurélie Bugaïska, and Patrick Bonin. The stability-plasticity dilemma: In-  
450 vestigating the continuum from catastrophic forgetting to age-limited learning effects. *Frontiers  
451 in psychology*, 4:504, 2013.
- 452 [37] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Razvan Pascanu, and Hassan Ghasemzadeh. Under-  
453 standing the role of training regimes in continual learning. *Advances in Neural Information  
454 Processing Systems*, 33:7308–7320, 2020.

- 455 [38] Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning.  
456 In *International Conference on Learning Representations*, 2020.
- 457 [39] Mohammad Mahdi Derakhshani, Xiantong Zhen, Ling Shao, and Cees Snoek. Kernel continual  
458 learning. In *International Conference on Machine Learning*, pages 2621–2631. PMLR, 2021.
- 459 [40] Jonathan Lorraine, Paul Vicol, and David Duvenaud. Optimizing millions of hyperparameters  
460 by implicit differentiation. In *International Conference on Artificial Intelligence and Statistics*,  
461 pages 1540–1552. PMLR, 2020.
- 462 [41] Alex Nichol and John Schulman. Reptile: a scalable metalearning algorithm. *arXiv preprint*  
463 *arXiv:1803.02999*, 2(3):4, 2018.
- 464 [42] Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE*  
465 *Signal Processing Magazine*, 29(6):141–142, 2012.
- 466 [43] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for  
467 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 468 [44] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.  
469 Reading digits in natural images with unsupervised feature learning. 2011.
- 470 [45] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
471 2009.
- 472 [46] Kai Wang, Bo Zhao, Xiangyu Peng, Zheng Zhu, Shuo Yang, Shuo Wang, Guan Huang, Hakan  
473 Bilen, Xinchao Wang, and Yang You. Cafe learning to condense dataset by aligning features. In  
474 *IEEE/CVF Conference on Computer Vision and Pattern Recognition 2022*, 2022.
- 475 [47] Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient  
476 lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018.
- 477 [48] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins,  
478 Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al.  
479 Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of*  
480 *sciences*, 114(13):3521–3526, 2017.

## 481 Checklist

- 482 1. For all authors...
- 483 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
484 contributions and scope? [Yes]
- 485 (b) Did you describe the limitations of your work? [Yes]
- 486 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 487 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
488 them? [Yes]
- 489 2. If you are including theoretical results...
- 490 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 491 (b) Did you include complete proofs of all theoretical results? [N/A]
- 492 3. If you ran experiments...
- 493 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
494 perimental results (either in the supplemental material or as a URL)? [Yes] See  
495 supplementary materials.
- 496 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were  
497 chosen)? [Yes] Details are in main paper and more in the supplementary materials.
- 498 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
499 ments multiple times)? [Yes]

- 500 (d) Did you include the total amount of compute and the type of resources used (e.g.,  
501 type of GPUs, internal cluster, or cloud provider)? [Yes] More details are in the  
502 supplementary materials.
- 503 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 504 (a) If your work uses existing assets, did you cite the creators? [Yes]  
505 (b) Did you mention the license of the assets? [N/A]  
506 (c) Did you include any new assets either in the supplemental material or as a URL? [No]  
507 (d) Did you discuss whether and how consent was obtained from people whose data you're  
508 using/curating? [N/A]  
509 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
510 information or offensive content? [N/A]
- 511 5. If you used crowdsourcing or conducted research with human subjects...
- 512 (a) Did you include the full text of instructions given to participants and screenshots, if  
513 applicable? [N/A]  
514 (b) Did you describe any potential participant risks, with links to Institutional Review  
515 Board (IRB) approvals, if applicable? [N/A]  
516 (c) Did you include the estimated hourly wage paid to participants and the total amount  
517 spent on participant compensation? [N/A]