

REMIND Your Neural Network to Prevent Catastrophic Forgetting

Tyler L. Hayes^{1(⊠)}, Kushal Kafle², Robik Shrestha¹, Manoj Acharya¹, and Christopher Kanan^{1,3,4}

Rochester Institute of Technology, Rochester, NY 14623, USA {tlh6792,rss9369,ma7583,kanan}@rit.edu
Adobe Research, San Jose, CA 95110, USA kkafle@adobe.com
³ Paige, New York, NY 10036, USA
⁴ Cornell Tech, New York, NY 10044, USA

Abstract. People learn throughout life. However, incrementally updating conventional neural networks leads to catastrophic forgetting. A common remedy is replay, which is inspired by how the brain consolidates memory. Replay involves fine-tuning a network on a mixture of new and old instances. While there is neuroscientific evidence that the brain replays compressed memories, existing methods for convolutional networks replay raw images. Here, we propose REMIND, a brain-inspired approach that enables efficient replay with compressed representations. REMIND is trained in an online manner, meaning it learns one example at a time, which is closer to how humans learn. Under the same constraints, REMIND outperforms other methods for incremental class learning on the ImageNet ILSVRC-2012 dataset. We probe REMIND's robustness to data ordering schemes known to induce catastrophic forgetting. We demonstrate REMIND's generality by pioneering online learning for Visual Question Answering (VQA) (https://github.com/tylerhayes/REMIND).

Keywords: Online learning · Brain-inspired · Deep learning

1 Introduction

The mammalian brain engages in continuous online learning of new skills, objects, threats, and environments. The world provides the brain a temporally structured stream of inputs, which is not independent and identically distributed (iid). Enabling online learning in artificial neural networks from non-iid data is known as lifelong learning. While conventional networks suffer from

T. L. Hayes, K. Kafle, R. Shrestha—Equal Contribution.

Electronic supplementary material The online version of this chapter (https://doi.org/10.1007/978-3-030-58598-3_28) contains supplementary material, which is available to authorized users.

© Springer Nature Switzerland AG 2020 A. Vedaldi et al. (Eds.): ECCV 2020, LNCS 12353, pp. 466–483, 2020. https://doi.org/10.1007/978-3-030-58598-3_28 catastrophic forgetting [1,57], with new learning overwriting existing representations, a wide variety of methods have recently been explored for overcoming this problem [13,15,27,47,53,58,64,77]. Some of the most successful methods for mitigating catastrophic forgetting use variants of replay [13,22,27,45,64,77], which involves mixing new instances with old ones and fine-tuning the network with this mixture. Replay is motivated by how the brain works: new experiences are encoded in the hippocampus and then these compressed memories are re-activated along with other memories so that the neocortex can learn them [51,60,72]. Without the hippocampus, people lose the ability to learn new semantic categories [48]. Replay occurs both during sleep [31] and when awake [41,74].

For lifelong learning in convolutional neural networks (CNNs), there are two major gaps between existing methods and how animals learn. The first is that replay is implemented by storing and replaying raw pixels, which is not biologically plausible. Based on hippocampal indexing theory [75], the hippocampus stores compressed representations of neocortical activity patterns while awake. To consolidate memories, these patterns are replayed and then the corresponding neocortical neurons are re-activated via reciprocal connectivity [51,60,72]. The representations stored in the hippocampus for replay are not veridical (e.g., raw pixels) [31,56], and its visual inputs are high in the visual processing hierarchy [29] rather than from primary visual cortex or retina.

The second major gap with existing approaches is that animals engage in streaming learning [20,21], or resource constrained online learning from non-iid (temporally correlated) experiences throughout life. In contrast, the most common paradigm for incremental training of CNNs is to break the training dataset into M distinct batches, where for ImageNet each batch typically has about 100000 instances from 100 classes that are not seen in later batches, and then the algorithm sequentially loops over each batch many times. This paradigm is not biologically plausible. There are many applications requiring online learning of non-iid data streams, where batched learning will not suffice, such as immediate on-device learning. Batched sys-

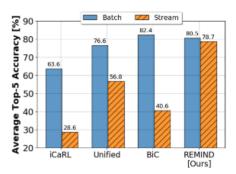


Fig. 1. Average top-5 accuracy results for streaming and incremental batch versions of state-ofthe-art models on ImageNet.

tems also take longer to train, further limiting their utility on resource constrained devices, such as smart appliances, robots, and toys. For example, BiC, a state-of-the-art incremental batch method, requires 65 hours to train in that paradigm whereas our proposed streaming model trains in under 12 hours. The incremental batch setting can be transformed into the streaming learning scenario by using very small batches and performing only a single pass through the dataset; however, this results in a large decrease in performance. As shown in

Fig. 1, state-of-the-art methods perform poorly on ImageNet in the streaming setting, with the best method suffering an over 19% drop in performance. In contrast, our model outperforms the best streaming model by 21.9% and is only 1.9% below the best batch model.

Here, we propose REMIND, or replay using memory indexing, a novel method that is heavily influenced by biological replay and hippocampal indexing theory. Our main contributions are:

- We introduce REMIND, a streaming learning model that implements hippocampal indexing theory using tensor quantization to efficiently store hidden representations (e.g., CNN feature maps) for later replay. REMIND implements this compression using Product Quantization (PQ) [30]. We are the first to test if forgetting in CNNs can be mitigated by replaying hidden representations rather than raw pixels.
- 2. REMIND outperforms existing models on the ImageNet ILSVRC-2012 [68] and CORe50 [52] datasets, while using the same amount of memory.
- 3. We demonstrate REMIND's robustness by pioneering streaming Visual Question Answering (VQA), in which an agent must answer questions about images and cannot be readily done with existing models. We establish new experimental paradigms, baselines, and metrics and subsequently achieve strong results on the CLEVR [33] and TDIUC [35] datasets.

2 Problem Formulation

There are multiple paradigms in which incremental learning has been studied [61]. In incremental batch learning, at each time step t an agent learns a data batch B_t containing N_t instances and their corresponding labels, where N_t is often 1000 to 100000. While much recent work has focused on incremental batch learning [13,14,18,27,44,45,64,77,81], streaming learning, or online learning from non-iid data streams with memory and/or compute constraints, more closely resembles animal learning and has many applications [20,21,49]. In streaming learning, a model learns online in a single pass, i.e., $N_t = 1$ for all t. It cannot loop over any portion of the (possibly infinite) dataset, and it can be evaluated at any point rather than only between large batches. Streaming learning can be approximated by having a system queue up small, temporally contiguous, mini-batches for learning, but as shown in Fig. 1, batch methods cannot easily adapt to this setting.

3 Related Work

Parisi et al. [61] identify three main mechanisms for mitigating forgetting in neural networks, namely 1) replay of previous knowledge, 2) regularization mechanisms to constrain parameter updates, and 3) expanding the network as more data becomes available. Replay has been shown to be one of the most effective methods for mitigating catastrophic forgetting [4,5,13,22,27,44,45,50,59,64,

77]. For ImageNet, all recent state-of-the-art methods for incremental class learning use replay of raw pixels with distillation loss. The earliest was iCaRL [64], which stored 20 images per class for replay. iCaRL used a nearest class pro-model [13] extended iCaRL to use the outputs of the CNN directly for classification, instead of a nearest class mean classifier. Additionally, End-to-End used more data augmentation and a balanced fine-tuning stage during training to improve performance. The Unified classifier [27] extended End-to-End by using a cosine normalization layer, a new loss constraint, and a margin ranking loss. The Bias Correction (BiC) [77] method extended End-to-End by training two additional parameters to correct bias in the output layer due to class imbalance. iCaRL, End-to-End, the Unified classifier, and BiC all: 1) store the same number of raw replay images per class, 2) use the same herding procedure for prototype selection, and 3) use distillation loss to prevent forgetting. REMIND, however, is the first model to demonstrate that storing and replaying quantized mid-level CNN features is an effective strategy to mitigate forgetting.

Regularization methods vary a weight's plasticity based on how important it is to previous tasks. These methods include Elastic Weight Consolidation (EWC) [47], Memory Aware Synapses (MAS) [3], Synaptic Intelligence (SI) [81], Riemannian Walk (RWALK) [14], Online Laplace Approximator [66], Hard Attention to the Task [70], and Learning without Memorizing [16]. The Averaged Gradient Episodic Memory (A-GEM) [15] model extends Gradient Episodic Memory [53], which uses replay with regularization. Variational Continual Learning [58] combines Bayesian inference with replay, while the Meta-Experience Replay model [65] combines replay with meta-learning. All of these regularization methods are typically used for incremental task learning, where batches of data are labeled as different tasks and the model must be told which task (batch) a sample came from during inference. When task labels are not available at test time, which is often true for agents operating in real-time, many methods cannot be used or they will fail [14, 17, 45]. While our main experiments focus on comparisons against state-of-the-art ImageNet models, we compare REMIND against several regularization models in Sect. 7, both with and without task labels. Some regularization methods also utilize cached data, e.g., GEM and A-GEM.

Another approach to mitigating forgetting is to expand the network as new tasks are observed, e.g., Progressive Neural Networks [69], Dynamically Expandable Networks [79], Adaptation by Distillation [26], and Dynamic Generative Memory [59]. However, these approaches also use task labels at test time, have growing memory requirements, and may not scale to thousand-category datasets.

4 REMIND: Replay Using Memory Indexing

REMIND is a novel brain-inspired method for training the parameters of a CNN in the streaming setting using replay. Learning involves two steps: 1) compressing the current input and 2) reconstructing a subset of previously compressed representations, mixing them with the current input, and updating the *plastic*

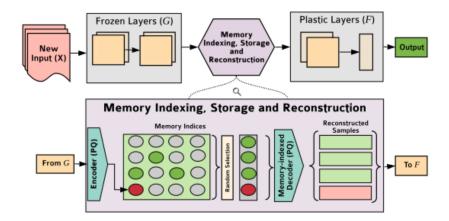


Fig. 2. REMIND takes in an input image and passes it through frozen layers of the network (G) to obtain tensor representations (feature maps). It then quantizes the tensors via product quantization and stores the indices in memory for future replay. The decoder reconstructs tensors from the stored indices to train the plastic layers (F) of the network before a final prediction is made.

weights of the network with this mixture (see Fig. 2). While earlier work for incremental batch learning with CNNs stored raw images for replay [13,27,64,77], by storing compressed mid-level CNN features, REMIND is able to store far more instances with a smaller memory budget. For example, iCaRL [64] uses a default memory budget of 20 K examples for ImageNet, but REMIND can store over 1M compressed instances using the same budget. This more closely resembles how replay occurs in the brain, with high-level visual representations being sent to the hippocampus for storage and re-activation, rather than early visual representations [29]. REMIND does not have an explicit sleep phase, with replay more closely resembling that during waking hours [41,74].

Formally, our CNN $y_i = F(G(\mathbf{X}_i))$ is trained in a streaming paradigm, where \mathbf{X}_i is the input image and y_i is the predicted output category. The network is composed of two nested functions: $G(\cdot)$, parameterized by θ_G , consists of the first J layers of the CNN and $F(\cdot)$, parameterized by θ_F , consists of the last L layers. REMIND keeps θ_G fixed since early layers of CNNs have been shown to be highly transferable [80]. The later layers, $F(\cdot)$, are trained in the streaming paradigm using REMIND. We discuss how $G(\cdot)$ is initialized in Sect. 4.2.

The output of $G(\mathbf{X}_i)$ is a tensor $\mathbf{Z}_i \in \mathbb{R}^{m \times m \times d}$, where m is the dimension of the feature map and d is the number of channels. Using the outputs of $G(\cdot)$, we train a vector quantization model for the \mathbf{Z}_i tensors. As training examples are observed, the quantization model is used to store the \mathbf{Z}_i features and their labels in a replay buffer as an $m \times m \times s$ array of integers using as few bits as necessary, where s is the number of indices that will be stored. For replay, we uniformly select r instances from the replay buffer, which was shown to work well in [14], and reconstruct them. Each of the reconstructed instances, $\hat{\mathbf{Z}}_i$, are mixed with the current input, and then θ_F is updated using backpropagation on this set of r+1 instances. Other selection strategies are discussed in Sect. 8.

During inference, we pass an image through $G(\cdot)$, and then the output, \mathbf{Z}_i , is quantized and reconstructed before being passed to $F(\cdot)$.

Our main version of REMIND uses PQ [30] to compress and store Z_i . For high-dimensional data, PQ tends to have much lower reconstruction error than models that use only k-means. The tensor Z_i consists of $m \times m$ d-dimensional tensor elements, and PQ partitions each d-dimensional tensor element into s subvectors, each of size d/s. PQ then creates a separate codebook for each partition by using k-means, where the codes within each codebook correspond to the centroids learned for that partition. Since the quantization is done independently for each partition, each sub-vector of the d-dimensional tensor element is assigned a separate integer, so the element is represented with s integers. If s is equal to one, then this approach is identical to using k-means for vector quantization, which we compare against. For our experiments, we set s = 32 and c = 256, so that each integer can be stored with 1 byte. We explore alternative values of s and c in supplemental materials (Fig. S4) and use the Faiss PQ implementation [32].

Since lifelong learning systems must be capable of learning from infinitely long data streams, we subject REMIND's replay buffer to a maximum memory restriction. That is, REMIND stores quantization indices in its buffer until this maximum capacity has been reached. Once the buffer is full and a new example comes in, we insert the new sample and randomly remove an example from the class with the most examples, which was shown to work well in [14,77]. We discuss other strategies for maintaining the replay buffer in Sect. 8.

4.1 Augmentation During Replay

To augment data during replay, REMIND uses random resized crops and a variant of manifold mixup [76] on the quantized tensors directly. For random crop augmentation, the tensors are randomly resized, then cropped and bilinearly interpolated to match the original tensor dimensions. To produce more robust representations, REMIND mixes features from multiple classes using manifold mixup. That is, REMIND uses its replay buffer to reconstruct two randomly chosen sets, \mathcal{A} and \mathcal{B} , of r instances each ($|\mathcal{A}| = |\mathcal{B}| = r$), which are linearly combined to obtain a set \mathcal{C} of r mixed instances ($|\mathcal{C}| = r$), i.e., a newly mixed instance, ($\mathbf{Z}_{\text{mix}}, y_{\text{mix}}$) $\in \mathcal{C}$, is formed as:

$$(\mathbf{Z}_{\text{mix}}, y_{\text{mix}}) = (\lambda \mathbf{Z}_a + (1 - \lambda) \mathbf{Z}_b, \lambda y_a + (1 - \lambda) y_b), \qquad (1)$$

where (\mathbf{Z}_a, y_a) and (\mathbf{Z}_b, y_b) denote instances from \mathcal{A} and \mathcal{B} respectively and $\lambda \sim \beta(\alpha, \alpha)$ is the mixing coefficient drawn from a β -distribution parameterized by hyperparameter α . We use $\alpha = 0.1$, which we found to work best in preliminary experiments. The current input is then combined with the set \mathcal{C} of r mixed samples, and θ_F is updated using this new set of r+1 instances.

4.2 Initializing REMIND

During learning, REMIND only updates $F(\cdot)$, i.e., the top of the CNN. It assumes that $G(\cdot)$, the lower level features of the CNN, are fixed. This implies

that the low-level visual representations must be highly transferable across image datasets, which is supported empirically [80]. There are multiple methods for training $G(\cdot)$, including supervised pre-training on a portion of the dataset, supervised pre-training on a different dataset, or unsupervised self-taught learning using a convolutional auto-encoder. Here, we follow the common practice of doing a 'base initialization' of the CNN [13,27,64,77]. This is done by training both θ_F and θ_G jointly on an initial subset of data offline, e.g., for class incremental learning on ImageNet we use the first 100 classes. After base initialization, θ_G is no longer plastic. All of the examples \mathbf{X}_i in the base initialization are pushed through the model to obtain $\mathbf{Z}_i = G(\mathbf{X}_i)$, and all of these \mathbf{Z}_i instances are used to learn the quantization model for $G(\mathbf{X}_i)$, which is kept fixed once acquired.

Following [13,27,64,77], we use ResNet-18 [25] for image classification, where we set $G(\cdot)$ to be the first 15 convolutional and 3 downsampling layers, which have 6,455,872 parameters, and $F'(\cdot)$ to be the remaining 3 layers (2 convolutional and 1 fully connected), which have 5,233,640 parameters. These layers were chosen for memory efficiency in the quantization model with ResNet-18, and we show the memory efficiency trade-off in supplemental materials (Fig. S1).

5 Experiments: Image Classification

5.1 Comparison Models

While REMIND learns on a per sample basis, most methods for incremental learning in CNNs do multiple loops through a batch. For fair comparison, we train these methods in the streaming setting to fairly compare against REMIND. Results for the incremental batch setting for these models are included in Fig. 1 and supplemental materials (Table S2 and Fig. S2-S3). We evaluate the following:

- REMIND Our main REMIND version uses PQ and replay augmentation.
 We also explore a version that omits data augmentation and a version that uses k-means rather than PQ.
- Fine-Tuning (No Buffer) Fine-Tuning is a baseline that fine-tunes θ_F of a CNN one sample at a time with a single epoch through the dataset. This approach does not use a buffer and suffers from catastrophic forgetting [45].
- ExStream Like REMIND, ExStream is a streaming learning method, however, it can only train fully connected layers of the network [22]. ExStream uses rehearsal by maintaining buffers of prototypes. It stores the input vector and combines the two nearest vectors in the buffer. After the buffer gets updated, all samples from its buffer are used to train the fully connected layers of a network. We use ExStream to train the final layer of the network, which is the only fully connected layer in ResNet-18.
- SLDA Streaming Linear Discriminant Analysis (SLDA) is a well-known streaming method that was shown to work well on deep CNN features [23]. It maintains running means for each class and a running tied covariance matrix. Given a new input, it assigns the label of the closest Gaussian in feature space. It can be used to compute the output layer of a CNN.

- iCaRL iCaRL is an incremental batch learning algorithm for CNNs [64].
 iCaRL stores images from earlier classes for replay, uses a distillation loss to preserve weights, and uses a nearest class mean classifier in feature space.
- Unified The Unified Classifier builds on iCaRL by using the outputs from the network for classification and introducing a cosine normalization layer, a constraint to preserve class geometry, and a margin ranking loss to maximize inter-class separation [27]. Unified also uses replay and distillation.
- BiC The Bias Correction (BiC) method builds on iCaRL by using the output layer of the network for classification and correcting the bias from class imbalance during training, i.e., more new samples than replay samples [77]. The method trains two additional bias correction parameters on the output layer, resulting in improved performance over distillation and replay alone.
- Offline The offline model is trained in a traditional, non-streaming setting and serves as an upper-bound on performance. We train two variants: one with only θ_F plastic and one with both θ_F and θ_G plastic.

Our main experiments focus on comparing state-of-the-art methods on ImageNet and we provide additional comparisons in Sect. 7. Although iCaRL, Unified, and BiC are traditionally trained in the incremental batch paradigm, we conduct experiments with these models in the streaming paradigm for fair comparison against REMIND. To train these streaming variants, we set the number of epochs to 1 and the batch size to r + 1 instances to match REMIND.

5.2 Model Configurations

In our setup, all models are trained instance-by-instance and have no batch requirements, unless otherwise noted. Because methods can be sensitive to the order in which new data are encountered, all models receive examples in the same order. The same base CNN initialization procedure is used by all models. For ExStream and SLDA, after base initialization, the streaming learning phase is re-started from the beginning of the data stream. All of the parameters except the output layer are kept frozen for ExStream and SLDA, whereas only $G(\cdot)$ is kept frozen for REMIND. All other comparison models do not freeze any layers and incremental training commences with the first new data sample. All models, except SLDA, are trained using cross-entropy loss with stochastic gradient descent and momentum. More parameter settings are in supplemental materials.

5.3 Datasets, Data Orderings, and Metrics

We conduct experiments with ImageNet and CORe50 by dividing both datasets into batches. The first batch is used for base initialization. Subsequently, all models use the same batch orderings, but they are sequentially fed individual samples and they cannot revisit any instances in a batch, unless otherwise noted. For ImageNet, the models are evaluated after each batch on all trained classes. For CORe50, models are evaluated on all test data after each batch.

ImageNet ILSVRC-2012 [68] has 1000 categories each with 732–1300 training samples and 50 validation samples, which we use for testing. During the base initialization phase, the model is trained offline on a set of 100 randomly selected classes. Following [13,27,64,77], each incremental batch then contains 100 random classes, which are not contained within any other batch. We study class incremental (class iid) learning with ImageNet.

CORe50 [52] contains sequences of video frames, with one object in each frame. It has 10 classes, and each sequence is acquired with varied environmental conditions. CORe50 is ideal for evaluating streaming learners since it is naturally non-iid and requires agents to learn from temporally correlated video streams. For CORe50, we follow [22] and sample at 1 frame per second, obtaining 600 training images and 225 test images per class. We use the bounding box crops and splits from [52]. Following [22], we use four training orderings to test the robustness of each algorithm under different conditions: 1) iid, where each batch has a random subset of training images, 2) class iid, where each batch has all of the images from two classes, which are randomly shuffled, 3) instance, where each batch has temporally ordered images from 80 unique object instances, and 4) class instance, where each batch has all of the temporally ordered instances from two classes. All batches have 1200 images across all orderings. Since CORe50 is small, CNNs are first initialized with pre-trained ImageNet weights and then fine-tuned on a subset of 1200 samples for base initialization.

We use the $\Omega_{\rm all}$ metric [22, 24, 45] for evaluation, which normalizes incremental learning performance by offline performance: $\Omega_{\rm all} = \frac{1}{T} \sum_{t=1}^{T} \frac{\alpha_t}{\alpha_{\rm offline},t}$, where T is the total number of testing events, α_t is the accuracy of the model for test t, and $\alpha_{\rm offline}$, is the accuracy of the optimized offline learner for test t. If $\Omega_{\rm all} = 1$, then the incremental learner's performance matched the offline model. We use top-5 and top-1 accuracies for ImageNet and CORe50, respectively. Average accuracy results are in supplemental materials (Table S2-S3).

5.4 Results: ImageNet

For ImageNet, we use the pre-trained PyTorch offline model with 89.08% top-5 accuracy to normalize $\Omega_{\rm all}$. We allow the iCaRL, Unified, and BiC models to store 10,000 (224×224 uint8) raw pixel image prototypes in a replay buffer, which is equivalent to 1.51 GB in memory. This allows REMIND to store indices for 959665 examples in its replay buffer. We set r=50 samples. We study additional buffer sizes in Sect. 7. Results for incremental class learning on ImageNet are shown in Table 1 and a learning curve for all models is shown in Fig. 3. REMIND outperforms all other comparison models, with SLDA achieving the sec-

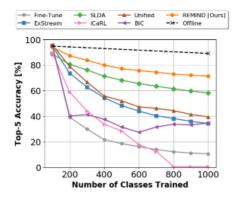


Fig. 3. Performance of streaming ImageNet models.

ond best performance. This is remarkable since REMIND only updates θ_F , whereas iCaRL, Unified, and BiC all update θ_F and θ_G .

REMIND is intended to be used for online streaming learning; however, we also created a variant suitable for incremental batch learning which is described in supplemental materials. Incremental batch results for REMIND and recent methods are given in Fig. 1 and supplemental materials (Table S2 and Fig. S2). While incremental batch methods train much more slowly, REMIND achieves comparable performance to the best methods.

Table 1. ResNet-18 streaming classification results on ImageNet and CORe50 using $\Omega_{\rm all}$. For CORe50, we explore performance across four ordering schemes and report the average of 10 permutations. Upper bounds are at the bottom.

Model	ImageNet	CORe50			
	CLS IID	IID	CLS IID	INST	CLS INST
Fine-Tune (θ_F)	0.288	0.961	0.334	0.851	0.334
ExStream	0.569	0.953	0.873	0.933	0.854
SLDA	0.752	0.976	0.958	0.963	0.959
iCaRL	0.306		0.690		0.644
Unified	0.614	_	0.510	_	0.527
BiC	0.440	_	0.410	_	0.415
REMIND	0.855	0.985	0.978	0.980	0.979
Offline (θ_F)	0.929	0.989	0.984	0.985	0.985
Offline	1.000	1.000	1.000	1.000	1.000

5.5 Results: CORe50

We use the CoRe50 dataset to study models under more realistic data orderings. Existing methods including iCaRL, Unified, and BiC assume that classes from one batch do not appear in other batches, making it difficult for them to learn the iid and instance orderings without modifications. To compute $\Omega_{\rm all}$, we use an offline model that obtains 93.11% top-1 accuracy. The iCaRL, Unified, and BiC models use replay budgets of 50 images, which is equivalent to 7.3 MB. This allows REMIND to store replay indices for 4465 examples. Results for other buffer sizes are in supplemental materials (Fig. S3). REMIND replays r=20 samples. $\Omega_{\rm all}$ results for CORe50 are provided in Table 1. For CORe50, REMIND outperforms all models for all orderings. In fact, REMIND is only 2.2% below the full offline model in the worst case, in terms of $\Omega_{\rm all}$. Methods that only trained the output layer performed well on CORe50 and poorly on ImageNet. This is likely because the CNNs used for CORe50 experiments are initialized with ImageNet weights, resulting in more robust representations. REMIND's remarkable performance on these various orderings demonstrate its versatility.

6 Experiments: Incremental VQA

In VQA, a system must produce an answer to a natural language question about an image [8,36,54], which requires capabilities such as object detection, scene understanding, and logical reasoning. Here, we use REMIND to pioneer streaming VQA. During training, a streaming VQA model receives a sequence of temporally ordered triplets $\mathcal{D} = \{(X_t, Q_t, A_t)\}_{t=1}^T$, where X_t is an image, Q_t is the question (string), and A_t is the answer. If an answer is not provided at time t, then the agent must use knowledge from time 1 to t-1 to predict A_t . To use REMIND for streaming VQA, we store each quantized feature along with a question string and answer, which can later be used for replay. REMIND can be used with almost any existing VQA system (e.g., attention-based [6,46,78], compositional [7,28], bi-modal fusion [10,19,71]) and it can be applied to similar tasks like image captioning [12] and referring expression recognition [43,63,67].

6.1 Experimental Setup

For our experiments, we use the TDIUC [35] and CLEVR [33] VQA datasets. TDIUC is composed of natural images and has over 1.7 million QA pairs organized into 12 question types including simple object recognition, complex counting, positional reasoning, and attribute classification. TDIUC tests for generalization across different underlying tasks required for VQA. CLEVR consists of over 700000 QA pairs for 70000 synthetically generated images and is organized into 5 question types. CLEVR specifically tests for multi-step compositional reasoning that is very rarely encountered in natural image VQA datasets. We combine REMIND with two popular VQA algorithms, using a modified version of the stacked attention network (SAN) [42,78] for TDIUC, and a simplified version of the Memory Attention and Control (MAC) [28,55] network for CLEVR. A ResNet-101 model pre-trained on ImageNet is used to extract features for both TDIUC and CLEVR. REMIND's PQ model is trained with 32 codebooks each of size 256. The final offline mean per-type accuracy with SAN on TDIUC is 67.59% and the final offline accuracy with MAC on CLEVR is 94.00%. Our main results with REMIND use a buffer consisting of 50% of the dataset and r = 50. Results for other buffer sizes are in supplemental materials (Table S4).

For both datasets, we explore two orderings of the training data: iid and question type (q-type). For iid, the dataset is randomly shuffled and the model is evaluated on all test data when multiples of 10% of the total training set are seen. The q-type ordering reflects a more interesting scenario where QA pairs for different VQA 'skills' are grouped together. Models are evaluated on all test data at the end of each q-type. We perform base initialization by training on the first 10% of the data for the iid ordering and on QA pairs belonging to the first q-type for the q-type ordering. Then, the remaining data is streamed into the model one sample at a time. The buffer is then incrementally updated with PQ encoded features and raw question strings. We use simple accuracy for CLEVR and mean-per-type accuracy for TDIUC.

We compare REMIND to ExStream [22], SLDA [23], an offline baseline, and a simple baseline where models are fine-tuned without a buffer, which causes catastrophic forgetting. To adapt ExStream and SLDA for VQA, we use a variant of the linear VQA model in [34], which concatenates ResNet-101 image features to question features extracted from a universal sentence encoder [73] and then trains a linear classifier. Parameter settings are in supplemental materials.

6.2 Results: VQA

Streaming VQA results for REMIND with a 50% buffer size are given in Table 2. Variants of REMIND with other buffer sizes are in supplemental materials (Table S4). REMIND outperforms the streaming baselines for both datasets, with strong performance on both TDIUC using the SAN model and CLEVR using the MAC model. Interestingly, for CLEVR the results are much greater for q-type

Table 2. $\Omega_{\rm all}$ results for streaming VQA.

Ordering	TDIUC		CLEVR		
	IID	Q-TYPE	IID	Q-TYPE	
Fine-Tune	0.716	0.273	0.494	0.260	
${\bf ExStream}$	0.676	0.701	0.477	0.375	
SLDA	0.624	0.644	0.518	0.496	
REMIND	0.917	0.919	0.720	0.985	
Offline	1.000	1.000	1.000	1.000	

than for iid. We hypothesize that the q-type ordering may be acting as a natural curriculum [11], allowing our streaming model to train more efficiently. Our results demonstrate that it is possible to train complex, multi-modal agents capable of attention and compositional reasoning in a streaming manner. Learning curves and qualitative examples are in supplemental materials (Fig. S5-S6).

7 Additional Classification Experiments

In this section, we study several of REMIND's components. In supplemental materials, we study other factors that influence REMIND's performance (Fig. S4), e.g., where to quantize, number of codebooks, codebook size, and replay samples (r). In supplemental materials, we also explore the performance of iCaRL, Unified, and BiC when only θ_F is updated (Sec. S3.2).

REMIND Components. REMIND is impacted by the size of its overall buffer, using augmentation, and the features used to train $F(\cdot)$. We study these on ImageNet and results are given in Table 3. REMIND (Main) denotes the variant of REMIND from our main experiments that uses augmentation with a buffer size of 959665 and 32 codebooks of size 256. PQ is critical to performance, with PQ (32 codebooks)

Table 3. REMIND variations on ImageNet with their memory (GB).

Variant	$\Omega_{ m all}$	Метогу
REMIND (Main)	0.855	1.51
100%Buffer	0.856	2.01
No Augmentation	0.818	1.51
k-Means	0.778	0.12
Real Features	0.868	24.08

outperforming k-means (1 codebook) by 7.7% in terms of $\Omega_{\rm all}$. Augmentation

is the next most helpful component and improves performance by 3.7%. Storing the entire dataset (100% Buffer) does not yield significant improvements. Using real features yields marginal improvements (1.3%) while requiring nearly 16 times more memory.

Replay Buffer Size. Since REMIND and several other models rely on a replay buffer to mitigate forgetting, we studied performance on ImageNet as a function of buffer size. We compared the performance of iCaRL, Unified, and BiC on ImageNet at three different buffer sizes (5K exemplars = 0.75 GB, 10K exemplars = 1.51 GB, and 20 K exemplars = 3.01 GB). To make the experiment fair, we compared REMIND to these models at equivalent buffer sizes, i.e., 479665 compressed samples = 0.75 GB, 959665 compressed samples = 1.51 GB, and 1281167 compressed samples (full dataset) = 2.01 GB. In

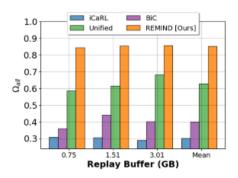


Fig. 4. $\Omega_{\rm all}$ as a function of buffer size for streaming ImageNet models.

Fig. 4, we see that more memory generally results in better performance. Overall, REMIND has the best performance and is nearly unaffected by buffer size. A plot with incremental batch models is in supplemental materials (Fig. S2), and follows the same trend: larger buffers yield better performance.

Regularization Comparisons. In Table 4, we show the results of REMIND and regularization methods for combating catastrophic forgetting on CORe50 class orderings. These regularization methods constrain weight updates to remain close to their previous values and are trained on batches of data, where each batch resembles a *task*. At test time, these models are provided with task labels, denoting which task an unseen sample came from. In our experiments, a *task* consists of several

Table 4. $\Omega_{\rm all}$ for regularization models averaged over 10 runs on CORe50 with and without Task Labels (TL).

Model	CLS IID		CLS INST		
	TL	No TL	TL	No TL	
SI	0.895	0.417	0.905	0.416	
EWC	0.893	0.413	0.903	0.413	
MAS	0.897	0.415	0.905	0.421	
RWALK	0.903	0.410	0.912	0.417	
A- GEM	0.925	0.417	0.916	0.421	
REMIND	0.995	0.978	0.995	0.979	
Offline	1.000	1.000	1.000	1.000	

classes, and providing task labels makes classification easier. We analyze performance when task labels are provided and when they are withheld. To evaluate REMIND and Offline with task labels, we mask off probabilities during test time for classes not included in the specific task. Consistent with [14,17,45], we find that regularization methods perform poorly when no task labels are provided. Regardless, REMIND outperforms all comparisons, both with and without task labels.

8 Discussion and Conclusion

We proposed REMIND, a brain-inspired replay-based approach to online learning in a streaming setting. REMIND achieved state-of-the-art results for object classification. Unlike iCaRL, Unified, and BiC, REMIND can be applied to iid and instance ordered data streams without modification. Moreover, we showed that REMIND is general enough for tasks like VQA with almost no changes.

REMIND replays compressed (lossy) representations that it stores, rather than veridical (raw pixel) experience, which is more consistent with memory consolidation in the brain. REMIND's replay is more consistent with how replay occurs in the brain during waking hours. Replay also occurs in the brain during slow wave sleep [9,31], and it would be interesting to explore how to effectively create a variant that utilizes sleep/wake cycles for replay. This could be especially beneficial for a deployed agent that is primarily engaged in online learning during certain hours, and is engaged in offline consolidation in other hours.

Several algorithmic improvements could be made to REMIND. We initialized REMIND's quantization model during the base initialization phase. For deployed, on-device learning this could instead be done by pre-training the codebook on a large dataset, or it could be initialized with large amounts of unlabeled data, potentially leading to improved representations. Another potential improvement is using selective replay. REMIND randomly chooses replay instances with uniform probability. In early experiments, we also tried choosing replay samples based on distance from current example, number of times a sample has been replayed, and the time since it was last replayed. While none performed better than uniform selection, we believe that selective replay still holds the potential to lead to better generalization with less computation. Because several comparison models used ResNet-18, we also used ResNet-18 for image classification so that we could compare against these models directly. The ResNet-18 layer used for quantization was chosen to ensure REMIND's memory efficiency, but co-designing the CNN architecture with REMIND could lead to considerably better results. Using less memory, REMIND stores far more compressed representations than competitors. For updating the replay buffer, we used random replacement, which worked well in [14,77]. We tried a queue and a distance-based strategy, but both performed nearly equivalent to random selection with higher computational costs. Furthermore, future variants of REMIND could incorporate mechanisms similar to [62] to explicitly account for the temporal nature of incoming data. To demonstrate REMIND's versatility, we pioneered streaming VQA and established strong baselines. It would be interesting to extend this to streaming chart question answering [37,38,40], object detection, visual query detection [2], and other problems in vision and language [39].

Acknowledgements. This work was supported in part by the DARPA/MTO Lifelong Learning Machines program [W911NF-18-2-0263], AFOSR grant [FA9550-18-1-0121], NSF award #1909696, and a gift from Adobe Research. We thank NVIDIA for the GPU donation. The views and conclusions contained herein are those of the authors and should not be interpreted as representing the official policies or endorsements of

any sponsor. We thank Michael Mozer, Ryne Roady, and Zhongchao Qian for feedback on early drafts of this paper.

References

- Abraham, W.C., Robins, A.: Memory retention-the synaptic stability versus plasticity dilemma. Trends Neurosci. 28, 73–78 (2005)
- Acharya, M., Jariwala, K., Kanan, C.: VQD: visual query detection in natural scenes. In: NAACL (2019)
- 3. Aljundi, R., Babiloni, F., Elhoseiny, M., Rohrbach, M., Tuytelaars, T.: Memory aware synapses: learning what (not) to forget. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11207, pp. 144–161. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01219-9_9
- 4. Aljundi, R., et al.: Online continual learning with maximal interfered retrieval. In: NeurIPS, pp. 11849–11860 (2019)
- 5. Aljundi, R., Lin, M., Goujaud, B., Bengio, Y.: Gradient based sample selection for online continual learning. In: NeurIPS, pp. 11816–11825 (2019)
- Anderson, P., et al.: Bottom-up and top-down attention for image captioning and visual question answering. In: CVPR (2018)
- Andreas, J., Rohrbach, M., Darrell, T., Klein, D.: Neural module networks. In: CVPR, pp. 39–48 (2016)
- 8. Antol, S., et al.: VQA: visual question answering. In: ICCV (2015)
- Barnes, D.C., Wilson, D.A.: Slow-wave sleep-imposed replay modulates both strength and precision of memory. J. Neurosci. 34(15), 5134–5142 (2014)
- Ben-Younes, H., Cadene, R., Cord, M., Thome, N.: Mutan: multimodal tucker fusion for visual question answering. In: ICCV (2017)
- Bengio, Y., Louradour, J., Collobert, R., Weston, J.: Curriculum learning. In: ICML, pp. 41–48 (2009)
- Bernardi, R., et al.: Automatic description generation from images: a survey of models, datasets, and evaluation measures. J. Artif. Intell. Res. 55, 409–442 (2016)
- Castro, F.M., Marín-Jiménez, M.J., Guil, N., Schmid, C., Alahari, K.: End-to-End incremental learning. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11216, pp. 241–257. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01258-8_15
- Chaudhry, A., Dokania, P.K., Ajanthan, T., Torr, P.H.S.: Riemannian walk for incremental learning: understanding forgetting and intransigence. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11215, pp. 556–572. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01252-6_33
- 15. Chaudhry, A., Ranzato, M., Rohrbach, M., Elhoseiny, M.: Efficient lifelong learning with A-GEM. In: ICLR (2019)
- Dhar, P., Singh, R.V., Peng, K.C., Wu, Z., Chellappa, R.: Learning without memorizing. In: CVPR, pp. 5138–5146 (2019)
- Farquhar, S., Gal, Y.: Towards robust evaluations of continual learning. arXiv:1805.09733 (2018)
- Fernando, C., et al.: Pathnet: evolution channels gradient descent in super neural networks. arXiv:1701.08734 (2017)
- Fukui, A., Park, D.H., Yang, D., Rohrbach, A., Darrell, T., Rohrbach, M.: Multi-modal compact bilinear pooling for visual question answering and visual grounding. In: EMNLP (2016)

- 20. Gama, J.: Knowledge Discovery from Data Streams. Chapman and Hall/CRC, Boca Raton (2010)
- Gama, J., Sebastião, R., Rodrigues, P.P.: On evaluating stream learning algorithms. Mach. Learn. 90(3), 317–346 (2013). https://doi.org/10.1007/s10994-012-5320-9
- Hayes, T.L., Cahill, N.D., Kanan, C.: Memory efficient experience replay for streaming learning. In: ICRA (2019)
- Hayes, T.L., Kanan, C.: Lifelong machine learning with deep streaming linear discriminant analysis. In: CVPRW (2020)
- 24. Hayes, T.L., Kemker, R., Cahill, N.D., Kanan, C.: New metrics and experimental paradigms for continual learning. In: CVPRW, pp. 2031 2034 (2018)
- He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
- Hou, S., Pan, X., Loy, C.C., Wang, Z., Lin, D.: Lifelong learning via progressive distillation and retrospection. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11207, pp. 452 467. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01219-9_27
- Hou, S., Pan, X., Wang, Z., Change Loy, C., Lin, D.: Learning a unified classifier incrementally via rebalancing. In: CVPR (2019)
- Hudson, D.A., Manning, C.D.: Compositional attention networks for machine reasoning. In: ICLR (2018)
- Insausti, R., et al.: The nonhuman primate hippocampus: neuroanatomy and patterns of cortical connectivity. In: Hannula, D.E., Duff, M.C. (eds.) The Hippocampus from Cells to Systems, pp. 3–36. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-50406-3_1
- Jegou, H., Douze, M., Schmid, C.: Product quantization for nearest neighbor search. TPAMI 33(1), 117–128 (2010)
- Ji, D., Wilson, M.A.: Coordinated memory replay in the visual cortex and hippocampus during sleep. Nat. Neurosci. 10(1), 100–107 (2007)
- 32. Johnson, J., Douze, M., Jégou, H.: Billion-scale similarity search with GPUs. IEEE Trans. Big Data (2019)
- Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, C.L., Girshick,
 R.: Clevr: a diagnostic dataset for compositional language and elementary visual reasoning. In: CVPR (2017)
- 34. Kafle, K., Kanan, C.: Answer-type prediction for visual question answering. In: CVPR, pp. 4976–4984 (2016)
- Kafle, K., Kanan, C.: An analysis of visual question answering algorithms. In: ICCV, pp. 1983–1991 (2017)
- Kafle, K., Kanan, C.: Visual question answering: datasets, algorithms, and future challenges. Comput. Vis. Image Underst. 163, 3–20 (2017)
- Kafle, K., Price, B., Cohen, S., Kanan, C.: DVQA: understanding data visualizations via question answering. In: CVPR, pp. 5648–5656 (2018)
- Kafle, K., Shrestha, R., Cohen, S., Price, B., Kanan, C.: Answering questions about data visualizations using efficient bimodal fusion. In: WACV, pp. 1498–1507 (2020)
- Kafle, K., Shrestha, R., Kanan, C.: Challenges and prospects in vision and language research. Front. Artif. Intell. 2, 28 (2019)
- Kahou, S.E., Michalski, V., Atkinson, A., Kádár, Á., Trischler, A., Bengio, Y.: Figureqa: An annotated figure dataset for visual reasoning. arXiv preprint arXiv:1710.07300 (2017)
- Karlsson, M.P., Frank, L.M.: Awake replay of remote experiences in the hippocampus. Nat. Neurosci. 12(7), 913 (2009)

- Kazemi, V., Elqursh, A.: Show, ask, attend, and answer: a strong baseline for visual question answering. arXiv:1704.03162 (2017)
- Kazemzadeh, S., Ordonez, V., Matten, M., Berg, T.: Referitgame: Referring to objects in photographs of natural scenes. In: EMNLP, pp. 787–798 (2014)
- Kemker, R., Kanan, C.: FearNet: brain-inspired model for incremental learning. In: ICLR (2018)
- Kemker, R., McClure, M., Abitino, A., Hayes, T.L., Kanan, C.: Measuring catastrophic forgetting in neural networks. In: AAAI (2018)
- Kim, J.H., Jun, J., Zhang, B.T.: Bilinear attention networks. In: NeurIPS, pp. 1564–1574 (2018)
- Kirkpatrick, J., et al.: Overcoming catastrophic forgetting in neural networks. In: PNAS (2017)
- 48. Konkel, A., Warren, D.E., Duff, M.C., Tranel, D., Cohen, N.J.: Hippocampal amnesia impairs all manner of relational memory. Front. Hum. Neurosci. 2, 15 (2008)
- Le, T., Stahl, F., Gaber, M.M., Gomes, J.B., Di Fatta, G.: On expressiveness and uncertainty awareness in rule-based classification for data streams. Neurocomputing 265, 127–141 (2017)
- Lee, K., Lee, K., Shin, J., Lee, H.: Overcoming catastrophic forgetting with unlabeled data in the wild. In: ICCV, pp. 312–321 (2019)
- Lewis, P.A., Durrant, S.J.: Overlapping memory replay during sleep builds cognitive schemata. Trends Cogn. Sci. 15(8), 343–351 (2011)
- Lomonaco, V., Maltoni, D.: Core50: a new dataset and benchmark for continuous object recognition. In: CoRL, pp. 17–26 (2017)
- Lopez-Paz, D., Ranzato, M.: Gradient episodic memory for continual learning. In: NeurIPS, pp. 6467-6476 (2017)
- Malinowski, M., Fritz, M.: A multi-world approach to question answering about real-world scenes based on uncertain input. In: NeurIPS (2014)
- 55. Marois, V., Jayram, T., Albouy, V., Kornuta, T., Bouhadjar, Y., Ozcan, A.S.: On transfer learning using a mac model variant. arXiv:1811.06529 (2018)
- McClelland, J.L., Goddard, N.H.: Considerations arising from a complementary learning systems perspective on hippocampus and neocortex. Hippocampus 6(6), 654–665 (1996)
- 57. McCloskey, M., Cohen, N.J.: Catastrophic interference in connectionist networks: the sequential learning problem. Psychol. Learn. Motiv. 24, 109–165 (1989)
- Nguyen, C.V., Li, Y., Bui, T.D., Turner, R.E.: Variational continual learning. In: ICLR (2018)
- Ostapenko, O., Puscas, M., Klein, T., Jähnichen, P., Nabi, M.: Learning to remember: a synaptic plasticity driven framework for continual learning. In: CVPR (2019)
- O'Neill, J., Pleydell-Bouverie, B., Dupret, D., Csicsvari, J.: Play it again: reactivation of waking experience and memory. Trends Neurosci. 33(5), 220–229 (2010)
- 61. Parisi, G.I., Kemker, R., Part, J.L., Kanan, C., Wermter, S.: Continual lifelong learning with neural networks: a review. Neural Netw. 113, 54–71 (2019)
- Parisi, G.I., Tani, J., Weber, C., Wermter, S.: Lifelong learning of spatiotemporal representations with dual-memory recurrent self-organization. Front. Neurorobot. 12, 78 (2018)
- Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: collecting region-to-phrase correspondences for richer image-to-sentence models. In: ICCV, pp. 2641–2649 (2015)
- Rebuffi, S.A., Kolesnikov, A., Sperl, G., Lampert, C.H.: icarl: incremental classifier and representation learning. In: CVPR (2017)

- 65. Riemer, M., et al.: Learning to learn without forgetting by maximizing transfer and minimizing interference. In: ICLR (2019)
- Ritter, H., Botev, A., Barber, D.: Online structured Laplace approximations for overcoming catastrophic forgetting. In: NeurIPS, pp. 3738–3748 (2018)
- Rohrbach, A., Rohrbach, M., Hu, R., Darrell, T., Schiele, B.: Grounding of textual phrases in images by reconstruction. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) ECCV 2016. LNCS, vol. 9905, pp. 817–834. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46448-0_49
- Russakovsky, O., et al.: ImageNet large scale visual recognition challenge. IJCV 115(3), 211–252 (2015). https://doi.org/10.1007/s11263-015-0816-y
- 69. Rusu, A.A., et al.: Progressive neural networks. arXiv:1606.04671 (2016)
- Serra, J., Suris, D., Miron, M., Karatzoglou, A.: Overcoming catastrophic forgetting with hard attention to the task. In: ICML, pp. 4555–4564 (2018)
- Shrestha, R., Kafle, K., Kanan, C.: Answer them all! toward universal visual question answering models. In: CVPR (2019)
- 72. Stickgold, R., Hobson, J.A., Fosse, R., Fosse, M.: Sleep, learning, and dreams: off-line memory reprocessing. Science **294**(5544), 1052–1057 (2001)
- Subramanian, S., Trischler, A., Bengio, Y., Pal, C.J.: Learning general purpose distributed sentence representations via large scale multi-task learning. In: ICLR (2018)
- Takahashi, S.: Episodic-like memory trace in awake replay of hippocampal place cell activity sequences. Elife 4, e08105 (2015)
- 75. Teyler, T.J., Rudy, J.W.: The hippocampal indexing theory and episodic memory: updating the index. Hippocampus 17(12), 1158–1169 (2007)
- 76. Verma, V., et al.: Manifold mixup: better representations by interpolating hidden states. In: ICML (2019)
- 77. Wu, Y., et al.: Large scale incremental learning. In: CVPR, pp. 374–382 (2019)
- 78. Yang, Z., He, X., Gao, J., Deng, L., Smola, A.J.: Stacked attention networks for image question answering. In: CVPR (2016)
- 79. Yoon, J., Yang, E., Lee, J., Hwang, S.J.: Lifelong learning with dynamically expandable networks. In: ICLR (2018)
- 80. Yosinski, J., Clune, J., Bengio, Y., Lipson, H.: How transferable are features in deep neural networks? In: NeurIPS, pp. 3320–3328 (2014)
- 81. Zenke, F., Poole, B., Ganguli, S.: Continual learning through synaptic intelligence. In: ICML, pp. 3987–3995 (2017)