

Models for Continual Learning

monaco, 2019), of these methods, towards continuing CL approaches (generalized probability) and trained models. In fact, we argue that training an unlabeled model can significantly improve classification as learning is treated as a regression task. Training the energy of an input sample while the energy of unlabeled samples is increased. An unlabeled sample offers freedom to continual learning and the negative reduces the energy when a new input, unlabeled, is being added.

Given a multi-label pair to a model to select the most important towards the unlabeled sample, it enables EBMs to learn from previous knowledge by learning a new energy function which serves as an auxiliary information. In the same way, the energy function can be learned.

In this paper, we address two main challenges in learning. 1) First, we propose incremental learning strategies for continually learning. 2) Secondly, successively learning either from labeled or unlabeled data has disadvantages in terms of efficiency. We propose a continual learning strategy based on stored data. This strategy can address issues such as memory usage. Typically, a continual learning system needs to store all the data of distinct tasks. This is a major problem in learning from unlabeled data. In this paper, we propose a continual learning strategy that does not require storing unlabeled data. This strategy can be naturally applied to the problems of learning from unlabeled data.

There are three main contributions. First, we introduce energy-based models for continual learning problems. We show that they work well with challenging problems in an agnostic setting and class-invariant replay. Secondly, we propose an objective that is simple and general for different types of models, with significant performance. This contrastive divergence can naturally handle the dynamics of classes and significantly reduce the number of epochs. Lastly, we show that the proposed approach performs well on four standard CL benchmarks: MNIST, CIFAR-10, and CIFAR-100. Our results point towards EBMs as a class of models that are robust towards the CL regime and the choice upon which to build further research. The code is made public to facilitate further research.

tions of our work. First, we models for classification continual learning that EBMs can naturally deal with CL, including the boundary-incremental learning without us ing an energy-based training that is broadly applicable to different significant boosts on their performance. The divergence based training objective with a dynamically growing number of hidden units reduces catastrophic forgetting. Proposed EBMs perform strongly on benchmarks: split MNIST, permuted MNIST, and FAR-100. These observations suggest that EBMs are less sensitive to the choice of models naturally inclined to catastrophic forgetting. This work is an important new baseline for future developments. The code are available at model.github.io/.
for Continual Learning.

other important distinction in CL is between incremental learning (*Task-IL*) and class-incremental learning (*Class-IL*) (van de Ven & Tolias, 2020). In *Task-IL*, also referred to as the *task-incremental* versus *class-incremental* (Arulkumar & Gal, 2018), models predict the data by choosing only from the labels in the data come from. In *Class-IL*, also referred to as the *class-incremental* setting, models chose between the labels so far when asked to predict the label. *Class-IL* is more challenging than *Task-IL* because it requires the model to select the correct labels from the previously learned classes. Generally, to perform well on *Class-IL*, methods need to store data, use replay, or learn another large dataset (Rebuffi et al., 2019; Belouadah et al., 2020; Ma et al., 2019; Hayes & Kanan, 2020).

2. Continual learning approaches

Numerous methods have been proposed for CL. Here we broadly partition them into task-specific, regularization, and replay-based methods.

Task-specific methods.

One way to reduce forgetting between tasks is by using different parts of the network for different tasks. For a fixed-size network, this could be achieved by learning a separate mask (Fernando et al., 2017; Serra et al., 2018) defining a different, random mask for each task (Zeng et al., 2019; Hu et al., 2019). This lets models grow or recruit new resources for new tasks, such as progressive neural networks (PNNs) (Kirkpatrick et al., 2016) and dynamically expandable networks (DENs) (Xie et al., 2017). Although these methods are generally good at reducing catastrophic forgetting, a key downside is that they require knowledge of task identities during learning. They are therefore not suitable for open-set learning.

Regularization-based methods.

Regularization-based methods seek to encourage the stability of those aspects of the model that are important for previous tasks. A common way to do this is to add a regularization loss to penalize the values of certain parameters. EWC (Kirkpatrick et al., 2017) and online EWC (Schwarz et al., 2018) evaluate the values of model parameters using the diagonal elements of the Hessian matrix. The Hessian matrix is a square matrix of second-order partial derivatives of the loss function with respect to the model parameters. The diagonal elements of the Hessian matrix represent the second-order approximation of the loss function with respect to the model parameters. The off-diagonal elements represent the interactions between the model parameters. By regularizing the values of the model parameters, EWC and online EWC prevent the model from changing too much when learning new tasks. This helps to reduce catastrophic forgetting.

incremental learning. It is between task-incremental learning (Prabhu et al., 2019; Rajasegaran et al., 2017) and multi-head setting (Lomonaco & Itoni, 2018). In task-incremental learning, the label of an input data is the label of the task where the data belongs to. This is referred to as the single-task setting, as it only uses one classes from all the tasks to predict the label of an input data. Multi-head setting is similar, as it requires modeling multiple heads to mixtures of new and old data. In contrast to *Class-II*, existing work on continual learning either pretrain models (Lomonaco & Itoni, 2018; Rajasegaran et al., 2017) or fine-tune them (Lomonaco & Itoni, 2018). The former approach is based on the assumption that the model has been trained on all the data before it is used for continual learning. The latter approach is based on the assumption that the model has been trained on all the data before it is used for continual learning. Both approaches have their own advantages and disadvantages. The former approach is more efficient, as it does not require retraining the model on all the data. The latter approach is more effective, as it can handle changes in the data distribution more easily. However, both approaches have their own disadvantages. The former approach is less effective, as it cannot handle changes in the data distribution. The latter approach is less efficient, as it requires retraining the model on all the data.

Replay methods. To prevent catastrophic forgetting, one approach is to periodically rehearse previous tasks (Robins, 1995). Exact or approximate replay can store data from previous tasks and reuse them for learning on new tasks. Although such methods have been shown to mitigate critical non-trivial questions about what information needs to be stored and how to do so (Kirkpatrick et al., 2017; Hou et al., 2019; Vinyals et al., 2017), they are computationally expensive. An alternative is to generate samples of previous tasks (Leibo et al., 2017), a generative model that can mitigate forgetting, an idea that has been shown to work well in training a generative model (Lesort et al., 2019; van Hasselt et al., 2019).

In contrast, we propose a simple method to reduce catastrophic forgetting while maintaining task-identity, without requiring significant changes to the learning capabilities, and without using a generative model.

3. CL with softmax

The most common way to train neural networks is to use a softmax-based classifier (SBC) (Goodfellow et al., 2016; Tolias et al., 2019; Pellegrini et al., 2019). In this section, we start by reviewing softmax-based classifiers and then introduce a family of softmax-based classifiers that can mitigate catastrophic forgetting.

3.1. Softmax-based classifiers

Given an input $\mathbf{x} \in \mathbb{R}^D$, a softmax-based classifier defines the conditional probability of N possible class labels as

$$p_{\theta}(y|\mathbf{x}) = \frac{\exp([f_{\theta}(\mathbf{x})]_y)}{\sum_{i \in \mathcal{Y}} \exp([f_{\theta}(\mathbf{x})]_i)},$$

where $f_{\theta}(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^N$ is a function that takes an input \mathbf{x} and returns a N -dimensional vector of log-odds ratios, also known as a vector of scores. A schema of SBC is shown in Figure 1(a).

reinforce knowledge, replay methods store previous information during training and experience replay based methods can re-learn tasks and revisit them when training is straightforward, such methods face challenges, such as how to select the data to use them (Lopez-Paz & Ranzato, 2019; Mundt et al., 2020). Replay the replayed data. In (Shin et al., 2017) it is sequentially trained to generate samples. While both types of replay can be important disadvantage is that they are relatively expensive. Additionally, storage may be possible while incrementally learning is a challenging problem in itself (van de Ven et al., 2020).

EBMs for continual learning that forget without requiring knowledge are gradually restricting the model's recall without using stored data.

-based classifiers

To do classification with deep neural networks we use softmax output layer in combination with cross-entropy loss function.

In continual learning, virtually all classification approaches are based on the softmax-max approach (Li & Hoiem, 2017; van de Ven & Welling, 2019; Zenke et al., 2017). We describe classification with a simple example. We then describe the challenges of classification in CL.

Classification

Given a discrete set $\mathcal{Y} = \{1, \dots, N\}$ and a continuous set \mathcal{X} , a traditional softmax-based classifier outputs the conditional probabilities of those labels as:

$$\text{softmax}(\theta(\mathbf{x}))_y = \frac{\exp(\theta_y(\mathbf{x}))}{\sum_i \exp(\theta_i(\mathbf{x}))}, \quad \text{for all } y \in \mathcal{Y}, \quad (1)$$

$\theta: \mathbb{R}^N \rightarrow \mathbb{R}^N$ is a feed-forward neural network, that maps an input \mathbf{x} to a N -dimensional vector. $[.]_i$ indicates the i th element of a vector. A 3-class classification task is shown in Figure 1 left.

Classifiers that are typically trained

used for class-incremental learning, classifiers face several challenges. One is softmax-based classifiers compute over all classes (or sometimes over been seen so far). As a result, when the likelihood of the currently observed but the likelihood of old classes is too since they are not encountered in the operation introduces competitive, which make the classifier catastrophic. We show such phenomenon of SBC.

4. CL with Energy-Based Models

In this section, we propose a simple based training objective that can successfully catastrophic forgetting in CL. We do in section 4.1 and then show how EBM classification in section 4.2 and continual learning.

4.1. Energy-based models

EBMs (LeCun et al., 2006) are a class of likelihood models that define the likelihood $x \in \mathcal{X} \subseteq \mathbb{R}^D$ using the Boltzmann distribution:

$$p_\theta(x) = \frac{\exp(-E_\theta(x))}{Z(\theta)}, \quad Z(\theta) = \sum_x \exp(-E_\theta(x)),$$

where $E_\theta(x) : \mathbb{R}^D \rightarrow \mathbb{R}$, known as the energy function. It maps each data point x to a scalar energy value. The denominator is the partition function. In deep learning, the energy function E_θ is a neural network.

EBMs are powerful models that have been applied to many different domains, such as structured prediction (McCallum, 2016; Gygli et al., 2017; Tu & Gimpel, 2019), machine translation (Sennrich et al., 2016), text generation (Deng et al., 2020), image captioning (Haarnoja et al., 2017), image generation (Hinton, 2009; Xie et al., 2016; 2018), memory modeling (Bartunov et al., 2018; Grathwohl et al., 2019), and biology (Scellier & Bengio, 2017). As far as we know, for continual learning has so far remained unexplored.

4.2. Energy-based models for classification

To solve the classification tasks, we apply the EBM framework to the incremental learning problem. Specifically, given a sequence of training data $\{(x_i, y_i)\}_{i=1}^n$, we want to learn a model that can correctly classify the new data while maintaining the performance on the old data. This is achieved by minimizing the following loss function:

$$\mathcal{L}(\theta) = -\sum_{i=1}^n \log p_\theta(x_i | y_i) + \lambda \sum_{i=1}^n \log p_\theta(x_i | \bar{y}_i),$$

where \bar{y}_i is a set of incorrect labels. The first term encourages the model to assign high probability to the correct label, while the second term discourages the model from assigning high probability to incorrect labels. The parameter λ controls the trade-off between the two terms.

softmax-based classifier. An important issue is that the cross-entropy loss favors all classes that have been trained on a new task, so the weight assigned to old classes is increased, which can lead to heavily suppressed weights for old tasks. The softmax dynamics also forget past tasks. This is discussed in Section 5.1.3.

Energy-based models

Energy-based models are simple but efficient energy-minimization models that successfully mitigate the overfitting problem. We first introduce EBMs and then show how EBMs are used for classification and learning in sections 4.3. and 5.1.3.

EBMs are a class of maximum likelihood models that model the probability distribution of a data point x given its energy distribution:

$$\exp(-E_\theta(x)) \quad (2)$$

where $E_\theta(x)$ is the energy function, θ is the energy parameter, and $Z(\theta)$ is the partition function. In learning applications, the energy function is often parameterized by θ .

EBMs have been applied to different tasks such as image prediction (Belanger & Brockschmidt, 2019), sequence generation (Tu et al., 2020), reinforcement learning (Salakhutdinov & Mnih, 2014; Nijkamp et al., 2019), and language modeling (Liu et al., 2019), classification (Krause et al., 2019), and generative modeling (Goodfellow et al., 2014). While EBMs have been applied to many tasks, they remain unexplored in the context of learning from unlabeled data.

Classification

We now adapt the above general framework to the classification task. We start by defining the energy function for a data point x given a class label y :

$$E_\theta(x, y) = -\log p_\theta(x|y) \quad (3)$$

where $p_\theta(x|y)$ is the conditional probability of observing data point x given class label y . The energy function is a measure of how well the data point x fits the class label y . The lower the energy, the more likely it is that the data point x belongs to the class y .

The energy function is composed of two parts: a feature-based part and a bias-based part. The feature-based part is a weighted sum of features of the data point x , where the weights are learned from the training data. The bias-based part is a constant term that shifts the energy function. The energy function is then normalized by the partition function $Z(\theta)$:

$$Z(\theta) = \sum_y \exp(-E_\theta(x, y)) \quad (4)$$

The final probability of observing data point x given class label y is then given by:

$$p_\theta(x|y) = \frac{\exp(-E_\theta(x, y))}{Z(\theta)} \quad (5)$$

Training. We want the model to learn the data distribution by maximizing the negative log likelihood:

$$\mathcal{L}_{\text{ML}}(\theta) = -\sum_{(x,y) \sim p_D} \log p(y|x; \theta)$$

with the expanded form:

$$\mathcal{L}_{\text{ML}}(\theta) = \mathbb{E}_{(x,y) \sim p_D} \log p(y|x; \theta)$$

Equation 5 minimizes the negative log likelihood loss, which maximizes the label y and minimizes the label \bar{y} . This is equivalent to decreasing the energy of the label y .

Inference. Given an input x , the inference of our EBMs is the class with the highest energy:

$$\hat{y} = \arg \max_y p(y|x; \theta)$$

4.3. Energy-based models

4.3.1. EBM TRAINING

We notice that in Equation 5, the labels y and \bar{y} given data x are mixed. This is because across all labels raised to the power of β , the higher the energy, the more likely the label is to be selected. This is a problem for a standard based classifier model, as it can lead to catastrophic forgetting (Kirkpatrick et al., 2016; Mordatch et al., 2019; Xie et al., 2019). To mitigate this problem and prevent catastrophic forgetting, we define the following loss function:

$$\mathcal{L}_{\text{CD}}(\theta; x, y) = \mathbb{E}_{(x,y) \sim p_D} \log p(y|x; \theta)$$

where y is the ground truth label and \bar{y} is the negative class label raised to the power of β . The labels in the current training set are mixed.

Different from the standard maximum likelihood estimation, we do not maximize likelihood nor minimize cross-entropy loss. Instead, by contrastive divergence, we update the parameters between the ground truth label y and the negative class label \bar{y} for a given data point. The negative class label \bar{y} is updated with previous classes \bar{y}_t to prevent catastrophic forgetting.

Importantly, such a self-supervised learning approach can be naturally applied to datasets with missing labels.

on function for normalization. The distribution defined by E_θ to model p_D , which we do by minimizing the loss of the data

$$= \mathbb{E}_{(x,y) \sim p_D} [-\log p_\theta(y|x)] \quad (4)$$

term:

$$\left[E_\theta(\mathbf{x}, y) + \log \left(\sum_{y' \in \mathcal{Y}} e^{-E_\theta(\mathbf{x}, y')} \right) \right]. \quad (5)$$

reduces the energy of \mathbf{x} at the ground truth label and increases the overall partition function by introducing energy of \mathbf{x} at other labels y' .

Given input \mathbf{x} , the class label predicted by the network is the one with the smallest energy at \mathbf{x} :

$$= \arg \min_{y' \in \mathcal{Y}} E_\theta(\mathbf{x}, y'), \quad (6)$$

models for continual learning

LEARNING OBJECTIVE

In Equation 5, the energy over all class labels is maximized. Directly maximizing energy over all classes causes the same problem as the softmax cross-entropy loss: that the old classes are suppressed when learning new classes and thus cause the model to forget. Inspired by (Hinton, 2002; Du & Hinton, 2012; Chen et al., 2016), we find that the contrastive approximation of Equation 5 can mitigate this problem and lead to a simpler equation. To do so, we introduce a contrastive divergence loss:

$$\mathbb{E}_{(x,y) \sim p_D} [E_\theta(\mathbf{x}, y) - E_\theta(\mathbf{x}, y^-)], \quad (7)$$

where y^- is a randomly sampled label from the set of classes in the training batch \mathcal{Y}_B such that $y^- \neq y$.

As a softmax-based classifier, EBMs maximize the probability of the correct label by normalizing over all classes but the correct one. By increasingly increasing the energy difference between the correct label and another negative label for each data point, this operation causes less interference between classes and enables EBMs to suffer less from forgetting.

The sampling strategy also allows our EBMs to learn without any modification to different parts of the network.

Figure 1: Schematic of the model-based classifier (SBC) and energy-based classifier (EBC). The SBC takes an image \mathbf{x} as input and outputs a fixed vector to represent the probabilities. The EBC takes a data \mathbf{x} and a class y as input and outputs a probability $m(\mathbf{x}, y)$. The dash lines are optional skip connections. The negative sample(s) from the training set do not require knowledge of the underlying tasks. In the application of EBMs in the *bouquet* dataset, which the underlying tasks are not explicitly labeled, we find that the proposed EBC is sufficient enough to achieve good performance on datasets (Table 1 and Table 7). It is also possible to use other strategies to partition classes in the partition function, where we explore alternative strategies to handle current training batch \mathcal{Y}_B as negative samples for all classes seen so far as negative samples in the EBM training. This provides more freedom for models to choose what to learn, which is important for preventing overfitting and continual learning.

4.3.2. ENERGY NETWORK

Another important difference from SBC is that the choice of model architecture is more flexible in EBMs. Traditional EBMs feed in \mathbf{x} as input. In contrast, EBC has two ways to combine \mathbf{x} and y in the network, with the only requirement that $E_\theta(\mathbf{x}, y)$ is differentiable, we can treat y as an attention filter to filter relevant information between \mathbf{x} and y .

To compute the energy of any image pair, we use y to influence a convolutional layer that serves as an attention filter (Xu et al., 2015). To capture the most relevant information between \mathbf{x} and y (right), we first send \mathbf{x} into a small linear feature $f(\mathbf{x})$. The label y is mapped into a feature $g(y)$ using a small learned linear projection. We use the gating mechanism to filter relevant information between \mathbf{x} and y :

$$m(\mathbf{x}, y) = G(f(\mathbf{x}), g(y))$$

The output is finally sent to a fully connected layer to obtain the energy value $E_\theta(\mathbf{x}, y)$. See Supplementary Figure 1 for the detailed architecture of the EBC.

architectures of the softmax-based models (EBM). SBC takes pre-defined N -dimensional of N different classes. EBM and outputs their energy value. connections.

current batch, our EBMs task on hand. This allows *boundary-agnostic* setting, in which boundaries are not given.

The training objective is performance on different CL. We note however that it is for choosing the negative in Equation 5. In Table 3, we compare two ways: 1) using all classes in the negative classes, and 2) using negative classes. The usage of training objective provides which classes to train on avoiding catastrophic forgetting in

In softmax-based classifiers architectures becomes more complex. Classification models only need to define an energy function with the mapping $f: (\mathbb{R}^D, \mathbb{N}) \rightarrow \mathbb{R}$. In EBMs, we use a linear or gate to select the most likely class y .

Given input data \mathbf{x} and class label y , we can obtain additional gain on \mathbf{x} , which is used by (Krahenbuhl et al., 2015) to select the most likely class between \mathbf{x} and y . In Figure 1 (a), we show a neural network to generate the feature map \mathbf{x} mapped into a same dimension feature map $\mathbf{g}(\mathbf{x})$. Then a learned network or a random weight block G to select the most likely class y :

$$f(\mathbf{x}, \mathbf{g}(\mathbf{x})). \quad (8)$$

Figure 1 (b) shows a neural network to generate the feature map \mathbf{x} mapped into a same dimension feature map $\mathbf{g}(\mathbf{x})$. Then a learned network or a random weight block G to select the most likely class y :

$$f(\mathbf{x}, \mathbf{g}(\mathbf{x})). \quad (8)$$

Figure 1 (c) shows a neural network to generate the feature map \mathbf{x} mapped into a same dimension feature map $\mathbf{g}(\mathbf{x})$. Then a learned network or a random weight block G to select the most likely class y :

$$f(\mathbf{x}, \mathbf{g}(\mathbf{x})). \quad (8)$$

any number of classes in new batches by defining a new conditional gain $g(y)$ and generating its energy value with our formulation gives us freedom to learn more-defining their number in advance. Standard classifiers have to change their model architecture when adding new class heads to the softmax output along with new classes. In contrast, EBM only modify the model architecture or resize the embedding space when adding new classes.

CE

Finally, because we evaluate according to a standard learning scenario (van de Ven & Tack, 2019; He et al., 2018), the model must be able to predict a label by choosing from all classes seen so far. At each time step, we choose a point from a batch \mathcal{B}_k with an associated label y_k , where \mathcal{Y}_k contains the class labels seen so far. The labels are $\mathcal{Y} = \bigcup_{k=1}^K \mathcal{Y}_k$ different class labels seen across all the batches. The MAP estimate is given by

$$\arg \min_y E_{\theta_K}(\mathbf{x}_k, y), \quad y \in \bigcup \mathcal{Y}_k,$$

$E_{\theta_K}(\mathbf{x}, y)$ is the energy function with parameters learned from training on the batches $\{\mathcal{B}_1, \dots, \mathcal{B}_K\}$. Our formulation can compute an energy for any distribution over labels, including unseen classes, which avoids needlessly defining the number of classes in advance.

IV. EBM TRAINING OBJECTIVE

We are committed to modeling the conditional distribution $p_D(x, y)$ as shown in Equation 3. Another key component of Boltzmann distribution is to define the joint energy $E_{\theta}(x, y)$. Then the training objective becomes

$$\mathbb{E}_{(\mathbf{x}, y) \sim p_D} [E_{\theta}(\mathbf{x}, y) - E_{\theta}(\mathbf{x}', y')],$$

where (\mathbf{x}, y) and (\mathbf{x}', y') are randomly sampled from the current distribution p_D . Since the main focus of this paper is to show that our formulation is both efficient and effective EBM training objective, we leave the detailed analysis of the results of using Equation 7 in training EBM to future work. In the supplement Section A, we show that the training objective is well-defined and satisfies the properties of a probability distribution.

Figure 3: Predicted label distribution for the split MNIST dataset. The SBC is the current task, while our EBM predicts the future task.

EBM training objective proposed energy-based training can be directly applied to existing methods by modifying the training of our proposed energy function. The softmax normalization operation is applied to each training batch. In Table 2, we can see that EBM significantly improves the performance of existing methods. This is because EBM does not suppress the probability of old classes while increasing the probability of new classes. EBM provides an orthogonal direction to the standard SBC model shown in Fig. 1. This is simple to implement on top of existing SBC models.

5.1.3. QUALITATIVE ANALYSIS

Most existing CL methods are based on SBC. To better understand why EBM can prevent catastrophic forgetting, we qualitatively compare it with the standard SBC model shown in Fig. 1.

Energy landscape. In Fig. 3, we show the energy landscapes after training on the first 100 labels of the MNIST dataset. For SBC, the energy landscape is flat at the negative value of the predicted probability. For EBM, the energy values (EBM) or energy differences (ΔEBM) between the 100 labels in the dataset are much larger than the SBC case. The values are normalized over the total number of labels. The diagonal elements in the energy matrix indicate correctly assigned labels. For the task T_9 , SBC assigns high energy values (around 80-90) to almost all the labels. For the task T_{10} , the highest probability label has an energy of around 100 (90-100). SBC tends to assign high energy values to both old and new classes for both old and new tasks. In contrast, EBM has low energy values for old classes and high energy values for new classes. This means that after training on the first 100 labels, EBM assigns low energies to the true labels of the first 100 labels. This shows that EBM is able to learn new tasks without catastrophic forgetting.

Predicted class distribution In Fig. 4, we show the predicted class distributions for the first 100 labels of the MNIST dataset, we plot the probability of the predicted classes. Only data from the first 100 labels are used for this figure. Taking the task T_9 as an example, it shows that EBM performs well on test data from the first 100 labels. The predicted class distributions for SBC and EBM are very similar. This shows that EBM is able to learn new tasks without catastrophic forgetting.

on existing approaches. Our proposed objective is simple and can also ingest CL approaches. We test this objective of baseline models to proxy objective, which computes the loss over class labels in the current task. We find that our proposed objective achieves the performance of different CL baselines. The new training objective does not rely on old classes when improving on new classes. Our training objective is designed to tackle the CL problem and outperform existing CL approaches.

ANALYSIS

EBMs are based on the softmax-based approach (Gershman et al., 2019; Zenke et al., 2017). To show that EBM suffer less from catastrophic forgetting, we directly compare our EBMs and the SBC model shown in Figure 1 in this part.

In Figure 2, we show the energy landscape of task 9 and task 10 of the permuted MNIST dataset. In task 9, the energy is given by the negative log-likelihood. Each datapoint has 100 predicted probabilities (SBC) corresponding to all 100 classes. For each datapoint, these values are sorted in descending order. Dark elements on the diagonal indicate correct predictions. After training on task 9, the predicted probabilities shift to classes from T_9 and away from classes from T_1 to T_8 . After learning task 10, the predicted probabilities shift to classes from task 10. This indicates that the model forgets to assign high probabilities to new classes when learning new data, indicating forgetting. In contrast, the energy landscape of our EBMs across tasks shows no such changes. Across the diagonal, the predicted probabilities remain high for all 100 classes. This indicates that the model does not forget old classes when learning new ones. In Figure 3, for the split MNIST dataset, the proportional distribution of predicted labels from the tasks seen so far was plotted. The second panel in the first row shows the distribution of predicted labels for two tasks after finishing training on them. The third panel in the first row shows the distribution of predicted labels for two tasks after learning the next task. The distribution of predicted labels for the first task remains stable, while the distribution for the second task changes significantly. This indicates that the model forgets old classes when learning new ones. In contrast, the distribution of predicted labels for the second task remains stable, while the distribution for the first task changes significantly. This indicates that the model does not forget old classes when learning new ones.

from each distribution, training one roughly uniformly (panel). It predicts class labels that fail to correctly classify tasks. In general, more uniform distributions in the continuing confusion matrix of model performance comparison of past and an evaluation of tasks in S.

5.1.4. IS THE TRAINING DATA UNIFORM?

Effect of non-uniformity on EBM training We compare three examples as described in section class. The first example is most similar to the optimized gradient descent batch training, randomly sampled negative examples. As shown in Table 3, we find the best results for the current state of our EBM, in terms of negative samples, same label interference, and seen class suppression.

Effect of label contamination on EBM training We modify the training

In class are similar, the ground truth probability should be uniform over those four classes. In the first task, the predictions of SBC are uniformly distributed over the first two classes. However, after learning new tasks, SBC forgets classes from the most recent task, indicating that it cannot correctly memorize classes from previous tasks. In contrast, the predictions of EBM are uniformly distributed over all classes seen so far.

To complement, we provide further analyses of the visual learning behavior of SBC and EBM. The confusion matrices in Supplement Section B.1 show the model capacity in Supplement Section B.2. The parameter importance in Supplement Section B.3 shows the evaluation of the interference with previous classes in Supplement Section B.4.

THE STRONG PERFORMANCE OF EBM IS DUE TO THE ENERGY TRAINING OBJECTIVE OVER THE LABEL CONDITIONING?

energy training objective. We conducted experiments on the CIFAR-100 dataset to investigate how different energy training objectives influence the CL results. We compare two different strategies for selecting the negative samples as described in Section 4.3.1. The first strategy selects samples so far as negative labels (**All Neg Samples**) which is similar to the way that the traditional EBM works [1]. The second one takes all the classes in the batch as negative labels (**All Neg Batch**). The first one selects one class from the current batch, while as we mentioned in Equation 7 (**1 Neg Sample**), we find using only one negative sample per class result, and using negatives sampled from the current batch is better than from all seen classes. The energy training objective aims at improving the energy samples while decreasing the energy of sampling negatives from the current batch. This balance with previous classes than sampling from the same classes. Using a single negative causes the model to focus on negative samples and thus has the best performance.

label conditioning. Next, we test whether label conditioning in our EBMs is important for learning. The results are shown in Table 2. As mentioned in Section 4.3.2, we compare baseline models using our training objective with the proposed EBM. The proposed EBM outperforms the baseline models in terms of accuracy and energy loss.

architecture also contributes to why EBM^{*} suffers from trophic forgetting.

Testing accuracy of each task as the training progresses is shown in Figure 4. We compare the standard classification error using our training objective (SBC^{*}) with the standard classification error using the permuted MNIST. We find that the accuracy of SBC drops sharply when learning new tasks, while the SBC drop is much less when learning new tasks using the new training objective used in SBC^{*}. The accuracy of EBM^{*} drops even slower than SBC^{*}, indicating that EBM^{*} can mitigate the forgetting problem.

It is interesting to note how the strong performance of EBM^{*} compares with the standard classification error of both the EBM training objective and the standard training architecture. Moreover, these results are obtained using a simple, and counterintuitively, directly opposite loss function to the cross-entropy loss used by existing approaches. This suggests a new way to approach continual learning.

TABLE 1: PERFORMANCE OF DIFFERENT ARCHITECTURES

We first evaluate the ability of EBM^{*} in integrating data information and learning features from the energy function. To investigate this, we propose to combine the information from data x and label y to construct a series of experiments on CIFAR-10 using four model architectures (V1-V4) that are trained sequentially at three stages: early, middle, and late stages, respectively (see Section C for more details). We find that V4 (trained in the late stage (V4)) performs the best in terms of learning a feature embedding of labels. This is due to the projection matrix which is randomly initialized from a uniform distribution $\mathcal{U}(0, 1)$. Even using this random projection can already generates better results than the baselines in Table 1. Note further that the “Fix” setting has much worse performance than the other baselines. Using a learned feature embedding can further improve the result. We may consider different normalization methods over the feature embedding. We find that Softmax (**End Softmax (V4)**) is better than L2 normalization (**End Norm2 (V4)**) and no normalization (**V4**).

TABLE 2: PERFORMANCE FOR DIFFERENT NUMBERS OF TASKS

Table 2 shows the accuracy of our proposed EBM^{*}, in comparison with the baseline EBM, for the task sequence of 5 tasks. We perform two types of experiments on CIFAR-100. In the first type of experiments, we train the network sequentially on 5 tasks. In the second type of experiments, we train the network sequentially on 5 tasks, but we also perform a few “refresh” experiments on the first 4 tasks. Refresh experiments are performed by training the network on the first 4 tasks, and then performing a few epochs of training on the 5th task. We find that the proposed EBM^{*} achieves better performance than the standard EBM across all the experiments. The proposed EBM^{*} also achieves better performance than the standard EBM in the refresh experiments. This indicates that the proposed EBM^{*} can mitigate the forgetting problem.

10 classes up into three classes EBM

5.1.7.

Although the *Classical* approach are following a proactive approach, they show that training from scratch from learners' buffer zone so far has been over a period of time in the range of 10 days.

We recall that the might be a number of existing chitcotechniques split them into lines of similar sizes and sizes of EBM and EBM, and C EBM, after that.

Again, using the proposed tactic combination, we can combine the two approaches.

5.2. EBM

When we are not ever, we are present, when we are not

ases per task. Here we additionally split the tasks into 5 tasks (i.e., 20 classes per task), 10 tasks (i.e., 10 classes per task) and 50 tasks (i.e., 2 classes per task). Our method substantially outperforms the baseline.

CAN EBMS USE REPLAY?

Although EBMs already achieve good performance on *IL*-*IL* without using any replay, we found that EBMs are flexible enough to be combined with replay techniques to further improve the performance. We compare the results of SBC and EBM using exact replay. When learning a new task, we mix the new data with old data stored in a memory buffer that stores examples from previous learned tasks to train the models. The examples in the buffer are randomly selected from the classes in the buffer, and the available memory budget k is the number of all classes encountered so far. We provide more details in Supplement Section D.

We report the results on 4 datasets. Note that the results may be slightly different from results reported in previous works because of the usage of different architectures, different memory sizes, and different training protocols. In the datasets. In our experiments, we compare SBC and EBMs to have similar model architecture, the same number of model parameters and the same training time on each dataset. After using extra memory, both SBC and EBM have improvements. EBMs still outperform SBC especially on more challenging datasets, such as CIFAR-10 and CIFAR-100. Interestingly, on CIFAR-100, EBMs outperform SBC without using replay (30.28% vs. Table 1). This shows that EBM can learn better than SBC with using replay (28.57% vs. Table 1).

In our focus is to try to address the CL problem by using replay or stored data. Our results show that the proposed EBM formulation provides an orthogonal approach to tackle the CL problems. This approach can be combined with existing CL methods.

now that EBMs

5.2.1. DATASETS
For the *boundary-aware* experiments, we use the same datasets as in Section 5.1.1. We consider the “continual learning” proposed in [1] where the classes are continually changing in each subsequent epoch. All experiments are run under this scenario.

5.2.2. COMPARISON
We focus on the impact of different model architectures, memory footprint, and thus computational costs. Since there are many methods of continual learning, it is not surprising that many of them are generally able to take the core action of learning a new task. However, a comparison of the performance of most algorithms, especially in terms of computational complexity, is missing. Online EWC, SI, and our proposed methods use a similar strategy, namely *boundary-aware sampling*, but sometimes with different weights. The results are shown in Figure 5. The error bars represent the standard deviation of the mean. The error bars are omitted for brevity.

The results are shown in Figure 5. It can be seen that all three methods have a significant advantage over the baseline. The experiments show that the proposed methods are able to handle data from different domains more effectively than the baseline. The proposed methods also achieve better performance than the baseline in terms of accuracy and computational complexity.

6. Discussion
In this paper, we propose a new framework for promising class-incremental learning. We evaluate our framework on a continual learning setting and show that it achieves many desirable characteristics. In particular, we show that getting in contact with a new class does not significantly affect the performance of the model. This is because the proposed framework is able to learn new classes incrementally without forgetting previously learned classes. Additionally, we show that the proposed framework is able to learn new classes in a more efficient way than the baseline. This is because the proposed framework uses a more efficient learning algorithm than the baseline. Finally, we show that the proposed framework is able to learn new classes in a more robust way than the baseline. This is because the proposed framework is able to handle data from different domains more effectively than the baseline.

lary-aware and be found that EBMs

ON PROTOCOLS
 described in Section 3.2, “continuous task-agnostic” setting (Li et al., 2018) to generate training data. In this setting, the frequency of task changes increases and decreases linearly according to the *Clock* protocol.

TUNING METHODS
 The performance of our model size and model capacity with replay-based training depends on the number of trials required for learning that rely on task switching. One trivial adaptation is to switch task at each step instead of waiting for a full epoch. This adaptation is impractical because of the large number of tasks. Instead, we managed to reduce the number of task switches in this way. Although the proposed model architectures as well as the experimental setup were performed on multiple benchmarks, the results are representative.

RESULTS
 We observe that EBM outperforms baselines on all the datasets. This indicates that the generalization of EBMs to new problems as EBMs can learn without task boundary information.

Work
 Energy-based models have been applied to a variety of different domains. We demonstrate that EBMs can learn to prevent catastrophic forgetting while we experimentally show that they can enhance the challenging task of learning to do multiple benchmarks simultaneously. Boundary-agnostic solutions have larger model sizes than their task-specific counterparts.

capac-