

---

# TASKDROP: A COMPETITIVE BASELINE FOR CONTINUAL LEARNING OF SENTIMENT CLASSIFICATION

---

A PREPRINT

**Jianping Mei\***

Zhejiang University of Technology  
jpmei@zjut.edu.cn

**Yilun Zheng\***

Zhejiang University of Technology  
zhengyilun@foxmail.com

**Qianwei Zhou**

Zhejiang University of Technology  
zhouqianweischolar@gmail.com

**Rui Yan**

Zhejiang University of Technology  
Ryan@zjut.edu.cn

December 7, 2021

## ABSTRACT

In this paper, we study the multi-task sentiment classification problem in the continual learning setting, i.e., a model is sequentially trained to classifier the sentiment of reviews of products in a particular category. The use of common sentiment words in reviews of different product categories leads to large cross-task similarity, which differentiates it from continual learning in other domains. This knowledge sharing nature renders forgetting reduction focused approaches less effective for the problem under consideration. Unlike existing approaches, where task-specific masks are learned with specifically presumed training objectives, we propose an approach called Task-aware Dropout (TaskDrop) to generate masks in a random way. While the standard dropout generates and applies random masks for each training instance per epoch for effective regularization, TaskDrop applies random masking for task-wise capacity allocation and reuse. We conducted experimental studies on three multi-task review datasets and made comparison to various baselines and state-of-the-art approaches. Our empirical results show that regardless of simplicity, TaskDrop overall achieved competitive performances for all the three datasets, especially after relative long term learning. This demonstrates that the proposed random capacity allocation mechanism works well for continual sentiment classification.

## 1 Introduction

The capability of learning new tasks while maintaining performance on learned ones is required for many practical applications where tasks are learned sequentially. For example, retraining a robot or a sentiment classification system each time encountering a new task is cumbersome or even impossible when previous data are no longer accessible. The sequential learning or continual learning ability is also fundamental for advanced artificial intelligence systems to adapt to unknown tasks Legg and Hutter [2007], Thrun and Mitchell [1995]. When sequentially learning tasks with very limited or no access to data of previous tasks, the model tends to forget what has been learned, leading to degraded performance on previous tasks McCloskey and Cohen [1989], Ratcliff [1990].

Most of the recent efforts in continual learning are made to deal with this so-called Catastrophic Forgetting problem for learning with deep neural network models. Some of them follow the joint training idea by replying some forms of previous information that are stored Rebuffi et al. [2017], Lopez-Paz and Ranzato [2017], Chaudhry et al. [2019], Rolnick et al. [2019], Han et al. [2020] or synthesized with a generative model Shin et al. [2017], Nguyen et al. [2018] when learning new tasks. Others focus on preserving the learned model by strictly freezing or penalizing large changes to the subset of model parameters which are regarded to be important to previous tasks French [1991], He and Jaeger

---

\*The authors contribute equally to this paper.

Method	Fixed model size	Way for mask obtaining	Mask sparsity control	Mask overlap
PackNet Mallya and Lazebnik [2018]	no	pruning with assigned ratio	pruning ratio	no
HAT Serrà et al. [2018]	yes	minimizing the weighted $l_1$ norm of attentions	regularization weight $c$	yes
KAN Ke et al. [2020a]	yes	minimizing the training loss of the current task	no direct control	yes
TaskDrop	yes	random sampling	retention ratio $p$	yes

Table 1: Comparison of characteristics of masking-based continual learning approaches.

[2018], Nicolas Y. Masse and Freedman, Mallya and Lazebnik [2018], Rusu et al. [2016], Fernando et al. [2017], Kirkpatrick et al. [2017], Li et al. [2019], Rajasegaran et al. [2019].

Although freezing-based approaches work in a more direct way for model preservation and hence forgetting reduction than regularization-based ones, the constant demand on additional capacity for new tasks becomes an issue when the number of tasks is large. In this work, we focus on the continual learning of sentiment classification tasks with a size-fixed encoder without re-accessing to data of previous tasks.

Given the overall model capacity, preserving existing knowledge inevitably restricts the learning of new tasks. Facing the tradeoff between remembering what have been learned and exploring new knowledge with a risk of forgetting, it is reasonable to favour the former if these tasks have low transferability because of their difference in distribution and other aspects, such as classification of different image datasets Li and Hoiem [2017], Serrà et al. [2018] or learning different types of relations Han et al. [2020]. For the sentiment classification scenario, typical sentiments and words used to express these sentiments are quite similar across tasks, e.g., “amazing” is a commonly used positive word when one makes comments on books or anything else. This high cross-task similarity nature increases the chance for backward transfer, i.e., learning from new tasks helps to strengthen the common knowledge base shared by previous tasks. For such kind of problems, preserving-oriented approaches usually become less effective as observed in Ke et al. [2020a] as well as in our experiments.

Instead of completely avoid reactivation of those preserved parameters, masking-based approaches learn non-exclusive masks so that frozen units still have the opportunity to be re-activated Fernando et al. [2017], Ke et al. [2020b,a]. A key problem for masking-based approaches is the way to obtain masks. Existing approaches learn masks based on different objectives, which are designed for specific type of problems. In this paper, we investigate a simple randomization-based approach called Task-aware Dropout (TaskDrop). At the beginning of a task, we generate random binary vectors as masks for each of the layers, and apply the corresponding mask to units of each layer during forward and backward passes. This random unit masking operation is also adopted in the well known dropout, an important trick for deep learning. While the standard dropout applies random masks to sample a large number of subnetworks for each forward pass without considering task boundaries, masking in TaskDrop works as random capacity allocation for each coming task. Capacity reuse or re-activation across tasks is simply controlled by the dropout rate, rather than guided by any presumed objectives as existing approaches. This gives the flexibility for TaskDrop to adjust its conservativeness to a proper degree, leading to better adaptiveness to different problems.

We carried out sentiment classification experiments in the continual learning setting on three multi-task datasets with different levels of transfer accuracy. Each task performs sentiment polarity prediction of reviews of products from one category. Along with state-of-the-art approaches, we also include the results of four baselines to account for various naive solutions as well as an upper-bound solution. Main contributions of this work are summarized as below:

- We investigate a simple masking-based approach for sequential learning of sentiment classification tasks. Masking with randomly generated task specific masks results in random capacity reuse, which is controlled by the dropout rate.
- Analytical discussions and empirical results show that the random capacity reuse mechanism works competitively compared to other elaborately designed approaches.
- We carried out experimental study on the robustness of continual learning approaches with tasks of different levels of relevance, to investigate how they perform when the nature of the target problem does not well fit their assumptions.

## 2 Continual Learning with Random Task Masks

### 2.1 Masking-based Continual Classification

Masking-based approaches apply masks to parameters or unit activations to decide what to access and to update when learning a new task. The most important difference among these approaches is thus how these masks are generated or learned. To preserve existing knowledge, PackNet Mallya and Lazebnik [2018] generates exclusive masks with an

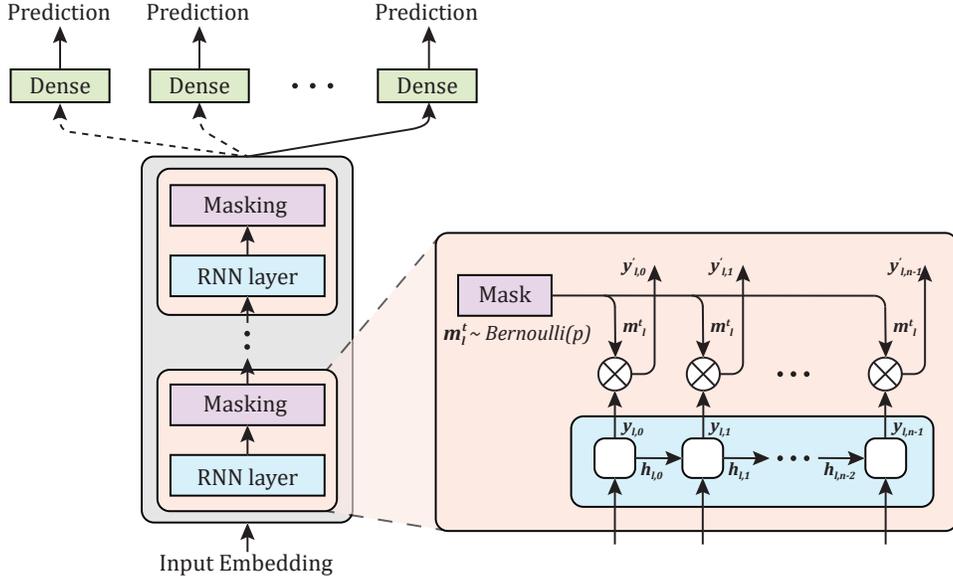


Figure 1: An overall model structure of TaskDrop consisting of a RNN encoder and multi-head dense layers. Each RNN layer is followed by the task-aware random masking operation with its details on the right.

assigned pruning ratio to directly freeze large-valued parameters, while HAT Serrà et al. [2018] only decreases the chance for reactivating units that have been already activated in previous tasks. In the most related approach KAN Ke et al. [2020a], which also works on continual sentiment classification, masks are learned to activate units to give optimized learning of the current task.

Inspired by the success of the dropout technique for deep learning, we explore a new solution based on random task-aware masks for the problem of continual learning. Instead of treating masks as additional parameters to be learned simultaneously with the network, our approach called Taskdrop directly generates masks like PackNet, but allows capacity reuse like those learnable masks. Table 1 compares characteristics of the proposed masking-based approach with those of existing ones. Next, we present the details on instantiating this idea into a solution for continual sentiment classification.

## 2.2 Task-Aware Dropout

**Overview.** The overall structure of TaskDrop is illustrated in Figure 1. It consists of a RNN backbone with a multi-MLP output layer, where each MLP corresponds to the classifier of a certain task. It is pretty much the same as the main network in KAN Ke et al. [2020a]. While in KAN an individual accessibility module, i.e., another RNN is trained in an alternating way with the main network for learning the masks, here in TaskDrop we only have the main network to learn.

Assume that the input embedding for the  $j$ th training case of task  $t$ , i.e., a piece of text with fixed length  $n$  is  $\mathbf{X}_j^t = [\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n]$ , and the output of the  $l$ th hidden layer at timestep  $i$  is  $y_{l,i}^t$ . For each hidden layer  $l$  with  $n_l$  units, a binary mask  $\mathbf{m}_l^t \in \{0, 1\}^{n_l}$  is element-wise multiplied to each of the  $n$  timestep outputs of this layer. At the beginning of a task, we generate masks for all the layers  $\{\mathbf{m}_l^t\}_l$  by randomly and independently drawing each element from the Bernoulli distribution with a probability of  $p$  for being 1, i.e.,

$$\mathbf{m}_l^t = [\text{Bernoulli}(p), \dots, \text{Bernoulli}(p)]. \quad (1)$$

These masks are stored and retrieved in every forward and backward pass during the whole learning process of this task. The probability  $p$  decides how likely that a hidden layer unit is retained or activated for this task. Simply speaking,  $p$  controls the sparsity of activated units of each layer, which is formulated with the  $l_1$  norm of parameter Yoon et al. [2018] or attention Serrà et al. [2018] of each layer. Increasing the value of  $p$  is likely leading to a model with more extensive unit sharing, or reuse.

After getting the class prediction, a cross-entropy loss is calculated for updating parameters of RNN and the  $t$ th classifier. Next, we elaborate more on forward passes and back-propagations with masking applied RNN.

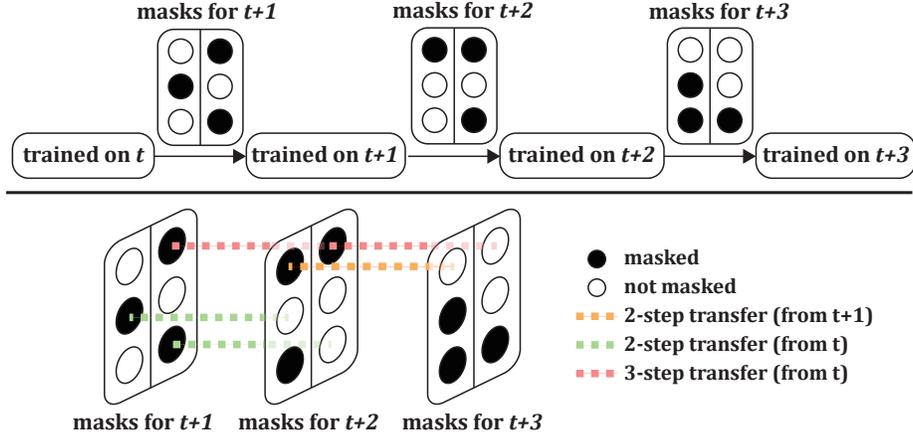


Figure 2: Illustration of skip-task transfer. Top: Masking-based learning of tasks in sequence. Bottom: different skip-task transfers produced by random masks of three subsequent tasks.

**Masking of RNN.** Following the work on modified dropout to RNN structures Zaremba et al. [2014], the masking operation as illustrated in Figure 1 only works on non-recurrent connections denoted by dashed arrows. To be more clear, given  $\mathbf{m}_i^t$  and the RNN output  $\mathbf{y}_{l,i} = \mathbf{h}_{l,i}$  at timestep  $i$  for a training case of task  $t$ , where  $\mathbf{h}_{l,i}$  is the recurrent state at this timestep, we apply the masking operation to  $\mathbf{y}_{l,i}$  only as below during feed-forward

$$\mathbf{y}'_{l,i} = \mathbf{y}_{l,i} \odot \mathbf{m}_i^t. \quad (2)$$

For recurrent connections as shown by solid arrows, information still flow through the units even they are masked out, so that the valuable memorization ability of GRU is not sacrificed. During back propagation, masked units receive no gradient.

Another detail is that we applied unit masking like the standard dropout, as well as HAT and KAN, while PackNet operates masking over weights.

### 2.3 Why Random Masks Work?

We now provide some analytical discussions on unique properties of TaskDrop to hopefully shed some light on understanding and explanation of the underlying working principle of this simple approach.

**Skip-task transfer.** The random masking operation gives TaskDrop the ability of  $s$ -step skip-task transfer with  $s > 1$  as illustrated in Figure 2. Specifically, parameters of the current task  $t$  may be preserved and used by the  $s$ th subsequent task with  $s \geq 2$  if the corresponding units happen to be masked during all the next  $s - 1$  tasks, i.e., a preservation duration of  $s - 1$  tasks. This mechanism allows direct knowledge transfer from one task to another that is  $s$  steps way on the task stream.

Given  $p$  the probability of being not masked, the probability of  $s$ -step sharing for each individual unit is  $P_p(s) = (1 - p)^{s-1}p$ . Figure 3 plots the curves of  $P_p(s)$  with respect to  $s$  for  $p = \{0.2, 0.4, 0.6, 0.8\}$  on the left and with respect to  $p$  for  $s = \{2, 3, 4\}$  on the right. As shown from the left plot of this figure,  $P_p(s)$  is monotonically decreasing with respect to  $s$ , and the larger the  $p$  is, the faster it decreases from with the initial value that equals to  $p$ . Let us focus on the conditions that lead to large  $P_p(s)$  with  $s \geq 2$  on the right. It is clear that  $s = 2$  has the largest probability for all  $p \in (0, 1)$ , and when  $p$  is either too small or too large, i.e.,  $0.3 \leq p \leq 0.7$ , we have  $P_p(s = 2) \geq 0.2$ . If increasing the steps to  $s = 3$ , the largest  $P_p(s = 3)$  is 0.15 when  $p$  is around 0.3. Although probabilities of skip-task transfers for  $s \geq 2$  are relative small compared to direct transfer, that is, re-activate in the next task, they actually help TaskDrop to achieve better results as demonstrated with empirical results.

Based on the above analysis that connects random masks with skip-task knowledge sharing, we next give some discussions by theoretical comparison between TaskDrop and other closely related approaches to further discuss the merits of random masks for continual learning.

**Parameter sharing: whole vs. portional.** When  $p \rightarrow 1$ , TaskDrop reduces to the No-masking baseline, which shares all the encoder parameters in learning the next task. In other words, No-masking is a hundred percent reuse variant of TaskDrop with no skip-task sharing. It may give good results if the next task is highly similar to the current one,

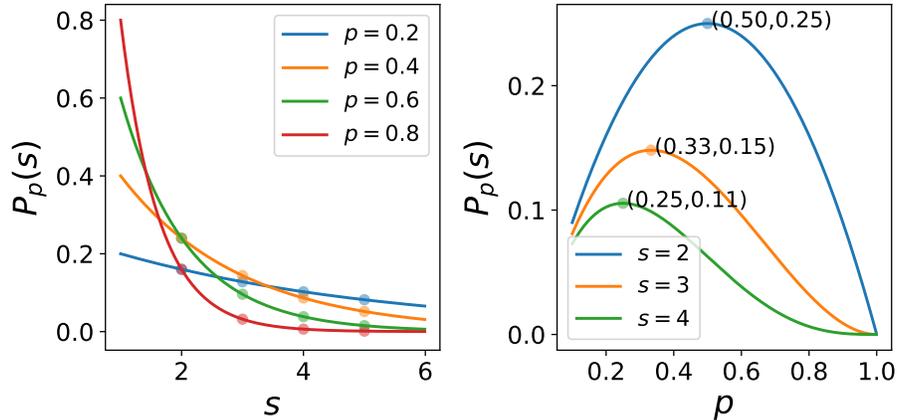


Figure 3: The probability of  $s$ -step skip-task transfer for each individual unit  $P_p(s)$  with respect to  $s$  on the left and with respect to retention ratio  $p$  on the right.

which means a smooth task boundary. Even though, masking a relative small portion of weights to allow some skip-task sharing makes the model more robust and usually performs better in relative long term learning as observed in our experiments. In cases that tasks are of little similarities, a large portion of parameter sharing could cause severe forgetting. For TaskDrop, we can adjust it to give more conservative solutions with a smaller  $p$ , leading to reduced activation overlap among representations French [1991].

**Preserving: temporary vs. permanent.** Freezing-based approaches permanently preserve parameters by prohibiting re-activation in all subsequent task learning. The  $s$ -step skip-task transfer in TaskDrop works as temporary preserving for a period of  $s - 1$  task learning. Although no-reactivation effectively solves the catastrophic forgetting problem, fresh capacity is constantly required in order to learn new tasks. That is to say, for a fixed model size, such kind of one-time using exhausts all the capacity after a certain number of tasks. Moreover, being over conservative results in unsatisfied results for the problem under consideration, where knowledge sharing is important.

**Dropout: per-sample each iteration vs. per-task.** Regardless the similarity in the way of mask generation and application, masks used in dropout and TaskDrop are different in both their lasting time and the number of instances to be applied for. The standard dropout generates random masks for each training sample each time it is processed, and the sampled network is learned upon a single instance. In TaskDrop, masks are generated at a per-task frequency and shared by all the samples during the whole learning process of this task. The above differences indicate that while dropout is an effective trick for improving model generalization ability, it is not designed for continual learning as task boundaries are unaware during the entire learning process.

## 2.4 Limitations

Since the skip-task transfer mechanism works properly after learning an enough number of tasks, TaskDrop may not be a good choice for learning only a very short task sequence. For example, when sequentially learning two tasks, random masks for the second task have little change of being better than those learned ones. Considering that TaskDrop incurs no additional parameters to be learned, we should have resources left for adopting a higher capacity classifier as a remedy to the effectiveness gap between random and learned masks if there is any.

## 3 Related Work

Given the long and diverse literature of continual learning, we only focus on recent approaches with a similar problem setting, i.e., continual task learning with neural network models without replaying data of previous tasks. Since dropout itself is not the focus of this paper, works on different dropout techniques are also excluded here.

**Regularization-based methods.** Structural regularization approaches penalize changes to learned knowledge when learning a new task Kirkpatrick et al. [2017], Lee et al. [2017], Li and Hoiem [2017], Zenke et al. [2017], Rahaf Aljundi and Tuytelaars [2018], He and Jaeger [2018], Wang et al. [2019], Yoon et al. [2018]. These approaches use different representations of the existing knowledge that they attempt to preserve, including output predictions Li and Hoiem

Metric	Dataset	Individual Networks	Classify-only	No-masking	TaskDrop	Multi-task
$A^{\leq 2}$	high-6	82.73 $\pm$ 1.41	84.03 $\pm$ 1.41	<b>86.03 <math>\pm</math> 1.54</b>	84.97 $\pm$ 1.38	86.75 $\pm$ 1.36
	mix-24	79.12 $\pm$ 2.90	80.71 $\pm$ 6.05	<b>82.90 <math>\pm</math> 3.20</b>	81.23 $\pm$ 3.16	84.13 $\pm$ 2.78
	low-6	74.33 $\pm$ 3.94	73.29 $\pm$ 4.53	78.84 $\pm$ 2.59	<b>80.04 <math>\pm</math> 3.64</b>	82.12 $\pm$ 2.77
	Average of three	78.73	79.34	<b>82.59</b>	82.08	84.33
$\rho^{\leq 2}$	high-6	-19.00 $\pm$ 4.15	-15.33 $\pm$ 4.56	<b>-10.64 <math>\pm</math> 6.24</b>	-13.30 $\pm$ 3.78	00.00
	mix-24	-10.57 $\pm$ 8.03	-12.02 $\pm$ 16.56	<b>-3.70 <math>\pm</math> 3.92</b>	-8.54 $\pm$ 3.49	00.00
	low-6	-50.44 $\pm$ 30.99	-52.14 $\pm$ 27.86	-33.54 $\pm$ 17.99	<b>-26.79 <math>\pm</math> 10.74</b>	00.00
	Average of three	-26.67	-26.50	<b>-15.96</b>	-16.21	00.00
$A^{\leq T}$	high-6	81.88 $\pm$ 0.55	83.58 $\pm$ 0.58	87.47 $\pm$ 0.89	<b>87.86 <math>\pm</math> 0.85</b>	90.83
	mix-24	78.33 $\pm$ 1.13	78.26 $\pm$ 3.72	87.47 $\pm$ 1.03	<b>87.87 <math>\pm</math> 0.82</b>	90.44
	low-6	72.37 $\pm$ 2.74	71.05 $\pm$ 4.54	79.06 $\pm$ 2.07	<b>80.83 <math>\pm</math> 0.88</b>	87.16
	Average of three	77.53	77.63	84.67	<b>85.52</b>	89.48
$\rho^{\leq T}$	high-6	-21.72 $\pm$ 1.28	-17.49 $\pm$ 1.42	-8.02 $\pm$ 2.19	<b>-7.06 <math>\pm</math> 2.05</b>	00.00
	mix-24	-32.44 $\pm$ 3.24	-33.06 $\pm$ 10.17	-8.01 $\pm$ 3.25	<b>-6.76 <math>\pm</math> 2.19</b>	00.00
	low-6	-52.74 $\pm$ 9.24	-55.81 $\pm$ 15.61	-30.70 $\pm$ 9.59	<b>-23.56 <math>\pm</math> 3.44</b>	00.00
	Average of three	-35.63	-35.45	-15.57	<b>-12.46</b>	00.00

Table 2: Comparison with reference approaches in averaged accuracy  $A^{\leq t}$  and forgetting ratio  $\rho^{\leq t}$  for two and all tasks.

[2017], hidden spaces Triki et al. [2017], or model parameters Kirkpatrick et al. [2017], Yoon et al. [2018]. The parameter importance may be measured with the diagonal of the Fisher information matrix Kirkpatrick et al. [2017] or based on the sensitivity of the learned function to their changes after convergence Rahaf Aljundi and Tuytelaars [2018], or computed during training in an online manner Zenke et al. [2017]. The work in Lee et al. [2017] further extends Kirkpatrick et al. [2017] with a separate model-merging step after learning a new task. A recent approach uses both regularization and memory-based example replay Huang et al. [2021].

**Freezing-based approaches.** Although regularization-based approaches work with well-defined objectives under certain assumptions, they only take implicit and indirect control over the forgetting problem. A straightforward way to overcome catastrophic forgetting is to immediately freeze parameters learned for a task and seek for additional capacity for new task learning Rusu et al. [2016], Fernando et al. [2017], Mallya and Lazebnik [2018], Yoon et al. [2018]. Some approaches pre-allocate certain capacity Rusu et al. [2016] or dynamically add capacity for a coming task Yoon et al. [2018], leading to models with a growing size. On the contrary, Mallya and Lazebnik [2018] starts with a size-fixed model, and frees part of capacity through parameter pruning. However, capacity distribution over tasks based on scheduled pruning ratios is not applicable when the number of tasks is unknown in advance. Moreover, allowing no reuse at all makes the model stop learning when the capacity limitation is reached.

**Masking-based approaches.** Masking-based approaches use masks for selective capacity reuse. Masks are learned to optimize specific objectives. In Serrà et al. [2018], masks with small sparsity and attention overlap are learned simultaneously with the network training by minimizing the weighted  $l_1$  norm of hard attentions. In the approach that is specifically designed for sentiment classification Ke et al. [2020a], masks are learned with an additional so-called accessibility module in an alternating way with the main network to optimize the performance of the current task. The work Rajasegaran et al. [2019] is an improvement over Fernando et al. [2017] for encouraging knowledge sharing and reuse.

In order to be conveniently trained with the rest of network using SGD, masks are relaxed to continuous attentions with a sigmoid function, which are gradually harden and become binary after an annealing process in Serrà et al. [2018]. Several later works follow this annealing strategy as well as masked gradient for training masks Rajasegaran et al. [2020], Ke et al. [2020a,b]. In a following work of Ke et al. [2020a], separate steps to measure whether two tasks are similar or not by comparing to reference models Ke et al. [2020b].

## 4 Experiments

We provide empirical results on three different subsets of review data to evaluate the performance of TaskDrop for continual sentiment classification. We compare TaskDrop with different reference and state-of-the-art approaches for their performance after short and long term sequential learning and visualization of learned representations. Comparison is also made between TaskDrop and the standard dropout. Finally, we study the impact of hyper-parameter. More details on experimental settings and additional results are included in the technical appendix due to space limitation.

Metric	Dataset	LwF	EWC	HAT	KAN	TaskDrop
$A^{\leq 2}$	high-6	<b>85.38 ± 1.35</b>	83.78 ± 1.83	84.68 ± 1.69	85.26 ± 1.79	84.97 ± 1.38
	mix-24	80.74 ± 4.94	<b>82.13 ± 2.65</b>	80.75 ± 3.04	81.39 ± 2.64	81.23 ± 3.16
	low-6	76.25 ± 3.36	79.74 ± 2.82	<b>80.38 ± 3.19</b>	80.34 ± 3.05	80.04 ± 3.64
	Average of three	80.79	81.88	81.94	<b>82.33</b>	82.08
$\rho^{\leq 2}$	high-6	<b>-11.71 ± 4.70</b>	-15.88 ± 3.48	-13.57 ± 4.01	-12.09 ± 4.67	-13.30 ± 3.78
	mix-24	-7.59 ± 9.10	-3.79 ± 13.58	-6.21 ± 4.79	<b>-3.76 ± 4.53</b>	-8.54 ± 3.49
	low-6	-44.07 ± 25.18	-28.76 ± 13.44	<b>-25.92 ± 11.30</b>	-27.30 ± 15.41	-26.79 ± 10.74
	Average of three	-21.12	-16.14	-15.23	<b>-14.38</b>	-16.21
$A^{\leq T}$	high-6	86.95 ± 0.65	85.40 ± 0.42	85.86 ± 0.45	87.19 ± 0.70	<b>87.86 ± 0.85</b>
	mix-24	83.38 ± 3.18	86.59 ± 0.65	82.99 ± 0.36	84.39 ± 1.04	<b>87.87 ± 0.82</b>
	low-6	76.70 ± 2.29	77.75 ± 3.69	78.97 ± 1.39	79.27 ± 1.54	<b>80.83 ± 0.88</b>
	Average of three	82.34	83.25	82.61	83.62	<b>85.52</b>
$\rho^{\leq T}$	high-6	-9.31 ± 1.59	-13.01 ± 1.09	-12.04 ± 1.16	-8.75 ± 1.67	<b>-7.06 ± 2.05</b>
	mix-24	-18.51 ± 8.36	-9.91 ± 1.73	-19.65 ± 1.15	-16.32 ± 3.11	<b>-6.76 ± 2.19</b>
	low-6	-38.30 ± 7.65	-35.11 ± 14.59	-28.85 ± 6.13	-30.57 ± 6.85	<b>-23.56 ± 3.44</b>
	Average of three	-22.04	-19.34	-20.18	-18.55	<b>-12.46</b>

Table 3: Comparison with state-of-the-art approaches in averaged accuracy  $A^{\leq t}$  and forgetting ratio  $\rho^{\leq t}$  for two and all tasks

#### 4.1 Setups

**Data.** We use the continual sentiment classification data from Ke et al. [2020a]. This dataset consists of Amazon reviews of 24 different categories of products, and each category makes up a task. In order to carry out more comprehensive evaluations on how each approach performs on tasks with different representational relevance, we further extract another two subsets with six categories each from this twenty-four dataset. Specifically, we define the Mutual Transfer Accuracy (MTA) between two tasks as the average of testing accuracy on one task of the model trained on the other, and then select six tasks that have the largest or smallest total MTA to form high-6 and low-6, respectively. The original dataset that has a mixed MTA level is referred as to mix-24.

**Baselines** We consider four state-of-the-art continual learning approaches that work with a fixed model with capacity reusing. LwF Li and Hoiem [2017] and EWC Kirkpatrick et al. [2017] are popular regularization-based approaches for handling catastrophic forgetting, and HAT Serrà et al. [2018] and KAN Ke et al. [2020a] are two masking-based approaches by learning task-specific masks with particular objectives. We also consider the three different reference approaches in continual learning. The results of multi-task learning where all tasks are jointly learned are also reported as an upper bound for sequential learning.

- **Individual Networks.** It trains an individual network from scratch for each task. This baseline gives the no-reusing solution for minimum forgetting at the cost of large capacity for each task.
- **Classify-only.** All tasks use the same encoder learned with the data of the first task and only train a task-specific classifier with the data of each task.
- **No-masking.** It has the same model structure as Classify-only, but all parameters of the encoder are sequentially learned on all tasks. It is a full knowledge sharing approach without any direct control over forgetting.

**Network and training.** We use the same network architectures and hyperparameter settings as Ke et al. [2020a], i.e., a single layer GRU for RNN encoder, a fully-connected layer for each classifier and the input embedding with pretrained BERT-base Devlin et al. [2018]. We set the retention ratio  $p$  in TaskDrop to 0.8, 0.5, and 0.6, respectively for high-6, mix-24, and low-6 when comparing with other continual learning approaches.

Each approach is evaluated in terms of averaged accuracy  $A^{\leq t}$  and forgetting ratio  $\rho^{\leq t}$  on the testing data of each task after a short-term learning of two tasks and long-term learning of  $T$  tasks.  $T$  is the total number of tasks of each dataset. Assume  $a^{\tau \leq t}$  is the testing accuracy for each task  $\tau$  after training  $t$  tasks, the averaged accuracy takes average over all the  $t$  tasks, i.e.,  $A^{\leq t} = \frac{1}{t} \sum_{\tau=1}^t a^{\tau \leq t}$ . Forgetting ratio Serrà et al. [2018] is an adjusted variant of accuracy by comparing to the random and joint learning accuracy.

Since the results are order-dependent, we generate 10 different sequences randomly for each dataset, and report the mean and standard deviations of these 10 sequences for all comparisons. All results are reproducible with settings used here. Codes will be published upon acceptance of the paper.

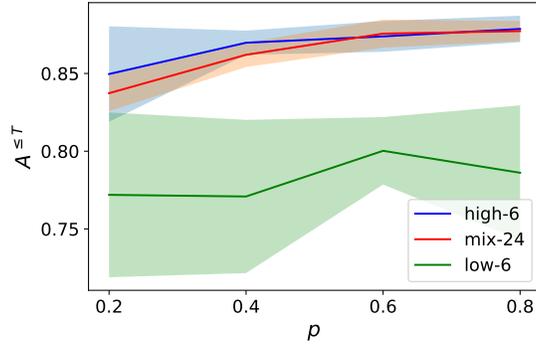


Figure 4: Performance in  $A^{\leq T}$  of TaskDrop on three datasets with respect to different retention ratio  $p$ .

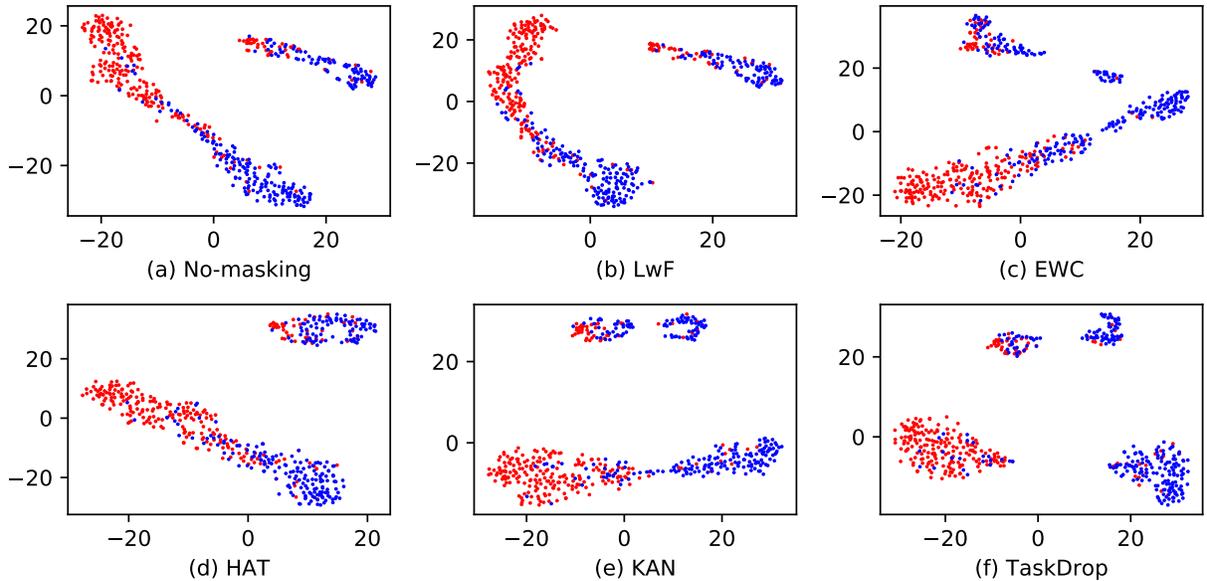


Figure 5: Representations of testing instances in *Instant Video* generated with different approaches, which are trained on other tasks in low-6. Representations are visualized by t-SNE Van der Maaten and Hinton [2008].

## 4.2 Comparison with Reference Approaches

Table 2 compares TaskDrop and four reference approaches on the three datasets. The average of the means over three datasets are also given to show the overall performance over three datasets. The results with respect to two metrics are pretty much consistent, but forgetting ratio is more sensitive than accuracy, i.e., two results with very close accuracy may have a large difference in forgetting ratio.

It is no surprise that Multi-task gives the best results for all the cases as it jointly learns data of all the tasks. For the rest approaches that have no access to data of previous tasks, TaskDrop and No-masking, the two which encourage more capacity reuse perform better. While the Individual Networks was reported to give the best results on image classification in Mallya and Lazebnik [2018], it performs much worse than all others for these sentiment datasets. By comparing No-masking’s long-term results with its short-term ones, it is seen that results of the three datasets are improved when learning more tasks. All the above observations confirm that forgetting is the main issue only for some continual learning problems, but not for sentiment classification.

Comparing the top approaches, TaskDrop outperforms No-masking from a long term, and No-masking gives good results on the two datasets with large cross-task relevance. As discussed early, random masking equips TaskDrop with the ability of skip-task transfer, which helps to reduce forgetting. However, the learning period of two tasks is too short

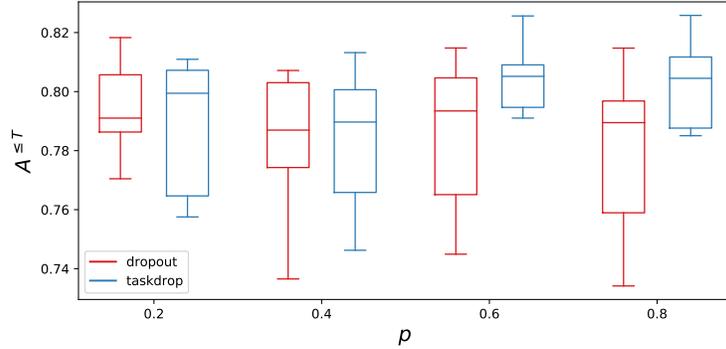


Figure 6: Comparison of TaskDrop and dropout in  $A^{\leq T}$  on low-6 with respect to different retention ratio  $p$ .

to allow any skip-task transfer. Nevertheless, TaskDrop achieves the best results on low-6 with respect to both short and long term learning. Because tasks of this dataset share less knowledge to each other, making full capacity sharing in No-masking less effective. The above observations indicate that even for sentiment classification tasks, randomly masking with a proper ratio still helps to learn a more robust and hence effective model from a relative long-term point of view.

### 4.3 Comparison with State-of-the-art Approaches

Table 3 compares TaskDrop with the four state-of-the-art continual learning approaches in the same way as above. Again, the results of TaskDrop after learning all the  $T$  tasks are the best on all three datasets.

The transfer-focused approach KAN outperforms the other three that dedicated for catastrophic forgetting. By comparing to results in Table 2, we found that No-Masking beats all of these state-of-the-art ones for most of the cases, except some cases on low-6, showing that the specific ways those approaches formulated for selective knowledge transfer are not suitable for datasets used here. We also compare different approaches via visualizing the learned representations in Figure 5, which plots the representations of test instances in the task of *Instant Video* with model trained on other tasks in low-6. It is seen that compact clusters are formed with features learned by TaskDrop. Two pairs of clusters with different class labels which are merged by other approaches are separated with a large margin by TaskDrop.

### 4.4 TaskDrop vs. Dropout

Now we compare the results of TaskDrop with the standard dropout. As seen from Figure 6 that TaskDrop performs better than dropout when  $p$  is not too small, i.e.,  $p \geq 0.4$ . This confirms that knowing the task boundary during training is important for the continual learning setting, which is not designed in dropout. The results of TaskDrop is slightly worse than dropout for  $p = 0.2$  as the ensemble of large number of subnetworks makes it less challenging for dropout to learning with a very sparse model.

### 4.5 Retention Ratio

Finally, we study the impact of the retention ratio  $p$ . Figure 4 plots the results of TaskDrop in  $A^{\leq T}$  with respect to different  $p$  for the three datasets. The three curves have a similar overall tendency, i.e., going up when increasing  $p$  from 0.2 and reaching the maximums with  $p$  between 0.6 and 0.8, and then slowly going down. Compared to the low-6 curve, the other two are much closer to each other in both shape and vertical position, indicating that low-6 has a more dissimilar nature from the other two datasets. The optimal value of  $p$  is smaller for low-6, which is not surprising as tasks of this dataset are less transferable. As a too small  $p$  causes the sampled encoder too simple to work reasonably, and also results in large standard deviations. Even though, with only 20% units of each layer, the results of TaskDrop are on par or better than those of Individual Networks.

## 5 Conclusion

We investigated a random masking-based approach for continual sentiment classification, where tasks are coming sequentially. The proposed approach gives competitive performance in our experiments without introducing extra

memory or learning modules, demonstrating the effectiveness of random masking based capacity allocation and reusing for the problem considered in this study. It will be interesting to further investigate this simple framework on problems from other domains to see whether the proposed random masking mechanism also works well for those problems that have different natures from sentiment classification tasks.

## References

- S. Legg and M. Hutter. Universal intelligence: a definition of machine intelligence. *Minds and Machines*, 17(4): 391–444, 2007.
- Sebastian Thrun and Tom M Mitchell. Lifelong robot learning. *Robotics and autonomous systems*, 15(1-2):25–46, 1995.
- M. McCloskey and N. Cohen. Catastrophic interference in connectionist networks: the sequential learning problem. *Psychology of Learning and Motivation*, pages 109–165, 1989.
- R. Ratcliff. Connectionist models of recognition memory: constraints imposed by learning and forgetting functions. *Psychological Review*, pages 285–308, 1990.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30:6467–6476, 2017.
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-GEM. In *ICLR*, 2019.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P.Lillicrap, and Greg Wayne. Experience replay for continual learning. In *NeurIPS*, 2019.
- Xu Han, Yi Dai, Tianyu Gao, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Continual relation learning via episodic memory activation and reconsolidation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6429–6440, Online, July 2020. Association for Computational Linguistics. doi:10.18653/v1/2020.acl-main.573. URL <https://aclanthology.org/2020.acl-main.573>.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *NIPS*, 2017.
- Cuong V Nguyen, Yingzhen Li, Thang D Bui, and Richard E Turner. Variational continual learning. 2018.
- Robert M French. Using semi-distributed representations to overcome catastrophic forgetting in connectionist networks. In *Proceedings of the 13th annual cognitive science society conference*, volume 1, pages 173–178, 1991.
- Xu He and Herbert Jaeger. Overcoming catastrophic interference using conceptor-aided backpropagation. In *International Conference on Learning Representations*, 2018.
- Gregory D. Grant Nicolas Y. Masse and David J. Freedman.
- Arun Mallya and Svetlana Lazebnik. Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7765–7773, 2018.
- Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra. Pathnet: Evolution channels gradient descent in super neural networks. *arXiv preprint arXiv:1701.08734*, 2017.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- X. Li, Y. Zhou, T. Wu, R. Socher, and C. Xiong. Learn to grow: A continual structure learning framework for overcoming catastrophic forgetting. 2019.
- Jathushan Rajasegaran, Munawar Hayat, Salman Khan, Fahad Shahbaz, and Khan Ling Shao. Random path selection for incremental learning. In *NeurIPS*, 2019.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.

- Joan Serra, Dídac Surís, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic forgetting with hard attention to the task. In *International Conference on Machine Learning*, pages 4548–4557. PMLR, 2018.
- Zixuan Ke, Bing Liu, Hao Wang, and Lei Shu. Continual learning with knowledge transfer for sentiment classification. In *ECML/PKDD (3)*, pages 683–698, 2020a.
- Zixuan Ke, Bing Liu, and Xingchang Huang. Continual learning of a mixed sequence of similar and dissimilar tasks. In *NeurIPS*, 2020b.
- Jaehong Yoon, Eunho Yang, Jeongtae Lee, and Sung Ju Hwang. Lifelong learning with dynamically expandable networks. *Proc. of the Int. Conf. on Learning Representations (ICLR)*, 2018.
- Wojciech Zaremba, Ilya Sutskever, and Oriol Vinyals. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.
- Sang-Woo Lee, Jin-Hwa Kim, Jaehyun Jun, Jung-Woo Ha, and Byoung-Tak Zhang. Overcoming catastrophic forgetting by incremental moment matching. *arXiv preprint arXiv:1703.08475*, 2017.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. *70(3987)*, 2017.
- Mohamed Elho seiny Marcus Rohrbach Rahaf Aljundi, Francesca Babiloni and Tinne Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, page 139–154, 2018.
- Hong Wang, Wenhan Xiong, Mo Yu, Xiaoxiao Guo, Shiyu Chang, and William Yang Wang. Sentence embedding alignment for lifelong relation extraction. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, June 2019.
- Amal Rannen Triki, Rahaf Aljundi, Matthew B. Blaschko, and Tinne Tuytelaars. Encoder based lifelong learning. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1329–1337. IEEE Computer Society, 2017. doi:10.1109/ICCV.2017.148. URL <https://doi.org/10.1109/ICCV.2017.148>.
- Y. Huang, Y. Zhang, J. Chen, X. Wang, and D. Yang. Continual learning for text classification with information disentanglement based regularization. In *NAACL*, 2021.
- J. Rajasegaran, S. Khan, M. Hayat, F. S. Khan, and M. Shah. iTAML: An incremental task-agnostic meta-learning approach. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9 (11), 2008.