

# MAP: Memory-Augmented Pre-trained Language Model for Domain Tasks

Anonymous EMNLP submission

## Abstract

Recently, domain-specific PLMs have been proposed to boost the task performance of specific domains (e.g., biomedical and computer science) by continuing to pre-train general PLMs with domain-specific corpora. However, this Domain-Adaptive Pre-Training (DAPT; Gururangan et al. (2020)) tends to forget the previous general knowledge acquired by general PLMs, which leads to a *catastrophic forgetting* phenomenon and sub-optimal performance. To alleviate this problem, we propose a new framework of **Memory-Augmented Pre-trained Language Model (MAP)**, which augments the domain-specific PLM by a memory representation built from the frozen general PLM without losing any general knowledge. Specifically, we propose a new memory-augmented layer, and based on it, different augmented strategies are explored to build the memory representation and then adaptively fuse it into the domain-specific PLM. We demonstrate the effectiveness of MAP on various domains (biomedical and computer science publications, news, and reviews) and different kinds (text classification, QA, NER) of tasks, and the extensive results show that the proposed MAP can achieve SOTA results on all tasks. The source code will be publicly available upon publication.

## 1 Introduction

Pre-trained Language models (PLMs), such as BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019b), have achieved promising performance on NLP tasks. Typically, these general models are firstly pre-trained on large unlabeled corpus and then directly fine-tuned on downstream tasks. However, there is an inherent gap in text distribution between unlabeled pre-training corpus and labeled task corpus, which leads to the distribution shift problem (Gururangan et al., 2020) and makes PLMs perform poorly on some domain tasks (Beltagy et al., 2019; Lee et al., 2020b). To

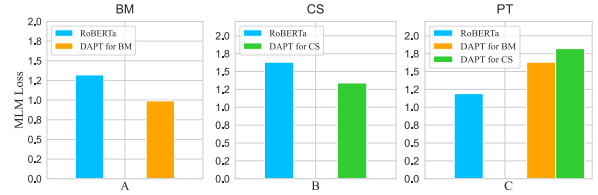


Figure 1: Masked LM (MLM) loss of RoBERTa on 50K randomly sampled documents from each domain before and after DAPT. Figure A and B denote the inference loss of general RoBERTa-base and domain-specific PLMs on the samples of biomedical (BM) and computer science (CS). Figure C means the loss of these models on the samples from the pre-training (PT) corpus of RoBERTa. We report the results of (Gururangan et al., 2020) and lower MLM loss is better.

address this shift problem, the domain-adaptive pre-training (DAPT) is proposed (Huang et al., 2019; Beltagy et al., 2019; Gururangan et al., 2020; Lee et al., 2020b) to further pretrain general PLMs on large-scale domain corpora, achieving better performance than general PLMs.

Although DAPT can effectively learn the domain distribution of the target task, its continual pretraining process updates the parameters of general PLMs, which inevitably leads to partial general knowledge being forgotten. This *catastrophic forgetting* (Goodfellow et al., 2014; Li and Hoiem, 2016; Thompson et al., 2019) phenomenon is verified in Figure 1, where we observe that the domain-specific PLMs show better results than general PLMs on domain corpus, but perform worse on the general corpus. We argue that this forgotten knowledge is beneficial for domain-specific PLMs and should be used to improve their generalization ability on domain tasks.

To alleviate the *catastrophic forgetting*, we propose a simple yet effective memory-augmented framework named general **Memory-Augmented Pre-trained model (MAP)**. In addition to the backbone domain-specific PLM, MAP introduces a new memory-augmented layer. It explicitly in-

incorporates the representation built from a frozen general PLM as the memory to make the backbone model access the complete general knowledge. Then, a new proposed memory-attention within the memory-augmented layer enables the domain-specific PLM adaptively combine the memory representation and the domain-specific representation. Using the memory built from the frozen general PLM has two advantages: (1) frozen PLM never suffers from forgetfulness since the parameters remain unchanged (Levine et al., 2022); (2) it doesn't require additional training for the general PLM during fine-tuning. However, building and fusing memory into a backbone model is essentially a many-to-many scenario, where we need to choose which layer output of the general PLM as the memory representation, and which layer in the domain-specific PLM should be fused. Thus, we propose several memory-augmented strategies for better building then combining the memory representation into domain-specific PLM.

We evaluate our MAP on text classification, Question Answering (QA) and Name Entity Recognition (NER) tasks covering four domains: biomedical science, computer science, news, and reviews. Experimental results demonstrate that MAP outperforms existing baselines on all tasks. We compare different memory-augmented strategies, and the results show that the proposed chunk-based gated memory transfer strategy achieves the best results. In addition, for the memory representation building, we empirically find that the freezing way is better than the unfreezing one, which also has better training efficiency. Furthermore, we apply the proposed framework to a small-scale domain pre-training setting and find that MAP is also practical in achieving lower MLM loss. Our contributions are summarized below:

- We empirically find that forgotten general knowledge due to *catastrophic forgetting* can benefit the domain-specific downstream tasks since it can improve PLMs' generalization ability.
- We propose a novel MAP framework, which introduces several memory-augmented strategies to construct the memory representation from the frozen general PLM effectively and then fuse it into the domain-specific PLM by a new memory-augmented layer.
- We conduct extensive experiments on various domain-specific tasks including text classification, QA, and NER, the results demonstrating

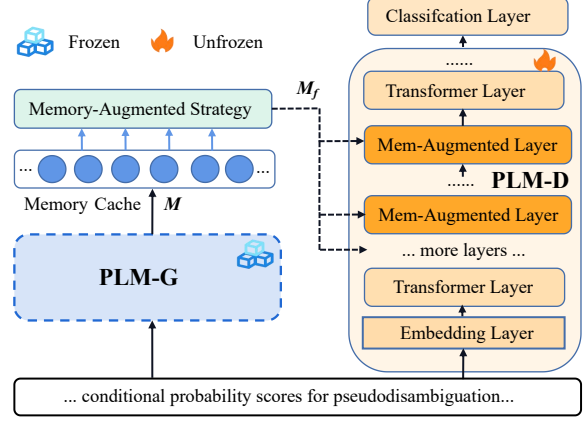


Figure 2: A framework of MAP with the cs-domain task input. PLM-G denotes the frozen general PLM, PLM-D denotes the domain-specific PLM.

that our MAP outperforms existing baselines.

## 2 The Method of MAP

In this section, we first overview the MAP framework. Then we detail the a new memory-augmented layer which fuses general knowledge into domain-specific PLMs. Finally, we propose different memory-augmented strategies including single-layer memory transfer, multiple-layer memory transfer, gated memory transfer and chunk-based gated memory transfer.

### 2.1 Overview

Our MAP framework aims to tackle the *catastrophic forgetting* of domain-specific PLMs by using the memory cache built from the frozen general PLMs, which is illustrated in Figure 2. Given an sequence  $\mathbf{x} = [x_1, x_2, \dots, x_t, \dots, x_n]$  with  $x_t$  denoting the  $t$ -th token, general PLMs output the contextual representations of the input tokens as the memory cache, which is fed into the domain-specific PLMs to build final representation for domain tasks:

$$\mathbf{M} = \text{PLM-G}(\mathbf{x}; \theta_g) \quad (1)$$

$$\mathbf{H} = \text{PLM-D}(\mathbf{x}, \mathbf{M}; \theta_d) \quad (2)$$

where  $\theta_g$  and  $\theta_d$  are the parameters of general and domain-specific PLMs respectively. We only update the  $\theta_d$  and the  $\theta_g$  is frozen when fune-tuning. The general PLM could be a BERT or RoBERTa, which contains  $l$  layers of Transformer (Vaswani et al., 2017) encoder blocks and outputs a set of hidden states denoted as a memory cache  $\mathbf{M} = \{\mathbf{M}^1, \mathbf{M}^2, \dots, \mathbf{M}^l\}$ . In the MAP

framework, the domain-specific PLM utilizes a new *memory-augmented layer* to adaptively incorporate the memory representation built from the memory cache  $M$  and enhance its generalization ability. Note that *memory-augmented layer* is built by only the parameters of original Transformer layer without adding new ones. Moreover, we explore different memory-augmented strategies and further propose the chunk-based memory transfer strategy, which fully uses the memories from different chunks of the general PLM.

## 2.2 Memory-Augmented Layer

Memory-augmented layer differs from traditional Transformer layer only in the multi-head self attention module. The new memory-augmented attention module is proposed to fuse the memory representation into the domain-specific PLM, denoted as memory-attention. Its main idea is to linearly transform the memory representation into new pairs of (keys, values) and concatenate them into the back of pairs produced by the domain-specific PLM. Then the multi-head self-attention is performed to adaptively fuse these new concatenated representation. The whole process reuses the parameters of the Transformer layer of domain-specific PLM and does not introduce any new parameters.

Specifically, if  $i$ -th Transformer layer is a memory-augmented one, it obtains the domain-specific representation  $H_{i-1}$  from previous layer and the memory representation  $M_f$  as the input and fuse them by the following way:

$$\text{Memory-Attention}(H_{i-1}, M_f) = \text{Concat}(\text{head}_1, \dots, \text{head}_k) W^o \quad (3)$$

where  $M_f$  is the memory representation directly extracted from the memory cache  $M$  or effectively constructed by some adaptive aggregation strategies, which has the same shape as the intermediate hidden state  $H_i$  of the domain-specific PLM,  $k$  means the number of heads and  $W^o$  is a trainable parameter matrix. Then,  $M_f$  is linearly transformed into new pairs of (keys, values) which were appended to the last of domain-specific ones:

$$\tilde{K}_{i,j} = \text{Concat}(K_{i,j}, M_f^k) \quad (4)$$

$$\tilde{V}_{i,j} = \text{Concat}(V_{i,j}, M_f^v) \quad (5)$$

$$Q_{i,j}, K_{i,j}, V_{i,j} = H_{i-1} W_{i,j}^q, H_{i-1} W_{i,j}^k, H_{i-1} W_{i,j}^v \quad (6)$$

$$M_f^k, M_f^v = M_f W_{i,j}^k, M_f W_{i,j}^v \quad (7)$$

where  $W_{i,j}^q$ ,  $W_{i,j}^k$  and  $W_{i,j}^v$  are trainable parameters to generate queries, keys, values respectively, and  $j$  refers to  $j$ -th attention head. Then the self-attention is performed on the queries and merged pairs of (keys, values) as follow:

$$\text{head}_j = \text{Softmax}\left(\frac{Q_{i,j} \tilde{K}_{i,j}^T}{\sqrt{d_k}}\right) \tilde{V}_{i,j} \quad (8)$$

where  $d_k$  is the head dimension acting as a scaling factor. Firstly, a unified attention matrix is computed by the standard scaled dot-product of each query against the keys of general memory and the domain-specific keys. Then, a softmax operation gets the normalized scores that weigh and sum these concatenated values. Without additional parameter update for the general PLM, domain-specific PLM can dynamically capture useful general knowledge and ignore noisy information through the memory-augmented layer.

## 2.3 Memory-Augmented Strategies

The remaining problem is how to build the memory representation  $M_f$  from the memory cache  $M$  and which layer of the domain-specific PLM should be the memory-augmented layer to fuse  $M_f$ . Essentially, it is a many-to-many layer assignment problem between the general PLM and the domain-specific PLM. To study the effect of layer assignment, we propose and compare different strategies, as shown in Figure 3.

**Single-Layer Memory Transfer** We first consider a single-layer memory transfer approach, where the last hidden state of the memory cache  $M$  is extracted as  $M_f$  and then it is fused into one layer of domain-specific PLM with memory-attention. We choose the layer near the top of the domain-specific PLM model as the memory-augmented layer which performs best in the experiment. This strategy does not require additional parameters.

**Multiple-Layer Memory Transfer** The single-layer memory transfer may ignore the knowledge learned from shallow layers of the general PLM. To perform layer-wise interaction between the general PLM and the domain-specific PLM, we propose a multiple-layer transfer strategy. This strategy leverages all hidden states from the memory cache  $M$  as the memory representations and then fuses them into the corresponding layers of the domain-specific PLM, which also does not introduce any new parameters.

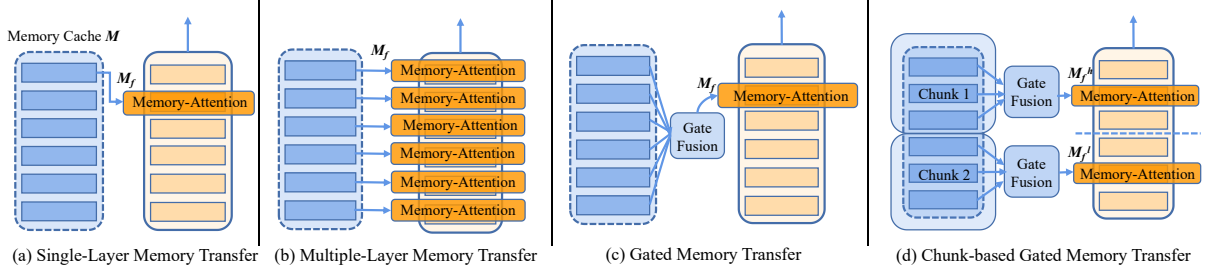


Figure 3: Memory-augmented strategies of the MAP framework. We take a 6-layer model as an example.

**Gated Memory Transfer** Multiple-layer memory transfer uses the hidden states output by all layers of the frozen general PLM as the memory representations, which inevitably introduces homogeneous and noise information. To avoid the problem, we further propose the strategy of gated memory transfer, which firstly exploits the token-level gate mechanism to adaptively weigh and sum representations of different layers into a memory representation, and then it will be fused into one layer of domain-specific PLM. We also choose the layer near the top of the domain-specific PLM as the memory-augmented one which achieves optimal performance in the experiment. The gate fusion mechanism is formulated as below:

$$\mathbf{m}_f^t = \sum_{l=1}^L \alpha_l^t \mathbf{m}_l^t \quad (9)$$

$$\alpha_l^t = \frac{\exp(g(\mathbf{m}_l^t))}{\sum_{i=1}^L \exp(g(\mathbf{m}_i^t))} \quad (10)$$

$$\mathbf{M}_f = \{\mathbf{m}_f^1, \mathbf{m}_f^2, \dots, \mathbf{m}_f^t, \dots, \mathbf{m}_f^n\} \quad (11)$$

where  $\mathbf{m}_l^t$  is the token representation,  $t$  denotes the token index,  $l$  is the layer index,  $n$  is the length of tokens and  $g$  refers to a linear layer. We utilize a softmax function to calculate the importance of tokens in different layers. Therefore, the output token representation  $\mathbf{m}_f^t$  is obtained by weighing the  $t$ -th token representations from different layers with their corresponding importance  $\alpha_l^t$ . Finally, the built memory representation  $\mathbf{M}_f$  is fused to the memory-augmented layer.

**Chunk-based Gated Memory Transfer** The previous work (Liu et al., 2019a; Phang et al., 2021) have observed that the representations from upper layers and lower layers of pre-trained language models are significantly different. Motivated by this observation, based on the gated memory transfer, we further propose a chunk-based variant, which separates the layers of general PLM

into a high-level chunk and a low-level chunk, and then apply the gate fusion strategy to get upper and lower-layer memory representations  $\mathbf{M}_f^h$  and  $\mathbf{M}_f^l$ , respectively. Finally, we fuse them into two memory-augmented layers in the domain-specific PLM, and the details of different layer selections for this strategy is presented in Section 4.3.

### 3 Experiments

In this section, we first introduce the evaluation tasks and metrics. Then we illustrate the baseline methods, implementation settings. Finally, we conduct the experimental analysis of MAP.

#### 3.1 Datasets and Metrics

**Datasets** We evaluate our model on three tasks: text classification, QA and NER. For text classification, we conduct experiments on eight tasks that cover four domains, including CHEMPROT (Kringelum et al., 2016) and RCT (Dernoncourt and Lee, 2017) in the biomedical domain, ACL-ARC (Jurgens et al., 2018) and SCI-ERC (Luan et al., 2018) in the computer science domain, HYPERPARTISAN (Kiesel et al., 2019) and AGNEWS (Zhang et al., 2015) in the news domain, HELPFULNESS (McAuley et al., 2015) and IMDB (Maas et al., 2011) in the reviews domain. In addition, we use micro-F1 as the metric for ChemProt and RCT, and use macro-F1 for the other datasets following (Gururangan et al., 2020). For NER, we use two datasets including NCBI-Disease (Dogan et al., 2014) in the biomedical domain, CoNNL-2003 (Sang and Meulder, 2003) in the news domain. We use the F1 score as the evaluation metric. For QA, we utilize two datasets including Medication (Pampari et al., 2018) in the biomedical domain, NewsQA (Trischler et al., 2017) in the news domain. We use the Exact-Match (EM) and the F1 score as the evaluation metrics. The detailed description and statistics of each tasks are shown in Appendix A.



Domain	BIOMED		CS		NEWS		REVIEWS	
Dataset	CP	RCT	CI	SE	HP	AG	AM	IMDB
<b>Fine-Tuning</b>	81.9 <sub>0.1</sub>	87.2 <sub>0.1</sub>	63.0 <sub>5.8</sub>	77.3 <sub>1.9</sub>	86.6 <sub>0.9</sub>	93.9 <sub>0.2</sub>	65.1 <sub>3.4</sub>	95.0 <sub>0.2</sub>
<b>DAPT</b>	84.2 <sub>0.2</sub>	87.6 <sub>0.1</sub>	75.4 <sub>2.5</sub>	80.8 <sub>1.5</sub>	88.2 <sub>5.9</sub>	93.9 <sub>0.2</sub>	66.5 <sub>1.4</sub>	95.4 <sub>0.1</sub>
<b>Logits Fusion</b>	84.4 <sub>0.3</sub>	87.6 <sub>0.1</sub>	77.5 <sub>2.5</sub>	82.9 <sub>1.3</sub>	91.8 <sub>2.5</sub>	93.8 <sub>0.1</sub>	67.6 <sub>0.2</sub>	95.1 <sub>0.2</sub>
<b>Ensemble LMs</b>	84.6 <sub>0.2</sub>	87.4 <sub>0.2</sub>	76.2 <sub>2.2</sub>	82.6 <sub>1.0</sub>	91.4 <sub>2.8</sub>	93.9 <sub>0.2</sub>	67.5 <sub>1.1</sub>	95.1 <sub>0.2</sub>
<i>MAP (Single-Layer Memory Transfer)</i>	84.7 <sub>0.3</sub>	87.8 <sub>0.1</sub>	78.6 <sub>2.3</sub>	83.2 <sub>0.6</sub>	93.3 <sub>1.6</sub>	94.0 <sub>0.1</sub>	68.5 <sub>0.6</sub>	95.4 <sub>0.1</sub>
<i>MAP (Multiple-Layer Memory Transfer)</i>	84.7 <sub>0.4</sub>	87.7 <sub>0.1</sub>	74.9 <sub>1.3</sub>	82.0 <sub>1.1</sub>	92.2 <sub>2.3</sub>	93.8 <sub>0.2</sub>	67.6 <sub>0.6</sub>	95.2 <sub>0.2</sub>
<i>MAP (Gated Memory Transfer)</i>	84.8 <sub>0.2</sub>	87.8 <sub>0.1</sub>	77.9 <sub>1.5</sub>	83.2 <sub>1.3</sub>	93.2 <sub>2.0</sub>	94.0 <sub>0.2</sub>	68.1 <sub>0.5</sub>	95.4 <sub>0.2</sub>
<i>MAP (Chunk-based Gated Memory Transfer)</i>	<b>85.0</b> <sub>0.3</sub>	<b>87.9</b> <sub>0.2</sub>	<b>79.4</b> <sub>1.0</sub>	<b>83.8</b> <sub>0.8</sub>	<b>95.2</b> <sub>0.0</sub>	<b>94.1</b> <sub>0.2</sub>	<b>69.0</b> <sub>0.8</sub>	<b>95.6</b> <sub>0.2</sub>

Table 1: The comparison against baselines on text classification tasks and performance of the proposed memory-augmented strategies. We report the averages across five random seeds, with standard deviations as subscripts. The best performance for each benchmark is marked in black bold. CP, CI, SE, HP, AG and AM denote CHEMPROT, ACL-ARC, SCIERC, HYPERPARTISAN, AGNEWS and AMAZON, respectively.

Domain	BIOMED		CS		NEWS		REVIEWS	
Dataset	CP	RCT	CI	SE	HP	AG	AM	IMDB
<b>TAPT</b>	82.6 <sub>0.4</sub>	87.7 <sub>0.1</sub>	67.4 <sub>1.8</sub>	79.3 <sub>1.5</sub>	90.4 <sub>5.2</sub>	94.5 <sub>0.1</sub>	68.5 <sub>1.9</sub>	95.5 <sub>0.1</sub>
<b>DAPT+TAPT</b>	84.4 <sub>0.4</sub>	87.8 <sub>0.1</sub>	75.6 <sub>3.8</sub>	81.3 <sub>1.7</sub>	90.0 <sub>6.6</sub>	94.6 <sub>0.2</sub>	68.7 <sub>1.8</sub>	95.6 <sub>0.1</sub>
<b>MAP*</b>	<b>85.1</b> <sub>0.3</sub>	<b>87.9</b> <sub>0.0</sub>	<b>79.6</b> <sub>2.2</sub>	<b>83.9</b> <sub>1.1</sub>	<b>95.2</b> <sub>0.0</sub>	<b>94.6</b> <sub>0.1</sub>	<b>69.9</b> <sub>0.3</sub>	<b>95.8</b> <sub>0.1</sub>

Table 2: The experimental results of MAP compared with TAPT and DAPT+TAPT on domain classification tasks. \* means MAP framework built with the backbone PLM pre-trained with the process of DAPT then TAPT.

### 3.2 Baselines

In our experiment, all the baselines are built on the RoBERTa-base. The details baselines of text classification tasks are described as follows:

- **Fine-Tuning**: directly fine-tuning the general PLM for downstream tasks.
- **DAPT** (Gururangan et al., 2020): pre-training the general PLM with large-scale domain unlabeled corpora to get the domain-specific PLM then fine-tuning it.
- **Logits Fusion**: a straight-forward method combines the frozen general PLM and the domain-specific PLM by adding their logits. This method is optimized end-to-end and does not include any memory-augmented strategies in the model.
- **Ensemble LMs**: an ensemble method that adds the predicted probabilities of the fine-tuned general and the domain-specific PLMs for final prediction.
- **TAPT** (Gururangan et al., 2020): task-adaptive pretraining continues to pre-train the PLM on the training dataset, and then we fine-tune it for the downstream tasks.

For the NER and QA tasks, in addition to these above baselines, we also introduce KALA (Kang et al., 2022). Since KALA is only verified in QA

and NER, we use it as a baseline for these two tasks.

- **KALA**: constructing an entity memory and knowledge graph on task-specific domain and then augmenting PLM by these additional knowledge.

### 3.3 Implementation

We implement the MAP framework based on RoBERTa-base. For the domain-specific PLM, we use the released pre-trained weights DAPT<sup>1</sup>. The more details on fine-tuning of the downstream tasks are shown in Appendix B.

### 3.4 Results and Analysis

Our experiment results on the domain-specific classification tasks are shown in Table 1 and 2, the results of QA and NER tasks are shown on Table 3.

#### Performance on Domain Classification Tasks

From Table 1, we can observe that MAP with the proposed chunk-based gated memory transfer can achieve better results than all the baselines, which proves that incorporating memory from the general frozen PLM is beneficial for the domain-specific PLM. Specifically, the strategy of chunk-based gate

<sup>1</sup><https://github.com/allenai/dont-stop-pretraining>

Domain	BIOMED		NEWS	
Dataset	Medication (QA)	NCBI-Disease (NER)	NewsQA (QA)	CoNNL-2003 (NER)
<b>Fine-Tuning</b>	26.9 <sub>0.2</sub>   71.5 <sub>0.5</sub>	86.9 <sub>1.0</sub>	57.2 <sub>0.6</sub>   71.9 <sub>0.4</sub>	95.6 <sub>0.2</sub>
<b>TAPT</b>	27.0 <sub>0.2</sub>   71.2 <sub>0.3</sub>	86.2 <sub>0.8</sub>	57.2 <sub>0.5</sub>   71.8 <sub>0.3</sub>	95.6 <sub>0.3</sub>
<b>DAPT</b>	27.2 <sub>0.3</sub>   71.4 <sub>0.4</sub>	87.2 <sub>0.4</sub>	58.7 <sub>0.6</sub>   72.4 <sub>0.4</sub>	95.8 <sub>0.2</sub>
<b>KALA</b>	27.3 <sub>0.4</sub>   71.1 <sub>0.5</sub>	87.7 <sub>0.3</sub>	58.0 <sub>0.6</sub>   72.7 <sub>0.3</sub>	95.4 <sub>0.3</sub>
<b>MAP</b>	<b>29.1<sub>0.3</sub>   72.2<sub>0.4</sub></b>	<b>88.7<sub>0.2</sub></b>	<b>59.9<sub>0.8</sub>   73.3<sub>0.4</sub></b>	<b>96.2<sub>0.2</sub></b>

Table 3: The experimental results of extractive QA and NER tasks in biomedical and news domains. We use Exact Match and F1 score as the metrics for the QA tasks: Medication and NewsQA, and F1 score for NER tasks: NCBI-Disease and CoNNL-2003.

memory transfer outperforms other strategies, we conjecture that it adaptively selects the token-level information across different layers and adequately utilizes the general knowledge from both the high-level and low-level chunks. However, we also observe that multiple-layer memory transfer has little improvement compared with the baselines, because it incorporates excessive redundant and noisy information from the general PLM without the proposed gated fusion. Besides, single-layer memory transfer is a simple yet effective strategy achieving better than the baselines and the non-gated fusion strategy of multiple-layer memory transfer.

Since the chunk-based gate memory transfer strategy achieves the best performance compared with the baselines, we use it as the default memory-augmented strategy within MAP in the following experiments.

**Comparison with Further TAPT** Further task-adaptive pre-training (TAPT) has been proven to improve the domain-adaptive pre-training (DAPT) (Gururangan et al., 2020). To demonstrate the effectiveness of MAP on TAPT, we build a MAP framework that replaces the domain-specific pre-trained PLM with the task-adaptive pre-trained PLM. From Table 2, we find that our MAP also outperforms DAPT+TAPT on all datasets, indicating that the proposed framework is general for different backbone models, including the domain-adaptive and the task-adaptive PLMs.

**Effectiveness for QA and NER** We also evaluate MAP on the tasks of QA and NER, and the experiment results are shown in Table 3. From the results, we see that our method achieves better results than the baselines on all datasets, especially KALA (Kang et al., 2022), which spends a considerable effort to construct entity memory and knowledge graph from the contexts. These results

Domain	BIOMED		CS	
Dataset	CP	RCT	CI	SE
<b>DAPT</b>	84.2 <sub>0.2</sub>	87.6 <sub>0.1</sub>	75.4 <sub>2.5</sub>	80.8 <sub>1.5</sub>
<b>MAP w/o Frozen</b>	84.3 <sub>0.5</sub>	87.3 <sub>0.3</sub>	78.1 <sub>2.2</sub>	82.4 <sub>1.1</sub>
<b>MAP</b>	<b>85.0<sub>0.3</sub></b>	<b>87.9<sub>0.2</sub></b>	<b>79.4<sub>1.0</sub></b>	<b>83.8<sub>0.8</sub></b>

Table 4: Results of utilizing the frozen and unfrozen general PLMs.

further demonstrate the effectiveness of MAP.

## 4 Further Discussion

In the following sections, we conduct some detailed analysis of MAP to demonstrate the effectiveness of the frozen general PLM and memory-attention. Moreover, we apply the proposed framework in the pre-training stage and study the effect of layer selection on the performance.

### 4.1 Effectiveness of Frozen Memory

We compare the frozen and unfrozen ways when building the general memory representation, and the results are shown in Table 4. From the table, we observe that the frozen method is better than the unfrozen model on all datasets, and both are better than the baseline DAPT. We argue that using the frozen memory has two advantages: (1) more efficient in model training without updating the parameters of the general PLM; (2) keeps general knowledge from PLM unchanged when fine-tuning, so it does not lead to a forgetting problem.

### 4.2 Effectiveness of Memory-Attention

To study the effectiveness of our proposed memory-attention module, we compare it with other attention-based variants within a memory-augmented layer. Specifically, cross-attention is an attention module widely-used in the multi-modal learning (Li et al., 2021; Zeng et al., 2021), we

Domain	BIOMED		CS	
Dataset	CP	RCT	CI	SE
Cross-Attention	82.8 <sub>0.4</sub>	86.6 <sub>0.1</sub>	73.5 <sub>2.0</sub>	82.2 <sub>1.0</sub>
Gate-Attention	84.7 <sub>0.3</sub>	87.6 <sub>0.1</sub>	78.1 <sub>1.4</sub>	83.4 <sub>0.7</sub>
Memory-Attention	<b>85.0<sub>0.3</sub></b>	<b>87.9<sub>0.2</sub></b>	<b>79.4<sub>1.0</sub></b>	<b>83.8<sub>0.8</sub></b>

Table 5: Performance of using the proposed memory-attention and other main-stream attention-based variants in MAP framework. We evaluate them on the datasets of biomedical and cs domains.

apply it to adaptively fuse the memory representation  $M_f$  and the output representation from the self-attention module. We also include the gate-attention (Wu et al., 2022) as the fusion baseline, which utilizes a gate mechanism to weigh and sum the local and external-memory for long-sequence modeling. As shown in Table 5, our memory-attention module outperforms other variants without additional trainable parameters.

### 4.3 Layer Selection for Memory-Attention

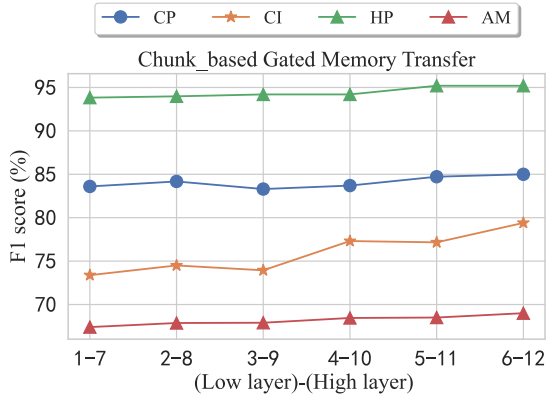


Figure 4: Performance of different layer selections in chunk-based gate memory transfer strategy.

Besides the strategy of multiple-layer memory transfer, other strategies need to do the layer selection. For the strategies of single and gated-memory transfer, we fuse the memory representation  $M_f$  into different layers  $\{3, 6, 9, 12\}$  in the domain-specific PLM and find that the layer 9 as the memory-augmented layer can achieve the best performance in both strategies. We present more detailed results in Appendix C. For chunk-based gate memory transfer strategy, we experiment with transferring the memory representation of the high-level chunk into layer 7 to 12, and the other one of the low-level chunk into layer 1 to 6, which keeps a same layer interval between the two memory-augmented layers in the domain-specific PLM. The

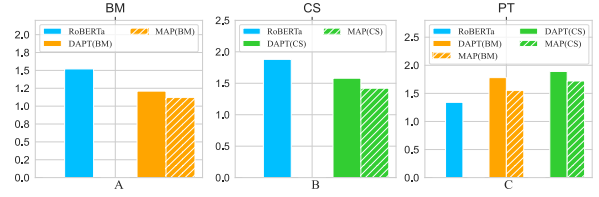


Figure 5: Masked LM loss for the pre-training stage (lower value is better). PT denotes the samples similar to RoBERTa’s pre-training corpus. DAPT(BM) denotes the domain-specific PLM for biomedical domain. MAP(BM) denotes the MAP framework with the biomedical-domain backbone. For instance, figure A represents further pre-training the models on the biomedical pre-training samples then inferring their MLM loss on the test samples.

experimental results are shown in Figure 4. The results show that there is an increasing tendency when placing memory-augmented layers to the top of the domain-specific PLM. Finally, we choose the layer 6 and 12 as the memory-augmented ones for the strategy of chunk-based gated memory transfer.

### 4.4 Apply MAP in the Pre-training Stage

In the previous experiment, we have incorporated the domain-specific PLM with the MAP framework in fine-tuning stage. In this section, we further study whether MAP is beneficial for pre-training stage. To this end, we randomly samples 50k documents from general<sup>2</sup>, biomedical and computer science domains (Lo et al., 2020), respectively. In addition, we randomly split 70% of the data from each domain as the pre-training samples and the rest data as the test samples. More details about pre-processing the domain samples are shown in Appendix B. Then, we pre-train the models on the pre-training samples and calculate the masked LM loss on the test samples. From the experiment results shown in Figure 5, compared with baseline DAPT, we observe that utilizing MAP can reduce masked LM loss on the biomedical, cs and general domain. These results demonstrate that the proposed MAP also mitigates catastrophic forgetting during the adaptive pre-training.

## 5 Related work

**Domain Adaptation for PLMs** Recently, the domain shift problem of PLMs has attracted increasing research (Beltagy et al., 2019; Huang et al., 2019; Lee et al., 2020b; Gururangan et al., 2020) since the domain discrepancies between

<sup>2</sup><https://github.com/soskek/bookcorpus>

the pre-training corpora and the downstream tasks can lead to a significant performance drop. To bridge the domain gap, SciBERT (Beltagy et al., 2019) and BioBERT (Lee et al., 2020b) further pre-train BERT with 1.14M scientific papers from Semantic Scholar corpus and biomedical documents, respectively, which can improve the performance of domain-specific NLU tasks compared with general BERT. Also, Gururangan et al. (2020) proposed domain-adaptive pre-training (DAPT) and task-adaptive pre-training (TAPT). Concretely, DAPT continues to pre-train the PLM with domain-specific corpora, while TAPT directly pre-trains the PLM on the task dataset. Moreover, BT-TAPT (Lee et al., 2021) inherits the crucial step of TAPT and leverages the back-translated strategy to augment the task data to improve the performance of PLM. TAPTER (Nishida et al., 2021) equips TAPT with domain-specific word embedding regularization to improve fine-tuning performance. However, above approaches suffer from catastrophic forgetting of general domain knowledge after adaptive pre-training, which leads to sub-optimal performance on downstream tasks.

**Catastrophic Forgetting** *Catastrophic forgetting* is a common phenomenon for continual learning, and it occurs when a training model forgets previously learned knowledge and over-fits to new tasks (McCloskey and Cohen, 1989). Typically, regularization-based methods (Goodfellow et al., 2014; Kirkpatrick et al., 2016; Serrà et al.) exploit regularization to constrain the parameter update to alleviate the forgetting problem, and the memory-based methods (Guo et al., 2020; Saha et al., 2021) mitigate forgetting by storing important samples from past tasks in the external memory and rehearsing them via some gradient transformation strategies. In addition, plenty of works have been proposed to address *catastrophic forgetting* for NLP tasks. Dakwale and Monz. (2017) subtly minimized KL-divergence of prediction losses as a regularization term between fine-tuning and general domain models. Lee et al. (2020a) introduced a new regularization technique to mix the PLM parameters with vanilla parameters instead of stochastic dropout. Chen et al. (2020) adopted multi-task learning to jointly learn pre-training and downstream tasks with less forgetting during fine-tuning. Xie et al. (2021) preserved the model neurons of general and language-specific parts during fine-tuning. However, our method is orthogonal

to the above approaches since we aim to effectively incorporate the domain-specific PLM with the memory representation built from the frozen general PLM to solve the forgetting issue without adding additional regularization terms in the model or using external memory to preserve samples from the past tasks.

**Knowledge-Enhanced PLMs** Knowledge-enhanced methods have shown effectiveness for PLMs via introducing internal or external knowledge. To improve the performance of fine-tuning tasks, REINA (Wang et al., 2022) retrieves the labeled training instances most similar to the input data and concatenates them before feeding them into PLMs. Besides, RETRO (Borgeaud et al., 2021) enhances the auto-regressive language model via leveraging a pre-trained frozen BERT model to retrieve related texts and then use a chunked cross-attention module to incorporate them. Memorizing transformer (Wu et al., 2022) leverages a learned gate to combine the attention results of the local context and the external context retrieved from previously seen sub-sequences. KALA (Kang et al., 2022) is the approach most relevant to our work. It incorporates intermediate hidden representations with domain-specific entities and their relational facts during task-specific fine-tune for domain tasks. However, our method doesn’t need to retrieve similar texts or construct additional knowledge graphs. We propose several memory-augmented strategies to build the memory representation and then transfer it into the domain-specific PLM to mitigate the forgetting of general knowledge.

## 6 Conclusion

In this work, we propose MAP, a novel framework that utilizes the memory-augmented layer to fuse the memory representation built from the frozen general PLM to mitigate *catastrophic forgetting* of general knowledge caused by domain-adaptive pre-training. We explore different memory-augmented strategies to construct the memory representation and empirically find that chunk-based gate memory transfer achieves the most optimal performance. We validate MAP on various domains of classification, QA, and NER tasks. The results show that our method consistently outperforms existing baselines on all datasets, implying that explicitly leveraging forgotten general knowledge is beneficial for domain-specific downstream tasks.



## 7 Limitations

Our MAP framework has been validated on domain-specific tasks and a small-scale domain pre-training experiment in Section 4.4. Due to the lack of large GPU resource, we have not validated our MAP framework in large-scale pre-training, a more challenging setting that we leave as future work. We also consider automatic layer selection to be an under-studied problem and believe that AutoML techniques (Pham et al., 2018; Tan and Le, 2019), such as evolutionary search (Deb et al., 2002; Chen et al., 2021), will be promising methods. Finally, the proposed framework is built on the encoder-only model, RoBERTa-base. In the future, we will apply our framework on the other types of architectures, such as decoder-only GPT (Radford et al., 2018) and encoder-decoder BART (Lewis et al., 2020).

## References

Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. *Scibert: A pretrained language model for scientific text*. In *EMNLP*.

Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. 2021. *Improving language models by retrieving from trillions of tokens*. *CoRR*, abs/2112.04426.

Huangke Chen, Ran Cheng, Witold Pedrycz, and Yaochu Jin. 2021. *Solving many-objective optimization problems via multistage evolutionary search*. *IEEE Trans. Syst. Man Cybern. Syst.*

Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. 2020. *Recall and learn: Fine-tuning deep pretrained language models with less forgetting*. In *EMNLP*.

Praveen Dakwale and Christof Monz. 2017. *Fine-tuning for neural machine translation with limited degradation across in-and out-of-domain data*. In *Proceedings of the XVI Machine Translation Summit*.

Kalyanmoy Deb, Samir Agrawal, Amrit Pratap, and T. Meyarivan. 2002. *A fast and elitist multiobjective genetic algorithm: NSGA-II*. *IEEE Trans. Evol. Comput.*

Franck Dernoncourt and Ji Young Lee. 2017. *Pubmed 200k rct: a dataset for sequential sentence classification in medical abstracts*. In *IJCNLP*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. *BERT: pre-training of deep bidirectional transformers for language understanding*. In *NAACL-HLT*.

Rezarta Islamaj Dogan, Robert Leaman, and Zhiyong Lu. 2014. *Ncbi disease corpus: A resource for disease name recognition and concept normalization*. *Journal of biomedical informatics*.

Ian J. Goodfellow, Mehdi Mirza, Xia Da, Aaron C. Courville, and Yoshua Bengio. 2014. *An empirical investigation of catastrophic forgetting in gradient-based neural networks*. In *ICLR*.

Yunhui Guo, Mingrui Liu, Tianbao Yang, and Tazjana Rosing. 2020. *Improved schemes for episodic memory-based lifelong learning*. In *NIPS*.

Suchin Gururangan, Ana Marasovic, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. *Don’t stop pretraining: Adapt language models to domains and tasks*. In *ACL*.

Kexin Huang, Jaan Altosaar, and Rajesh Ranganath. 2019. *Clinicalbert: Modeling clinical notes and predicting hospital readmission*. *CoRR*, abs/1904.05342.

David Jurgens, Srikanth Kumar, Raine Hoover, Daniel A. McFarland, and Dan Jurafsky. 2018. *Measuring the evolution of a scientific field through citation frames*. *Transactions of the Association for Computational Linguistics*.

Minki Kang, Jinheon Baek, and Sung Ju Hwang. 2022. *KALA: knowledge-augmented language model adaptation*. In *NAACL*.

Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, D. Corney, Benno Stein, and Martin Potthast. 2019. *Semeval-2019 task 4: Hyperpartisan news detection*. In *\*SEMEVAL*.

James Kirkpatrick, Razvan Pascanu, Neil C. Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. 2016. *Overcoming catastrophic forgetting in neural networks*. *CoRR*, abs/1612.00796.

Jens Kringelum, Sonny Kim Kjærulff, Søren Brunak, Ole Lund, Tudor I. Oprea, and Olivier Taboureaux. 2016. *Chemprot-3.0: a global chemical biology diseases mapping*. *Database J. Biol. Databases Curation*.

Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. 2020a. *Mixout: Effective regularization to fine-tune large-scale pretrained language models*. In *ICLR*.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020b. *Biobert: a pre-trained biomedical language representation model for biomedical text mining*. *Bioinform.*

714	Junghoon Lee, Joungee Kim, and Pilsung Kang. 2021.	Kosuke Nishida, Kyosuke Nishida, and Sen Yoshida.	768
715	<a href="#">Back-translated task adaptive pretraining: Improving</a>	2021. <a href="#">Task-adaptive pre-training of language models</a>	769
716	<a href="#">accuracy and robustness on text classification.</a> <i>CoRR</i> ,	<a href="#">with word embedding regularization.</a> In <i>Findings of</i>	770
717	<a href="#">abs/2107.10474.</a>	<i>ACL.</i>	771
718	Yoav Levine, Itay Dalmedigos, Ori Ram, Yoel Zeldes,	Anusri Pampari, Preethi Raghavan, Jennifer J. Liang,	772
719	Daniel Jannai, Dor Muhlgay, Yoni Osin, Opher	and Jian Peng. 2018. <a href="#">emrqa: A large corpus for</a>	773
720	Lieber, Barak Lenz, Shai Shalev-Shwartz, Amnon	<a href="#">question answering on electronic medical records.</a> In	774
721	Shashua, Kevin Leyton-Brown, and Yoav Shoham.	<i>EMNLP.</i>	775
722	2022. <a href="#">Standing on the shoulders of giant frozen lan-</a>	Hieu Pham, Melody Y. Guan, Barret Zoph, Quoc V. Le,	776
723	<a href="#">guage models.</a> volume <a href="#">abs/2204.10019.</a>	and Jeff Dean. 2018. <a href="#">Efficient neural architecture</a>	777
724	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan	<a href="#">search via parameter sharing.</a> In <i>ICML</i> , Proceedings	778
725	Ghazvininejad, Abdelrahman Mohamed, Omer Levy,	of Machine Learning Research.	779
726	Veselin Stoyanov, and Luke Zettlemoyer. 2020.	Jason Phang, Haokun Liu, and Samuel R. Bowman.	780
727	<a href="#">BART: denoising sequence-to-sequence pre-training</a>	2021. <a href="#">Fine-tuned transformers show clusters of</a>	781
728	<a href="#">for natural language generation, translation, and com-</a>	<a href="#">similar representations across layers.</a> In <i>Black-</i>	782
729	<a href="#">prehension.</a> In <i>ACL.</i>	<i>boxNLP@EMNLP.</i>	783
730	Junnan Li, Ramprasaath R. Selvaraju, Akhilesh	Alec Radford, Karthik Narasimhan, Tim Salimans, and	784
731	Gotmare, Shafiq R. Joty, Caiming Xiong, and	Ilya Sutskever. 2018. <a href="#">Improving language under-</a>	785
732	Steven Chu-Hong Hoi. 2021. <a href="#">Align before fuse: Vi-</a>	<a href="#">standing by generative pre-training.</a>	786
733	<a href="#">sion and language representation learning with mo-</a>	Gobinda Saha, Isha Garg, and Kaushik Roy. 2021. <a href="#">Gra-</a>	787
734	<a href="#">mentum distillation.</a> In <i>NIPS.</i>	<a href="#">dient projection memory for continual learning.</a> In	788
735	Zhizhong Li and Derek Hoiem. 2016. <a href="#">Learning without</a>	<i>ICLR.</i>	789
736	<a href="#">forgetting.</a> In <i>ECCV.</i>	Erik Tjong Kim Sang and Fien De Meulder. 2003. <a href="#">In-</a>	790
737	Nelson F. Liu, Matt Gardner, Yonatan Belinkov,	<a href="#">troduction to the conll-2003 shared task: Language-</a>	791
738	Matthew E. Peters, and Noah A. Smith. 2019a. <a href="#">Lin-</a>	<a href="#">independent named entity recognition.</a> In <i>CoNLL.</i>	792
739	<a href="#">guistic knowledge and transferability of contextual</a>	Joan Serrà, Didac Suris, Marius Miron, and Alexandros	793
740	<a href="#">representations.</a> In <i>NAACL-HLT</i> , pages 1073–1094.	Karatzoglou. <a href="#">Overcoming catastrophic forgetting</a>	794
741	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	<a href="#">with hard attention to the task.</a> In <i>ICML.</i>	795
742	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,	Mingxing Tan and Quoc V. Le. 2019. <a href="#">Efficientnet: Re-</a>	796
743	Luke Zettlemoyer, and Veselin Stoyanov. 2019b.	<a href="#">thinking model scaling for convolutional neural net-</a>	797
744	<a href="#">Roberta: A robustly optimized BERT pretraining</a>	<a href="#">works.</a> In <i>ICML</i> , Proceedings of Machine Learning	798
745	<a href="#">approach.</a> <i>CoRR</i> , <a href="#">abs/1907.11692.</a>	Research.	799
746	Kyle Lo, Lucy Lu Wang, Mark Neumann, Rod-	Brian Thompson, Jeremy Gwinnup, Huda Khayrallah,	800
747	ney Michael Kinney, and Daniel S. Weld. 2020.	Kevin Duh, and Philipp Koehn. 2019. <a href="#">Overcoming</a>	801
748	<a href="#">S2orc: The semantic scholar open research corpus.</a>	<a href="#">catastrophic forgetting during domain adaptation of</a>	802
749	In <i>ACL.</i>	<a href="#">neural machine translation.</a> In <i>NAACL-HLT.</i>	803
750	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh	Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris,	804
751	Hajishirzi. 2018. <a href="#">Multi-task identification of entities,</a>	Alessandro Sordoni, Philip Bachman, and Kaheer	805
752	<a href="#">relations, and coreference for scientific knowledge</a>	Suleman. 2017. <a href="#">Newsqa: A machine comprehension</a>	806
753	<a href="#">graph construction.</a> In <i>EMNLP.</i>	<a href="#">dataset.</a> In <i>Rep4NLP@ACL.</i>	807
754	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	808
755	Huang, A. Ng, and Christopher Potts. 2011. <a href="#">Learning</a>	Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz	809
756	<a href="#">word vectors for sentiment analysis.</a> In <i>ACL.</i>	Kaiser, and Illia Polosukhin. 2017. <a href="#">Attention is all</a>	810
757	Julian McAuley, Christopher Targett, Qinfeng Shi, and	<a href="#">you need.</a> In <i>NIPS.</i>	811
758	Anton van den Hengel. 2015. <a href="#">Image-based recom-</a>	Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu,	812
759	<a href="#">mendations on styles and substitutes.</a> <i>SIGIR.</i>	Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael	813
760	Michael McCloskey and Neil J. Cohen. 1989. <a href="#">Catas-</a>	Zeng. 2022. <a href="#">Training data is more valuable than you</a>	814
761	<a href="#">trophic interference in connectionist networks: The</a>	<a href="#">think: A simple and effective method by retrieving</a>	815
762	<a href="#">sequential learning problem.</a> <i>The Psychology of</i>	<a href="#">from training data.</a> In <i>ACL.</i>	816
763	<a href="#">Learning and Motivation.</a>	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	817
764	Mark Neumann, Daniel King, Iz Beltagy, and Waleed	Chaumond, Clement Delangue, Anthony Moi, Pier-	818
765	Ammar. <a href="#">Scispacy: Fast and robust models</a>	ric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz,	819
766	<a href="#">for biomedical natural language processing.</a> In	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	820
767	<i>BioNLP@ACL.</i>		

- Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *EMNLP*.
- Yuhuai Wu, Markus N. Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. [Memorizing transformers](#). *CoRR*, abs/2203.08913.
- Wanying Xie, Yang Feng, Shuhao Gu, and Dong Yu. 2021. [Importance-based neuron allocation for multi-lingual neural machine translation](#). In *ACL*.
- Yan Zeng, Xinsong Zhang, and Hang Li. 2021. [Multi-grained vision language pre-training: Aligning texts with visual concepts](#). In *ICML*.
- Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015. [Character-level convolutional networks for text classification](#). *ArXiv*, abs/1509.01626.

## A Dataset Descriptions and Statistics

This section describes the details and statistics of three tasks: domain classification, domain extractive question answering (QA), and named entity recognition (NER).

**For the text classification tasks**, we leverage the following datasets covering four domains, including biomedical science, computer science, news, and reviews. In the biomedical domain, CHEMPROT ([Kringelum et al., 2016](#)) is the relation classification dataset based on chemical-protein interaction. RCT ([Dernoncourt and Lee, 2017](#)) is the role classification task constructed from the abstract of the biomedical articles. In the computer science domain, ACL-ARC ([Jurgens et al., 2018](#)) is the task of annotated citations for articles’ functions. SCIERC ([Luan et al., 2018](#)) is constructed from scientific abstracts annotated with relation. In the news domain, HYPERPARTISAN ([Kiesel et al., 2019](#)) is the news text classification for determining partisan leanings. AGNEWS ([Zhang et al., 2015](#)) is the topic classification for news. In the reviews domain, AMAZON ([McAuley et al., 2015](#)) is a binary classification task consisting of feedback on products. IMDB ([Maas et al., 2011](#)) consists of movies reviews, which is a binary sentiment classification dataset.

**For the NER tasks**, we use two datasets involving the news and biomedical domain. Concretely, CoNLL-2003 ([Sang and Meulder, 2003](#)) consists of news stories from the Reuters Corpus. NCBI-Disease ([Dogan et al., 2014](#)) is annotated with disease mentions.

**For the QA tasks**, we utilize two domain-specific datasets. Specifically, NewsQA ([Trischler et al., 2017](#)) is a machine comprehension dataset consisting of news articles. Medication ([Pampari et al., 2018](#)) is constructed by electronic medical records from clinical text.

The detailed statistics of text classification tasks and NER tasks are shown in Table 7, QA in Table 8.

## B Implementation Details

We use the huggingface ([Wolf et al., 2020](#)) library to implement our MAP framework, which contains various transformer-based pre-trained language models (PLMs) and their saved checkpoints. We implement the DAPT, TAPT and DAPT+TAPT models of biomedical, cs, news and reviews do-

mains<sup>3</sup> from the library released by (Gururangan et al., 2020). All the experiments are implemented on Nvidia V100 GPUs with 32GB memory. We select the best checkpoint on the validation set during training to infer the test set.

**Configurations for Classification Tasks** In this section, we explain the setting of fine-tuning for domain-specific classification tasks. We fine-tune the domain-specific PLM with our MAP framework for 5 to 15 epochs, respectively, with the same learning rate of 4e-5 and the dropout rate of 0.5. The default classification layer of the model is 1 except for the IMDB dataset with 2, and the default maximum sequence length is 256 except for IMDB with 512. We leverage the Adam optimizer to schedule the learning rate, with the Adam epsilon of 1e-8, the Adam beta-1 of 0.9, and the Adam beta-2 of 0.999. We apply the grid-search method to find the optimal batch size and numbers of GPUs for all the classification datasets. The detailed settings are shown in Table 7.

**Configurations for QA and NER** For the extractive QA tasks, we fine-tune the domain-specific PLM with our MAP framework for 3 epochs, which can converge to optimal performance. Besides, we train the model with the maximum sequence length of 384 and the learning rate of 3e-5, the weight decay rate with 1e-2 and the warm-up rate of 6e-2. For the experiments on NER tasks, we fine-tune NCBI-Disease for 20 epochs and CoNNL-2003 for 15 epochs, with the same maximum sequence length of 128, the learning rate of 5e-5, and the weight decay rate and warm-up rate are set to 0. Different from domain classification tasks, we utilize AdamW as learning rate optimizer instead of Adam. In addition, we also adopt the grid-search method to find the optimal batch size and number of GPUs for all tasks. The detailed settings are shown in Table 10.

**Configurations for small-scale Pre-training** This part describes the experimental settings of adaptive pre-training with MAP framework. Our simple pre-training experiment needs some external domain-relative corpora of two domains: Biomedical and Computer Science. Following by (Gururangan et al., 2020), we adopt SCISPACY (Neumann et al.) as a sentence splitting tool to obtain abstract and body paragraphs from S2ORC (Lo et al.,

Hyperparameters	Domain corpus
Training epochs	5 to 10
Batch size per GPU	32
Number of GPUs	4
Maximum Sequence Length	512
The number of text lines (pre-training)	35K
The number of text lines (inference)	15K
Learning Rate	4e-5
Learning Rate Optimizer	Adam
Adam epsilon	1e-8
Adam beta 1	0.9
Adam beta 2	0.999

Table 6: Hyperparameters for simple adaptive pre-training with MAP framework on biomedical and cs corpora.

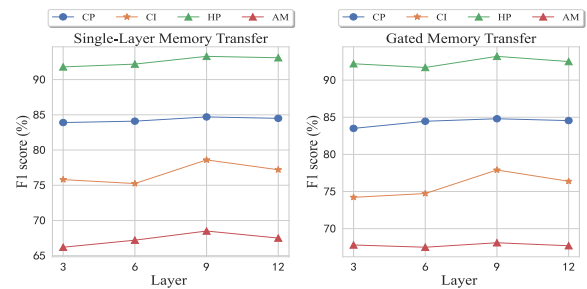


Figure 6: Performance of different layer-selection indexes of memory-attention for single-layer memory transfer and gated memory transfer strategies.

2020). After pre-processing for the corpora, we randomly sample 50K data for each domain and split 70% of them as pre-training sets and 30% as test sets. For the general corpus similar to ROBERTA’s pre-training corpus, we also randomly sample 50K data from BOOKCORPUS<sup>4</sup> and split them by using the method mentioned above. The detailed hyperparameter settings for this cross-domain adaptive pre-training are shown in Table 6.

## C Layer Selection Experiment

For single-layer memory transfer and gated memory transfer strategies, we experiment with adding the memory-attention to layer 3,6,9 and 12 in a 12-layer RoBERTa-base model, with result shown in Figure 6. We empirically find that adding memory-attention to the 9-th layer of the domain-specific model as the memory-augmented layer will obtain the best results for the two strategies. However, adding memory-attention to either too upper or too lower obtained fewer gains. Therefore, we adopt memory-attention on the 9-th layer as the default choice for the two strategies in the main experiment shown in Table 1.

<sup>3</sup><https://huggingface.co/allenai>

<sup>4</sup><https://github.com/soskek/bookcorpu>



Task	Domain	Dataset	Train	Dev.	Test	Classes
Classification	BIOMED	CHEMPROT	4169	2427	3469	13
	BIOMED	RCT <sup>†</sup>	180040	30,212	30135	5
	CS	ACL-ARC	1,688	114	139	6
	CS	SCIERC	3219	455	974	7
	NEWS	HYPERPARTISAN	515	65	65	2
	NEWS	AGNews <sup>†</sup>	115000	5000	7600	4
	REVIEWS	Amazon <sup>†</sup>	115251	5000	25000	2
	REVIEWS	IMDB <sup>†</sup>	20,000	5000	25000	2
Named Entity Recognition	News	CoNLL-2003	14,041	3,250	3,453	-
	Biomed	NCBI-Disease	5,433	924	941	-

Table 7: Statistics of Classification and NER tasks involving four domains including Biomedical, Computer Science, News and Reviews. <sup>†</sup> indicates high-resource settings.

Domain	Dataset	Train		Dev.		Test	
		Context	Question	Context	Question	Context	Question
Reviews	NewsQA	11428	74160	-	-	106	674
Biomed	Medication	182	7518	26	1858	53	4005

Table 8: Statistics of QA tasks including News and Biomedical domains. We report the number of contexts and questions of the two datasets.

	BIOMED		CS		NEWS		REVIEWS	
Hyperparameters	CP	RCT	CI	SE	HP	AG	AM	IMDB
Training Epochs	14	4	15	15	12	5	6	15
Batch Size per GPU	24	16	32	32	32	16	16	16
Number of GPUs	3	4	1	1	1	4	4	4
Maximum Sequence Length	256	256	256	256	256	384	512	512
Learning Rate					4e-5			
Dropout					0.5			
Classification Layer	1	1	1	1	1	2	1	2
Learning Rate Optimizer					Adam			
Adam Epsilon					1e-8			
Adam Beta 1					0.9			
Adam Beta 2					0.999			

Table 9: Hyperparameters for fine-tuning on eight classification tasks of four domains, we use these hyperparameters for reporting the performances of our proposed MAP framework in the main papers.

Hyperparameters	BIOMED		NEWS	
	Medication	NCBI-Disease	NewsQA	CoNNL-2003
Training epochs	3	20	3	15
Batch size per GPU	16	32	16	32
Number of GPUs	4	1	4	1
Maximum Sequence Length	384	128	384	128
Classification layer			1	
Learning rate optimizer			AdamW	
Learning rate	3e-5	5e-5	3e-5	5e-5
Weight Decay	1e-2	0	1e-2	0
LR decay Warm-up rate	6e-2	0	6e-2	0

Table 10: Hyperparameters for fine-tuning on QA and NER tasks of biomedical and news domains, we use them for reporting the performances of our proposed MAP framework in the main papers.