

LPC: A Logits and Parameter Calibration Framework on Continual Learning

Anonymous ACL submission

Abstract

Deep learning based pre-trained natural language processing (NLP) models typically pre-train on large unlabeled corpora first, then fine-tune on new tasks. When we execute such a paradigm on continuously sequential tasks, the model will suffer from the catastrophic forgetting problem (i.e., they forget the parameters learned in previous tasks when we train the model on newly emerged tasks). Inspired by the idea of how humans learn things, we aim to maintain the old knowledge when we transfer to novel contents and calibrate the old and new knowledge. We propose a Logits and Parameter Calibration (LPC) framework to reduce the catastrophic forgetting in the continual learning process. The proposed framework includes two important components, the Logits Calibration (LC) and Parameter Calibration (PC). The core idea is to reduce the difference between old knowledge and new knowledge by doing calibration on logits and parameters so that the model can maintain old knowledge while learning new tasks without preserving data in previous tasks. First, we preserve the parameters learned from the base tasks. Second, we train the existing model on novel tasks and estimate the difference between base logits and parameters and novel logits and parameters. Third, we drift from the base tasks to novel tasks gradually. Furthermore, we integrate the logits and parameter calibration into a brand-new optimization algorithm. Finally, we do experiments on 7 scenarios of the GLUE (the General Language Understanding Evaluation) benchmark. The experimental results show that our model achieves state-of-the-art performance on all 7 scenarios.

1 Introduction

Predicting labels for a large number of instances occurring continuously is a crucial problem in many real-world applications like online tweets/news summary, online product classification

in e-commerce systems, and online dialogue learning systems. In these scenarios, we not only require the model to learn from its own experiences, but also expect the model to be capable of continuously acquiring, fine-tuning, and transferring knowledge over time (Parisi et al., 2019), which is also known as continual learning. One of the most essential existing challenges we aim to solve in the continual learning is the catastrophic forgetting problem (McCloskey and Cohen, 1989; Kirkpatrick et al., 2017a). The forgetting typically happens when we apply the pre-trained model (e.g., BERT (Devlin et al., 2018)) on newly emerged tasks, the model usually forgets the parameters it learned from previous tasks when we train it on new incoming tasks.

Existing works trying to solve the catastrophic forgetting problem are varied, which can be divided into two categories: (1) storing exemplars from previous classes (Rebuffi et al., 2017; Rolnick et al., 2019); (2) regularizing the parameters when we fine-tune the model on new tasks (Kirkpatrick et al., 2017b; Li and Hoiem, 2017; Aljundi et al., 2018). Such methods aim to transfer or store the knowledge of previous tasks to the newly emerged tasks and preserve the knowledge learned previously. Memory Aware Synapses (MAS) (Aljundi et al., 2018) is an advanced approach by computing the importance of the neural network parameters in an unsupervised and online manner. MAS assigns more weights on the parameters that are most important to the model and allows the model to selectively forgets those weights that are not so essential. Also, Lee et al. (Lee et al., 2020) successfully reduce the catastrophic forgetting during the fine-tuning step by randomly mixing pre-trained parameters into a downstream model in a dropout-style.

The methods mentioned above address the catastrophic forgetting through multi-task learning. They typically require storing the data from old or pre-trained tasks, and replay them during the

084 fine-tuning. However, this learning pattern does
085 not consider the constraint of memory resource or
086 privacy issues, e.g., the data of old tasks is often
087 inaccessible or too large for the continual adapta-
088 tion setting. Unlike the multi-task learning strategy,
089 here we focus on the calibration of knowledge gap
090 between different tasks, which can reduce the catas-
091 trophic forgetting without any old data/task replay.
092 The proposed calibration framework is used for
093 both encoder parameters and output classifiers, we
094 first train the novel model and evaluate the base
095 model on novel tasks. Specifically, we add the log-
096 its calibration that can amplify the softmax output
097 of the base model, and overcome the bias towards
098 the novel category (Zhao et al., 2020), which can
099 simultaneously enforce the model to preserve pre-
100 vious knowledge via explicit weight constraints.
101 Also, for the calibration on encoder parameters, we
102 encourage the model to maintain previously learned
103 knowledge by simulating the training objective us-
104 ing the parameters of the base model. Then, during
105 the training on novel tasks, the model will calibrate
106 parameters with target drift from the base tasks to
107 the novel tasks to balance new task learning and
108 old knowledge maintenance. It allows the model
109 to focus on novel tasks by making the learning ob-
110 jective drifting from the base tasks to novel tasks
111 gradually.

112 Accordingly, we propose a Logits and Parame-
113 ter Calibration framework LPC (shown in Figure
114 1) on continual learning scenario, which is used
115 to reduce catastrophic forgetting without further
116 data storage. The proposed calibration mechanism
117 includes two components for both model encoder
118 parameters and output logits, we finally integrate
119 these two calibrations into a brand-new optimiza-
120 tion algorithm by decoupling them from the gradi-
121 ent updates in Adam optimizer. We do experiments
122 on the GLUE benchmark with pre-trained mod-
123 els BERT-base and ALBERT-xxlarge and achieve
124 state-of-the-art performance.

125 The contributions of our work are three folds.
126 First, we propose LPC, a novel continual learning
127 framework, which can reduce catastrophic forget-
128 ting effectively without storing previous instances.
129 Second, we develop a new mechanism by calibrat-
130 ing the logits and parameters with target drift from
131 base tasks to novel tasks, thereby alleviating the
132 catastrophic forgetting during the model updating.
133 Third, combining with a parameter regularization
134 based approach, our model achieves state-of-the-

art performance while addressing the old knowl-
edge forgetting without data storage. Therefore,
the newly proposed LPC is feasible for researchers
to use for further explorations in this field.

2 Related Works 139

2.1 Continual Learning 140

Continual learning is also named as life-long learn-
ing, sequential learning, or incremental learning.
As the name suggests, continual learning aims to
learn tasks in a sequential way. In the field of
biology, biological neural networks exhibit con-
tinual learning in which they acquire new knowl-
edge over a lifetime (Zenke et al., 2017). How-
ever, continual learning in deep neural networks
suffers from a phenomenon called *catastrophic for-
getting* (Shin et al., 2017). Thus, one of the most
essential goals of continual learning systems is to
achieve satisfying performance on all tasks in an
incremental way. Reducing catastrophic forgetting
plays a vital role to achieve it. Current continual
learning approaches can be classified into the fol-
lowing three families (De Lange et al., 2019): (1)
Replay methods, (2) Regularization-based meth-
ods, and (3) Parameter isolation methods. Replay
methods store samples in a raw format or generate
pseudo-samples. Regularization-based methods
eschews storing raw inputs, prioritizing privacy,
and alleviating memory requirements. Parameter
isolation methods dedicate different model parame-
ters to each task to prevent any possible forgetting.
Our method is an advanced regularization-based
method alleviating the catastrophic forgetting with-
out data storage.

2.2 Fine-tuning 168

Fine-tuning is a successful method in transfer
learning by the following four steps: (1) pre-
train a source neural network model on the source
datasets; (2) create a new neural network model
which copies all model designs and their parame-
ters on the source model except the output layer;
(3) add an output layer to the target model; (4) train
the target model on the target datasets. Girshick
et al. (Girshick et al., 2014) propose a R-CNN to
fine-tune all network parameters. Long et al. (Long
et al., 2015) propose DAN only fine-tuning the pa-
rameters of the last few layers. Li et al. (Li et al.,
2018) investigate several regularization schemes
that explicitly promote the similarity of the fine-
tuned model with the original pre-trained model.

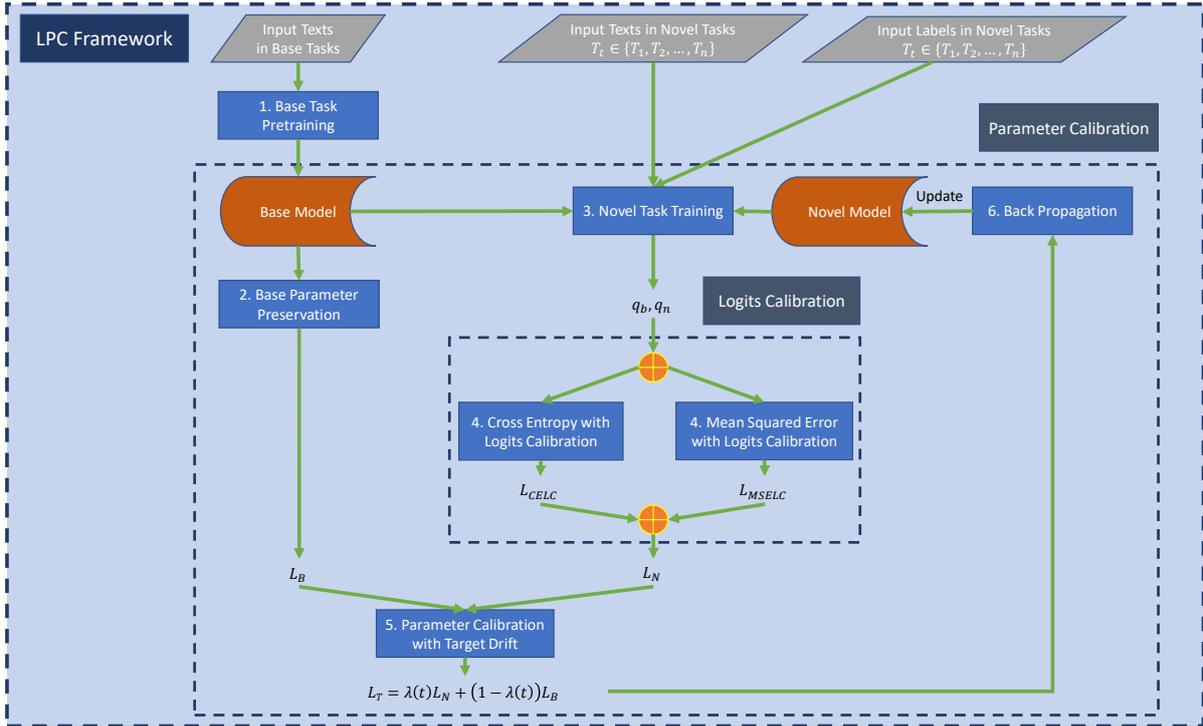


Figure 1: The Overview of LPC Framework. (1) Given the large-scale input texts from base tasks, we first pre-train a base model on these texts and initialize our model (for novel tasks) same as the base one. (2) We do base parameter preservation to preserve the parameters of the pre-trained model, and estimate the difference between base parameters and novel parameters during training, then compute the loss for the base model L_B . (3) During the novel task training, we compute logits q_b and q_n for the base model and the novel model, respectively. (4) We do logits calibration (e.g., cross entropy with logits calibration for classification tasks) given q_b and q_n using L_{CELC} or L_{MSELC} for regression tasks as the loss for the novel model L_N . (5) In the parameter calibration with target drift, the objective function drifts from L_B to L_N gradually with the annealing coefficient $\lambda(t)$. (6) Finally, we perform back propagation to update the parameters of the novel model with parameter calibration.

184 Guo et al. (Guo et al., 2019) propose an adaptive
 185 fine-tuning approach SpotTune which automatic-
 186 ally decides the optimal set of layers to fine-tune
 187 in a pre-trained model on a new task. Our method
 188 is a trade-off between multi-task learning and fine-
 189 tuning.

190 3 Proposed Approach

191 In this section, we introduce our proposed Logits
 192 and Parameter Calibration framework, LPC. The
 193 LPC framework includes two essential parts: (1)
 194 Logits Calibration (LC) that execute calibration
 195 on the logits to reduce the logits forgetting and
 196 increase the accuracy and (2) Parameter Calibra-
 197 tion (PC) to do calibration on the parameters to
 198 reduce the parameter forgetting. For the Logits
 199 Calibration, we apply the Cross Entropy with Log-
 200 its Calibration (CELC) for classification tasks (or
 201 the Means Squared Error with Logits Calibration
 202 (MSELC) for regression tasks). The Parameter Cal-
 203 ibration consists of three components: (1) Base

204 Parameter Preservation (BPP) that tries to preserve
 205 the parameters we learned from base tasks, (2)
 206 Novel Task Training (NTT) that trains the previous
 207 model on novel tasks, and (3) Parameter Calibra-
 208 tion with Target Drift (PCTD) that focuses on drift-
 209 ing from the base tasks to novel tasks gradually.
 210 What is more, we introduce LPC algorithm by inte-
 211 grating the Logits Calibration (CELC or MSELC)
 212 and all three parts of Parameter Calibration (BPP,
 213 NTT, and PCTD) into a brand-new optimization
 214 algorithm based on the well-known Adam (Kingma
 215 and Ba, 2014) optimization algorithm.

216 3.1 Logits Calibration

217 In this section, we introduce our proposed logits
 218 calibration, the Cross Entropy with Logits Calibra-
 219 tion (CELC) for classification tasks. Some other
 220 loss function (e.g., the Mean Squared Error for re-
 221 gression) can also be combined with the Logits
 222 Calibration, in the following paragraph.

Algorithm 1 LPC

```
1: given initial learning rate  $\alpha \in \mathbb{R}$ , momentum factors  $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ , pre-trained
   parameter vector  $\theta^* \in \mathbb{R}^n$ , hyperparameter for the regularizer  $\delta \in \mathbb{R}$ , coefficient of the quadratic
   penalty  $\gamma \in \mathbb{R}$ , hyperparameter controlling the annealing rate  $r \in \mathbb{R}$ , hyperparameter controlling the
   timesteps  $t_0 \in \mathbb{N}$ .
2: initialize timestep  $t \leftarrow 0$ , parameter vector  $\theta_{t=0} \in \mathbb{R}^n$ , importance weights  $\Omega \leftarrow \mathbf{1}$ , first moment
   vector  $m_{t=0} \leftarrow 0$ , second moment vector  $v_{t=0} \leftarrow 0$ , schedule multiplier  $\eta_{t=0} \in \mathbb{R}$ .
3: repeat
4:    $t \leftarrow t + 1$  ▷ update timestep
5:    $x \leftarrow \text{SelectBatch}(\mathbf{x})$  ▷ select batch data
6:    $q_{n,t} \leftarrow Q_{n,t}(x, \theta_{t-1})$  ▷ compute output logits for the novel model
7:    $q_{b,t} \leftarrow Q_{b,t}(x, \theta^*)$  ▷ compute output logits for the base model
8:    $\nabla(f_t(x; \theta_{t-1})) \leftarrow \nabla(L_{CELC}(q_{n,t}, q_{b,t}) \parallel L_{MSELC}(q_{n,t}, q_{b,t}))$  ▷ compute gradients
9:    $\Omega_t \leftarrow \Omega_{t-1}$ 
10:  for  $k \leftarrow 0$  to  $N$  do
11:     $g_t(x_k) \leftarrow \nabla l_2^2(f_t(x_k; \theta_{t-1}))$ 
12:     $\Omega_t \leftarrow \Omega_t + \|g_t(x_k)\|$ 
13:  end for
14:   $\Omega_t \leftarrow \Omega_t / N$  ▷ compute importance weights after each update epochs
15:   $\lambda(t) \leftarrow 1 / (1 + \exp(-r \cdot (t - t_0)))$  ▷ compute annealing coefficient
16:   $g_t \leftarrow \lambda(t) \nabla f_t(x; \theta_{t-1}) + 2(1 - \lambda(t))\delta\gamma\Omega_t(\theta_{t-1} - \theta^*)$  ▷ compute new gradients
17:   $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1)g_t$  ▷ update biased first moment estimate
18:   $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2)g_t^2$  ▷ update biased second raw moment estimate
19:   $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$  ▷ compute bias-corrected first moment estimate
20:   $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$  ▷ compute bias-corrected second raw moment estimate
21:   $\eta_t \leftarrow \text{SetScheduleMultiplier}(t)$  ▷ can be fixed, decay, or also be used for warm restarts
22:   $\theta_t \leftarrow \theta_{t-1} - \eta_t (\lambda(t) \alpha \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon) + 2(1 - \lambda(t))\delta\gamma\Omega_t(\theta_{t-1} - \theta^*))$  ▷ update parameters
23: until stopping criterion is met
24: return optimized parameters  $\theta_t$ 
```

3.1.1 Cross Entropy with Logits Calibration

The Cross Entropy (CE) Loss (Zhang and Sabuncu, 2018) is a widely-used loss for classification tasks in deep learning. It first applies a log softmax function on the output logits of the neural network. Then, it computes the negative log likelihood (nll) loss on the output of the log softmax function. Typically, the cross entropy loss can be defined as follows:

$$L_{CE}(q) = - \sum_{i=1}^{N_C} p_i \log \left(\frac{\exp(q_{n,i})}{\sum_{j=1}^{N_C} \exp(q_{n,j})} \right) \quad (1)$$

where N_C is the total number of classes in the novel tasks. $q_{n,i}$ represents the output logits for class i of the novel model on the novel tasks. p_i

can be considered as the binary label of class i . If the data input x belongs to class i , the value of p_i will be 1, otherwise, the value will be 0.

Nevertheless, the original cross entropy loss only concerns the performance of the novel model. Thus, the model will suffer the catastrophic forgetting problem with the step increasing. In order to reduce the catastrophic forgetting problem, we consider to simultaneously evaluate the base model on the novel tasks and compute the output logits of the base model q_b .

Inspired by the idea from (Kukleva et al., 2021) which revises the original cross entropy loss by adding the summation of the exponential logits of the base classes classifier to the denominator to change the normalization scale, we add the logits information of the base model into the cross entropy loss.

254 However, different from Kukleva’s method, we
 255 do logtis calibration by adding the difference be-
 256 tween each logtis of the novel model and the base
 257 model ($q_{n,i} - q_{b,i}$) to the corresponding output log-
 258 its $q_{n,i}$ of the novel model for class i . In this way,
 259 the model can preserve important output logtis in-
 260 formation of each class for the base model in an
 261 element-wise way. Our proposed Cross Entropy
 262 with Logtis Calibration (CELC) Loss is shown in
 263 Equation 2:

$$L = - \sum_{i=1}^{N_C} p_i \log \left(\frac{\exp(q_{n,i} + \mu(q_{n,i} - q_{b,i}))}{\sum_{j=1}^{N_C} \exp(q_{n,j} + \mu(q_{n,j} - q_{b,j}))} \right) \quad (2)$$

264 where we multiply the difference between the
 265 logtis for the novel model and the base model
 266 ($q_n - q_b$) by a weight item $\mu \in [0, 1]$ to control
 267 the calibration degree. By employing this new loss
 268 function, we can also increase the accuracy of the
 269 model through training process by giving a reward
 270 to the logtis for the correct class if $q_{n,i}$ is larger
 271 than $q_{b,i}$, otherwise, giving a penalty to the logtis
 272 for the correct class if $q_{n,i}$ is smaller than $q_{b,i}$.
 273

274 3.1.2 Mean Squared Error with Logtis 275 Calibration

276 Mean Squared Error (MSE) Loss (Fisher, 1922)
 277 is the most commonly-used loss function for re-
 278 gression tasks. It computes the squared L2 norm
 279 between output logtis and the true values and takes
 280 the mean of the full batch. Follow the idea of the
 281 logtis calibration on cross entropy loss, we evaluate
 282 the base model on novel tasks and take out the out-
 283 put logtis q_b . We measure the difference between
 284 the output logtis of the novel model and the base
 285 model by adding a squared L2 norm on the differ-
 286 ence between logtis of the novel model and the base
 287 model $(q_n - q_b)^2$ to the original function. The pro-
 288 posed Mean Squared Error with Logtis Calibration
 289 Loss L_{MSELC} is shown in Equation 3:

$$L_{MSELC}(q) = (q_n - p)^2 + \mu(q_n - q_b)^2 \quad (3)$$

291 3.2 Parameter Calibration

292 In this section, we introduce the second module
 293 of our model, Parameter Calibration (PC). Our
 294 proposed Parameter Calibration method can effec-
 295 tively reduce the catastrophic forgetting by giving
 296 a penalty to the prediction if the parameters of

297 the novel model are different from the base model
 298 by adding the squared difference between the pa-
 299 rameters of the novel model and the base model
 300 to the training loss. It includes three parts: (1)
 301 Base Parameter Preservation (BPP), (2) Novel Task
 302 Training (NTT), and (3) Parameter Calibration with
 303 Target Drift (PCTD).

304 3.2.1 Base Parameter Preservation

305 As shown in Figure 1, in Base Parameter Preser-
 306 vation, we try to maintain the parameters of the
 307 base model like BERT (Devlin et al., 2018). Here,
 308 we add a regularization to the posterior of param-
 309 eters given data. The Base Parameter Preservation
 310 method can be regarded as an improved method
 311 derived from EWC (Kirkpatrick et al., 2017b) and
 312 MAS (Aljundi et al., 2018). Different from EWC,
 313 BPP measures the importance of each parameter
 314 by introducing the importance weights Ω . During
 315 training, the novel model preserves the information
 316 of the most important parameters to a great extent
 317 by penalizing the changes to those important pa-
 318 rameters more severely. The detailed derivation of
 319 our proposed loss function is shown in Equation 4:

$$\begin{aligned} L_B &= -\log p(\theta|D_B) \\ &\approx -\log p(\theta^*|D_B) + \delta(\theta - \theta^*)^T H(\theta^*)\Omega(\theta)(\theta - \theta^*) \\ &\approx \delta(\theta - \theta^*)^T H(\theta^*)\Omega(\theta)(\theta - \theta^*) \\ &\approx \delta(\theta - \theta^*)^T (NF(\theta^*) + H_{prior}(\theta^*))\Omega(\theta)(\theta - \theta^*) \\ &\approx \delta N \sum_{ij} F_{ij}\Omega_{ij}(\theta_{ij} - \theta_{ij}^*)^2 \\ &\approx \delta NF \sum_{ij} \Omega_{ij}(\theta_{ij} - \theta_{ij}^*)^2 \\ &= \delta\gamma \sum_{ij} \Omega_{ij}(\theta_{ij} - \theta_{ij}^*)^2 \end{aligned} \quad (4)$$

320 where δ is a hyperparameter for the regularizer.
 321 $H(\theta^*)$ is the Hessian matrix of the optimization
 322 objective with respect to θ^* . We can approximate
 323 $H(\theta^*)$ with the empirical Fisher information ma-
 324 trix $F(\theta^*)$ (Martens, 2014). N is the total num-
 325 ber of data inputs in D_B . $H_{prior}(\theta^*)$ is the Hes-
 326 sian matrix of the negative log prior probability
 327 $-\log p(\theta)$. EWC ignores $H_{prior}(\theta^*)$ and approx-
 328 imates $H(\theta^*)$ by assigning the diagonal values of
 329 $F(\theta^*)$ to $H(\theta^*)$. Thus, we replace NF with a
 330 constant value γ at the end of the derivation. We
 331 can consider γ as a coefficient of the quadratic
 332 penalty. During the derivation, we can simply ig-
 333 nore $-\log p(\theta^*|D_B)$ as it is a constant term with
 334

respect to θ^* . $\Omega(\theta)$ is estimated by the sensitivity of the squared L2 norm of the function output to their changes. We can obtain Ω_{ij} by accumulating the gradients over the given data points by Equation 5:

$$\Omega_{ij} = \frac{1}{N} \sum_{k=1}^N \|g_{ij}(x_k)\| \quad (5)$$

where $g_{ij}(x_k) = \frac{\partial [l_2^2(f(x_k; \theta))]}{\partial \theta_{ij}}$ is the gradients of the squared L2 norm of the learned neural network with respect to the parameter θ_{ij} . The output of $f(x_k; \theta)$ is the loss of the network.

In Equation 4, θ_{ij} is the parameter of the novel model of the connections between pairs of neurons n_i and n_j in two consecutive layers. θ^* represents the pre-trained parameters, which can be assumed as a local minimum of the parameter space as shown in Equation 6:

$$\theta^* = \arg \min_{\theta} \{-\log p(\theta|D_B)\} \quad (6)$$

3.2.2 Novel Task Training

In the novel task training process, we train the novel model and evaluate the base model on novel tasks simultaneously. The function of the neural network whose output is the loss of the model can be represented as follows:

$$L_N = f_t(x; \theta_{t-1}) \quad (7)$$

where t is the timestep. We compute the loss by the proposed Cross Entropy with Logits Calibration (CELC) for classification tasks and Mean Squared Error with Logits Calibration (MSELC) for regression tasks as follows:

$$f_t = L_{CELC}(Q(x; \theta_{t-1})) \parallel L_{MSELC}(Q(x; \theta_{t-1})) \quad (8)$$

where $Q(x; \theta_{t-1})$ represents the function of the novel model and the base model whose output are logits with data inputs x and parameters θ_{t-1} of the model in timestep $t - 1$.

3.2.3 Parameter Calibration with Target Drift

Multi-task learning tries to achieve satisfying performance on both of the base tasks and novel tasks. However, one of the most essential problems of multi-task learning is that it is inconsistent with adaptation (Chen et al., 2020). To deal with this problem, we introduce Parameter Calibration with

Target Drift, a method allowing the objective function to gradually drift from L_B to L_N with the annealing coefficient $\lambda(t)$:

$$L_T = \lambda(t)L_N + (1 - \lambda(t))L_B \quad (9)$$

where t refers to the timestep during the training process. We compute $\lambda(t) = \frac{1}{1 + \exp(-r \cdot (t - t_0))}$ as the sigmoid annealing function (Kiperwasser and Ballesteros, 2018), where r is the hyperparameter controlling the annealing rate and t_0 is the hyperparameter controlling the timesteps.

When $t < t_0$, $-r \cdot (t - t_0)$ will be positive. In this case, if $r \rightarrow \infty$, then $\exp(-r \cdot (t - t_0)) \rightarrow \infty$, $\lambda(t) \rightarrow 0$, $L_T = L_B$. When $t > t_0$, $-r \cdot (t - t_0)$ will be negative. In this case, if $r \rightarrow \infty$, then $\exp(-r \cdot (t - t_0)) \rightarrow 0$, $\lambda(t) \rightarrow 1$, $L_T = L_N$. At this moment, our method can be regarded as fine-tuning. Otherwise, if $r \rightarrow 0$, then $-r \cdot (t - t_0) \rightarrow 0$, $\exp(-r \cdot (t - t_0)) \rightarrow 1$, $\lambda(t) \rightarrow 0.5$, $L_T = 0.5L_N + 0.5L_B$. In this case, our method can be regarded as multi-task learning. Finally, if $0 < r < \infty$, then $0 < \lambda < 1$. In this case, our method can be regarded as a trade off between fine-tuning and multi-task learning. With time goes by, the objective of the model drifts from base tasks to novel tasks gradually. Finally, by doing back propagation, we update parameters of the novel model with parameter calibration.

3.3 LPC Algorithm

In this section, we combine the Logits Calibration (CELC) with all three parts of Parameter Calibration (BPP, NTT, and PCTD) into a brand-new optimization algorithm as shown in Algorithm 1. The Logits Calibration (LC) part is shown from line 6 to line 8. The Parameter Calibration (PC) part is shown from line 9 to line 18 and line 24. Here, we introduce LPC Algorithm which integrates the quadratic penalty with importance weights and the annealing coefficient into a complete optimization algorithm by decoupling them from the gradient update in Adam optimization algorithm (Kingma and Ba, 2014). The orange part in Algorithm 1 depicts how LPC is different from Adam.

From line 9 to line 16, we show how we calculate Ω by initializing Ω as a tensor filled with the scalar value one. The size of Ω are the same as that of parameter size of the base model and the novel model. From line 11 to line 14, we accumulate the gradients of the squared L2 norm of the learned neural network over the given data inputs to obtain

Table 1: Experimental Results. All of results are the medians over 5 runs. The metric for CoLA is mcc (Matthew Correlation Coefficient). The metric for STS-B is corr (Average of Pearson and Spearman Correlation Coefficient). All other metrics are acc (Accuracy).

Model	CoLA	MRPC	QNLI	RTE	SST-2	STS-B	WNLI	Avg	Avg	Avg
	mcc	acc	acc	acc	acc	corr	acc	acc	mcc	corr
BERT-base + PALs (Stickland and Murray, 2019)	51.2	84.6	90.0	76.0	92.6	85.8	N/A	N/A	51.2	85.8
BERT-large + Adapters (Houlsby et al., 2019)	59.5	89.5	90.7	71.5	94.0	86.9	N/A	N/A	59.5	86.9
BERT-large + Diff pruning (Guo et al., 2020)	60.5	87.0	92.9	68.1	93.8	83.5	N/A	N/A	60.5	83.5
BERT-base + Adam (rerun) <i>Median</i>	57.1	81.3	91.0	63.9	93.1	89.2	56.3	77.1	57.1	89.2
BERT-base + RecAdam (rerun) <i>Median</i>	59.9	85.7	91.4	70.8	93.1	90.0	56.3	79.5	59.9	90.0
BERT-base + LPC <i>Median</i>	61.8	86.1	91.5	74.7	93.2	90.3	62.0	81.5	61.8	90.3
ALBERT-xxlarge + Adam (rerun) <i>Median</i>	70.5	88.8	93.7	72.9	91.1	92.2	69.0	83.1	70.5	92.2
ALBERT-xxlarge + RecAdam (rerun) <i>Median</i>	70.5	87.5	93.9	89.5	93.9	92.8	78.9	88.7	70.5	92.8
ALBERT-xxlarge + LPC <i>Median</i>	74.1	89.4	94.3	89.5	95.8	93.3	81.7	90.1	74.1	93.3

importance weights Ω_{ij} for parameter θ_{ij} . In line 15, we compute the mean value of Ω_{ij} by dividing it by N . Here, N is the total number of data inputs at a given phase. In line 18, we compute the gradients of the loss function as a weighted combination of the gradients of L_N and L_B . In line 24, we update the network parameters θ by the gradient descent method.

4 Evaluations

In this section, we evaluate LPC on the General Language Understanding Evaluation (GLUE) (Wang et al., 2018) benchmark. We compare our model with PALs (Stickland and Murray, 2019), Adapters (Houlsby et al., 2019), Diff Pruning (Guo et al., 2020), Adam (Kingma and Ba, 2014), and RecAdam (Chen et al., 2020).

4.1 Experimental Setup

We perform the experiments based on deep pre-trained language models BERT-base¹ (Devlin et al., 2018) and ALBERT-xxlarge (Lan et al., 2019), respectively. BERT aims to learn a Transformer encoder for representing texts. BERT is a multi-layer bidirectional Transformer encoder. In BERT, the Transformer uses bidirectional self-attention. ALBERT is an advanced deep pre-trained language model with lower memory consumption and faster training speed than BERT. ALBERT improves BERT using parameter reduction techniques and employing self-supervised loss for sentence-order prediction (SOP).

We evaluate our approach LPC on the GLUE benchmark. GLUE benchmark is a collection of resources for training, evaluating, and analyzing natu-

ral language understanding (NLU) systems (Wang et al., 2018). It contains the following 9 different scenarios: (1) Single-Sentence Scenarios: **CoLA** The Corpus of Linguistic Acceptability (Warstadt et al., 2019), and **SST-2** The Stanford Sentiment Treebank (Socher et al., 2013); (2) Similarity and Paraphrase Scenarios: **MRPC** The Microsoft Research Paraphrase Corpus (Dolan and Brockett, 2005), **QQP** The Quora Question Pairs dataset², and **STS-B** The Semantic Textual Similarity Benchmark (Cer et al., 2017); (3) Inference Scenarios: **MNLI** The Multi-Genre Natural Language Inference Corpus (Williams et al., 2017), **QNLI** The Stanford Question Answering Dataset (Rajpurkar et al., 2016), **RTE** The Recognizing Textual Entailment datasets (Dagan et al., 2005) (Haim et al., 2006) (Giampiccolo et al., 2007) (Bentivogli et al., 2009), and **WNLI** The Winograd Schema Challenge (Levesque et al., 2012).

4.2 Results

We perform experiments on 7 scenarios of the GLUE benchmark as shown in Table 1. From the experimental results with BERT-base model, we outperform BERT-base with Adam and BERT-base with RecAdam (Chen et al., 2020) models on 7 out of 7 scenarios of the GLUE benchmark and achieve 2.5% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 3.2% improvements measured by mcc on CoLA, and 0.3% improvements measured by corr on STS-B compared with RecAdam. Especially, we achieve significant improvements on WNLI corpus (+10.1%) and CoLA corpus (+5.5%). From the experimen-

¹https://huggingface.co/transformers/model_doc/bert.html

²<https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs>

Table 2: The Results of Ablation Study on Adam, Logits Calibration (LC), Parameter Calibration (PC), and Logits and Parameter Calibration (LPC). All of results are the medians over 5 runs. The metric for CoLA is mcc (Matthew Correlation Coefficient). The metric for STS-B is corr (Average of Pearson and Spearman Correlation Coefficient). All other metrics are acc (Accuracy).

Model	CoLA	MRPC	QNLI	RTE	SST-2	STS-B	WNLI	Avg	Avg	Avg
	mcc	acc	acc	acc	acc	corr	acc	acc	mcc	corr
BERT-base + Adam (rerun) <i>Median</i>	8.5k	3.7k	105k	2.5k	67k	7k	634	77.1	57.1	89.2
BERT-base + LC <i>Median</i>	57.1	81.3	91.0	63.9	93.1	89.2	56.3	77.1	57.1	89.2
BERT-base + LC <i>Median</i>	61.2	82.8	91.5	66.8	92.3	89.3	56.3	77.9	61.2	89.3
BERT-base + PC <i>Median</i>	61.4	85.3	91.5	72.2	92.8	90.2	57.7	79.9	61.4	90.2
BERT-base + LPC <i>Median</i>	61.8	86.1	91.5	74.7	93.2	90.3	62.0	81.5	61.8	90.3

tal results on ALBERT-xxlarge model, we outperform ALBERT-xxlarge with Adam and ALBERT-xxlarge with RecAdam models on 7 out of 7 scenarios of the GLUE benchmark and achieve 1.6% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 5.1% improvements measured by mcc on CoLA, and 0.5% improvements measured by corr on STS-B compared with RecAdam. Specifically, we achieve great improvements on CoLA corpus (+5.1%) and WNLI corpus (+3.5%). What is more, there is no obvious relationship between the size of the datasets and the results. Namely, our model performs well on both large datasets and small datasets.

4.3 Ablation Study

As we have mentioned, our model (LPC) has two important components, Logits Calibration (LC) and Parameter Calibration (PC). Thus, we do ablation study on these two components separately with BERT-base pre-trained model on the 7 scenarios of the GLUE benchmark. The results of ablation study is shown in Table 2. We can see both of LC and PC achieve better results than the baseline Adam. LPC achieves the best results among all three models. Compared with Adam, LC achieves 1.0% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 7.2% improvements measured by mcc on CoLA, and 0.1% improvements measured by corr on STS-B. Compared with Adam, PC achieves 3.6% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 7.5% improvements measured by mcc on CoLA, and 1.1% improvements measured by corr on STS-B. Compared with Adam, LPC achieves 5.7% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 8.2% improvements measured by mcc on CoLA, and 1.2% improvements measured by corr

on STS-B.

5 Conclusion

In this paper, we propose Logits and Parameter Calibration (LPC) framework on continual learning to deal with the catastrophic forgetting problem. The proposed framework includes two important components, Logits Calibration (LC) and Parameter Calibration (PC). We propose the Cross Entropy Loss with Logits Calibration (CELC) for classification tasks and the Mean Squared Error with Logits Calibration (MSELC) for regression tasks. The Parameter Calibration consists of three components: (1) Base Parameter Preservation (BPP), (2) Novel Task Training (NTT), and (3) Parameter Calibration with Target Drift (PCTD). What is more, we introduce LPC algorithm by integrating the Logits Calibration and all three parts of Parameter Calibration (BPP, NTT, and PCTD) into a brand-new optimization algorithm based on the well-known Adam optimization algorithm. We do experiments on 7 scenarios of GLUE benchmark and achieve state-of-the-art performance on all the 7 scenarios. The limitation of our work is that our work cannot handle the online learning settings. This means that when data comes in an online manner (sometimes without labels), we have no technique to handle it. Thus, our future direction is to make our model fit the online learning settings. We also release the open-source LPC Algorithm to further benefit the continual learning research community.

References

Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars. 2018. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 139–154.

565	Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. 2009. The fifth pascal recognizing textual entailment challenge. In <i>TAC</i> .	618
566		619
567		620
568	Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. 2017. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. <i>arXiv preprint arXiv:1708.00055</i> .	621
569		622
570		623
571		624
572		625
573	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. <i>arXiv preprint arXiv:2004.12651</i> .	626
574		627
575		628
576		629
577	Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The pascal recognising textual entailment challenge. In <i>Machine Learning Challenges Workshop</i> , pages 177–190. Springer.	630
578		631
579		632
580		633
581	Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Gregory Slabaugh, and Tinne Tuytelaars. 2019. A continual learning survey: Defying forgetting in classification tasks. <i>arXiv preprint arXiv:1909.08383</i> .	634
582		635
583		636
584		637
585		638
586	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .	639
587		640
588		641
589		642
590	William B Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In <i>Proceedings of the Third International Workshop on Paraphrasing (IWP2005)</i> .	643
591		644
592		645
593		646
594	Ronald A Fisher. 1922. On the mathematical foundations of theoretical statistics. <i>Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character</i> , 222(594-604):309–368.	647
595		648
596		649
597		650
598		651
599	Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and William B Dolan. 2007. The third pascal recognizing textual entailment challenge. In <i>Proceedings of the ACL-PASCAL workshop on textual entailment and paraphrasing</i> , pages 1–9.	652
600		653
601		654
602		655
603		656
604	Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 580–587.	657
605		658
606		659
607		660
608		661
609	Demi Guo, Alexander M Rush, and Yoon Kim. 2020. Parameter-efficient transfer learning with diff pruning. <i>arXiv preprint arXiv:2012.07463</i> .	662
610		663
611		664
612	Yunhui Guo, Honghui Shi, Abhishek Kumar, Kristen Grauman, Tajana Rosing, and Rogerio Feris. 2019. Spottune: transfer learning through adaptive fine-tuning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 4805–4814.	665
613		666
614		667
615		668
616		669
617		670
	R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini, and Idan Szpektor. 2006. The second pascal recognising textual entailment challenge. In <i>Proceedings of the Second PASCAL Challenges Workshop on Recognising Textual Entailment</i> .	671
		672
		673
	Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In <i>International Conference on Machine Learning</i> , pages 2790–2799. PMLR.	674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700

674	Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. 2015. Learning transferable features with deep adaptation networks. In <i>International conference on machine learning</i> , pages 97–105. PMLR.	Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual learning through synaptic intelligence. In <i>International Conference on Machine Learning</i> , pages 3987–3995. PMLR.	729
675			730
676			731
677			732
678	James Martens. 2014. New insights and perspectives on the natural gradient method. <i>arXiv preprint arXiv:1412.1193</i> .	Zhilu Zhang and Mert R Sabuncu. 2018. Generalized cross entropy loss for training deep neural networks with noisy labels. In <i>32nd Conference on Neural Information Processing Systems (NeurIPS)</i> .	733
679			734
680			735
681	Michael McCloskey and Neal J Cohen. 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In <i>Psychology of learning and motivation</i> , volume 24, pages 109–165. Elsevier.		736
682			
683		Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. 2020. Maintaining discrimination and fairness in class incremental learning. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 13208–13217.	737
684			738
685			739
686	German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. 2019. Continual lifelong learning with neural networks: A review. <i>Neural Networks</i> , 113:54–71.		740
687			741
688		A LPC Appendix	742
689		A.1 Hyperparameter Analysis	743
690	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> .	In this section, we analyze the most essential hyperparameters we set in the LPC model. δ is a hyperparameter controlling the level of regularization. Setting δ between 1 and 2 balances the level of regularization. Ω is a parameter measuring the importance of different parameters in the model. Initializing Ω as ones makes the importance of each parameter more balanced. The hyperparameter u_e controls the updating epochs of Ω . Typically, u_e is between 1 and 16. w_s is a hyperparameter controlling the number of steps of updating with low learning rate before/at the beginning of the training process. We set w_s as 0, 320 or 640. After these warmup steps, we will use the regular learning rate to train our model until convergence. In other words, we have a few steps adjustment before we actually train the model. From our experiments, we find that the hyperparameters δ , u_e , and w_s have great influences on the experimental results.	744
691			745
692			746
693			747
694	Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. 2017. icarl: Incremental classifier and representation learning. In <i>CVPR</i> , pages 2001–2010.		748
695			749
696			750
697			751
698	David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P Lillicrap, and Greg Wayne. 2019. Experience replay for continual learning. <i>NeurIPS</i> .		752
699			753
700			754
701	Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. 2017. Continual learning with deep generative replay. <i>arXiv preprint arXiv:1705.08690</i> .		755
702			756
703			757
704	Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In <i>Proceedings of the 2013 conference on empirical methods in natural language processing</i> , pages 1631–1642.		758
705			759
706			760
707			761
708			762
709			763
710			764
711	Asa Cooper Stickland and Iain Murray. 2019. Bert and pals: Projected attention layers for efficient adaptation in multi-task learning. In <i>International Conference on Machine Learning</i> , pages 5986–5995. PMLR.		765
712			766
713			767
714			768
715			769
716	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. <i>arXiv preprint arXiv:1804.07461</i> .		770
717			771
718			772
719			773
720			774
721	Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2019. Neural network acceptability judgments. <i>Transactions of the Association for Computational Linguistics</i> , 7:625–641.		775
722			776
723			777
724			778
725	Adina Williams, Nikita Nangia, and Samuel R Bowman. 2017. A broad-coverage challenge corpus for sentence understanding through inference. <i>arXiv preprint arXiv:1704.05426</i> .		779
726			
727			
728			

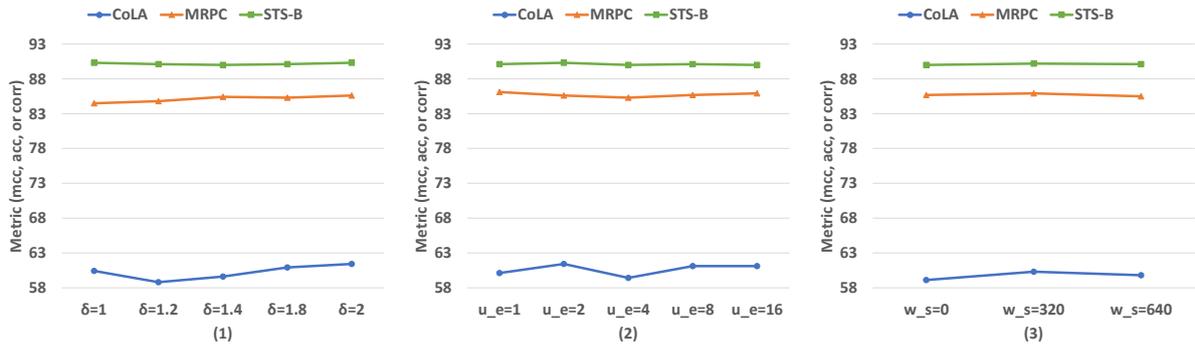


Figure 2: Comparison of Different Hyperparameter Initializations on CoLA, MRPC, and STS-B Corpora with BERT-base Pre-trained Model. The metric for CoLA, MRPC, and STS-B are mcc (Matthew Correlation Coefficient), acc (Accuracy), and corr (Average of Pearson and Spearman Correlation Coefficient), respectively.

Table 3: The Results of Ablation Study on Adam, Logits Calibration (LC), Parameter Calibration (PC), and Logits and Parameter Calibration (LPC) with ALBERT-xxlarge Pre-trained Model. All of results are the medians over 5 runs. The metric for CoLA is mcc (Matthew Correlation Coefficient). The metric for STS-B is corr (Average of Pearson and Spearman Correlation Coefficient). All other metrics are acc (Accuracy).

Model	CoLA	MRPC	QNLI	RTE	SST-2	STS-B	WNLI	Avg	Avg	Avg
	mcc	acc	acc	acc	acc	corr	acc	acc	mcc	corr
ALBERT-xxlarge + Adam (rerun) <i>Median</i>	8.5k	3.7k	105k	2.5k	67k	7k	634	82.9	70.5	92.2
ALBERT-xxlarge + LC <i>Median</i>	71.0	88.5	93.8	88.4	95.5	92.3	70.4	87.3	71.0	92.3
ALBERT-xxlarge + PC <i>Median</i>	74.1	88.6	94.0	88.4	95.7	92.9	74.6	88.3	74.1	92.9
ALBERT-xxlarge + LPC <i>Median</i>	74.1	89.4	94.3	89.5	95.8	93.3	81.7	90.1	74.1	93.3

to 8, the model performance increases again. When u_e increases from 8 to 16, there is a slight increase on the performance.

In Figure 2 chart (3), we set $\delta = 1$ and $u_e = 1$. We can see when w_s increases from 0 to 320, there is a big increase on all three corpora. However, when w_s increases from 320 to 640, the performance decreases slightly, instead.

A.2 Forgetting Analysis

In addition to computing accuracy, we also measure the forgetting by computing the euclidean distance between the parameters of the novel model and the base model on CoLA corpus. Figure 3 shows the comparison of parameter forgetting from the first epoch to the last epoch and the corresponding accuracy after convergence with epoch increasing among LPC, RecAdam and Adam. In Figure 3 chart (1), with epoch increasing, the euclidean distance of Adam increases a lot, which means the forgetting of Adam is huge with the epoch increasing. However, our model (LPC) reduces the forgetting in a large extent compared with Adam and achieves similar forgetting with RecAdam, another baseline trying to reduce catastrophic forgetting. Here, the

forgetting of our model is a little bit worse than RecAdam is because our model tries to remember the most important parameters while forget unimportant parameters. Furthermore, in Figure 3 chart (2), we can see our model (LPC) achieves the best accuracy compared with RecAdam and Adam all the time after convergence (Epoch 4).

A.3 Ablation Study with ALBERT-xxlarge Model

As we have mentioned, our model (LPC) has two important components, Logits Calibration (LC) and Parameter Calibration (PC). In addition to doing ablation study with BERT-base pre-trained model, we also do ablation study on these two components separately with ALBERT-xxlarge pre-trained model on the 7 scenarios of the GLUE benchmark. The results of ablation study with ALBERT-xxlarge pre-trained model is shown in Table 3. We can see with ALBERT-xxlarge pre-trained model, both of LC and PC achieve better results than the baseline Adam. LPC achieves the best results among all three models. Compared with Adam, LC achieves 5.3% improvements on average measured by acc on MRPC, QNLI, RTE, SST-2, and WNLI, 0.7% im-

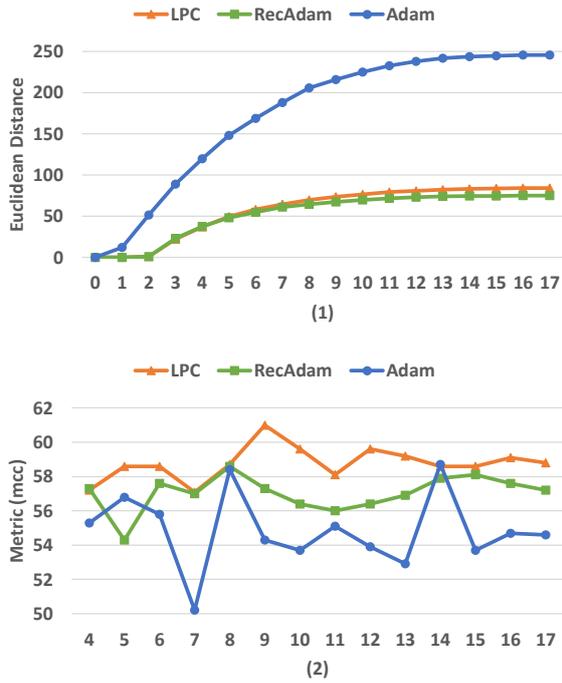


Figure 3: Comparison of Parameter Forgetting and Model Performance with the Epoch Increasing on CoLA Corpus with BERT-base Pre-trained Model.

828 improvements measured by mcc on CoLA, and 0.1%
 829 improvements measured by corr on STS-B. Compared with Adam, PC achieves 6.5% improvements
 830 on average measured by acc on MRPC, QNLI, RTE,
 831 SST-2, and WNLI, 5.1% improvements measured
 832 by mcc on CoLA, and 0.8% improvements measured
 833 by corr on STS-B. Compared with Adam,
 834 LPC achieves 8.7% improvements on average measured
 835 by acc on MRPC, QNLI, RTE, SST-2, and
 836 WNLI, 5.1% improvements measured by mcc on
 837 CoLA, and 1.2% improvements measured by corr
 838 on STS-B. Thus, we can conclude that our model
 839 can achieve state-of-the-art results with different
 840 pre-trained model. These results prove the scalability
 841 of our model.
 842