DATA EFFICIENT CONTINUAL LEARNING OF LARGE LANGUAGE MODEL

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

Paper under double-blind review

ABSTRACT

Continual Learning (CL) in large language models (LLMs) aims to enable models to learn from evolving data distributions while preserving previously acquired knowledge. However, existing CL methods primarily rely on statistical correlations from observed data, which are particularly vulnerable under limited data settings. This reliance results in two major drawbacks: (1) increased susceptibility to forgetting previously learned knowledge when data distribution shifts occur, and (2) a tendency to depend on spurious features instead of uncovering true causal relationships in new tasks. These issues become even more pronounced, especially when training data is limited. To address these challenges, we introduce a causality-guided CL approach that reinterprets CL through the lens of causal inference. Our method aims to mitigate the dependency of model parameters on the data inputs, leading to two key advantages: (1) reduced catastrophic forgetting, and (2) decreased dependence on spurious correlations, thereby improving generalization across both old and new tasks. Extensive experiments on pre-trained LLMs, including T5-large and Llama2, demonstrate that our approach significantly outperforms state-of-the-art (SOTA) CL methods in LLMs, particularly when the amount of training data is limited.

028 1 INTRODUCTION

Large Language Models (LLMs) have demonstrated remarkable performance across a wide range of 031 natural language processing (NLP) tasks, including text generation, translation, sentiment analysis, and question-answering (Brown et al., 2020; Achiam et al., 2024). These models have revolutionized 033 the field by achieving human-like proficiency in many applications, driven by their ability to process 034 and generate vast amounts of text with a deep understanding of context and semantics. However, despite these impressive achievements, a critical challenge remains: enabling LLMs to continually 035 learn and adapt to new information while preserving the knowledge acquired from previous tasks. This capability is essential for the development of truly autonomous and intelligent systems, as it 037 would allow LLMs to evolve over time Ke et al. (2023), expanding their knowledge base without suffering from catastrophic forgetting (CF) (McCloskey & Cohen, 1989) —a common problem where newly learned information interferes with previously acquired knowledge. Achieving this 040 level of continual learning (CL) is a crucial step toward realizing artificial general intelligence (AGI), 041 where a system can seamlessly integrate and apply knowledge across diverse domains, exhibiting a 042 level of adaptability and intelligence comparable to human cognition. 043

Current CL methods for LLMs primarily focus on capturing statistical correlations between input 044 data and labels. Although these approaches can be effective, it exposes critical vulnerabilities in 045 retaining past knowledge and acquiring new knowledge. (1) Firstly, these methods are particularly 046 susceptible to non-stationary data distributions during training. Most CL techniques are designed to 047 reinforce statistical relationships observed in the training data, which often lack resilience when en-048 countering non-stationary data distribution. This problem is especially pronounced in settings with limited labeled data. As training data distributions change, models that rely exclusively on statistical correlations struggle to generalize, resulting in degraded performance and intensified catastrophic 051 forgetting. (2) Second, during the learning of new tasks, CL models often memorize spurious features instead of identifying true predictive patterns in the data (Bombari & Mondelli, 2024). This 052 problem becomes more pronounced when training data is scarce, as the model may overemphasize coincidental patterns that appear to correlate with the output labels. In text classification, spurious 054 features are input patterns that seem related to the target label but lack a true causal connection to the 055 task. These misleading features can cause the model to make errors, especially when the true task-056 relevant features differ. For instance, consider a sentiment analysis model classifying movie reviews 057 as positive or negative. A spurious feature could be the word "Oscar," often seen in positive reviews but not necessarily a sign of sentiment. For example: "This movie was boring despite the Oscar 058 nomination." (Negative) "An amazing Oscar-winning performance!" (Positive). The model might mistakenly learn that "Oscar" always implies a positive sentiment, leading it to misclassify the first 060 sentence as positive despite its negative tone. Spurious features like this can make the model overfit 061 to irrelevant correlations in the training data, causing a significant drop in performance when faced 062 with new data that doesn't follow the same patterns. This challenge is worsened when training data 063 is limited, making the model overly dependent on superficial associations rather than genuine pre-064 dictive rules. Therefore, traditional CL approaches that rely on statistical correlations often struggle 065 to retain past knowledge and adapt effectively to new tasks, leading to suboptimal performance. 066

- Humans can learn efficiently from a limited number of examples by leveraging causal knowledge.
 Even with small amounts of data, understanding the causal relationships enables us to make accurate predictions and decisions. Causal inference aims to replicate this efficiency by identifying stable causal factors, which allows models to perform well even with limited data. Furthermore, humans intuitively understand that not all correlations imply causation and often disregard spurious relationships. Causal inference (Pearl, 2009), similarly, helps models avoid misleading correlations by adjusting for variables that might obscure the true causal effect.
- Inspired by how humans process new information while retaining old knowledge, and recognizing 074 the limitations of relying solely on statistical correlations in existing CL methods, we propose a 075 novel approach that models and leverages the causal relationships among different random variables 076 in CL for LLMs. Causal relationships go deeper than mere correlations. By identifying and mod-077 eling these causal connections, models can potentially become more robust to distributional shifts. They would not merely react to observed patterns but would instead infer the underlying mechanisms 079 driving those patterns. By grounding CL models in causal reasoning, we can enhance their ability to generalize across different contexts. When the data distribution shifts, a model that understands 081 the causal pathways can adjust more effectively rather than just memorizing surface-level correlations. This shift from a correlation-based approach to a causality-guided framework represents a fundamental change in how we approach CL in LLMs. 083
- 084 To achieve causal learning during CL in LLMs, we begin by analyzing the CL process through a 085 causal lens, modeling the relationships among key variables. Our findings reveal that the heavy dependence of CL model parameters on input data is a critical factor contributing to two major is-087 sues: (1) catastrophic forgetting of previously learned knowledge, as over-reliance on data inputs 880 causes significant parameter shifts when the training data distribution changes; and (2) the tendency to memorize spurious correlations in new tasks, leading to poor generalization. To address these 089 challenges, we propose a soft intervention approach that mitigates the influence of input data on model parameters. By reducing the impact of data inputs, our method offers two primary bene-091 fits: (1) diminished catastrophic forgetting of prior knowledge, as the model parameters become 092 less sensitive to shifts in data distribution; and (2) reduced memorization of spurious correlations when learning new tasks, enhancing the model's ability to identify true causal relationships. As a 094 result, our approach significantly improves performance across both previously learned and newly 095 introduced tasks, delivering a more robust and generalizable CL framework. 096
- To assess the effectiveness of our proposed method against SOTA approaches, we conduct extensive experiments on multiple datasets and pre-trained LLMs, including T5-large and Llama2, under limited labeled data settings. The results demonstrate that our method significantly outperforms existing SOTA CL methods for LLMs.
- 101 Our contributions can be summarized as follows:
- 102
- 103
- 104 105

- We present a novel and general causal machine learning framework applicable to a wide range of machine learning models and tasks. In particular, we tailor this general framework to CL in LLMs. By reinterpreting the CL process in LLMs through the lens of causality, we offer a fresh perspective on mitigating catastrophic forgetting.
- We develop a causality-guided CL algorithm that reduces the dependence of model parameters on input data, achieving two key advantages: minimizing forgetting of previously

learned knowledge and enhancing new task performance by mitigating the memorization
 of spurious correlations.

110 111

• Through extensive experiments on multiple datasets and large-scale language models, we demonstrate that our method significantly outperforms SOTA CL approaches in LLMs.

112 113

114 2 RELATED WORKS

115 2.1 TRADITIONAL CONTINUAL LEARNING

117 Traditional CL approaches can be broadly classified into the following categories: (1) regularization-118 based methods (Kirkpatrick et al., 2017; Li & Hoiem, 2017; Zenke et al., 2017; Chaudhry et al., 119 2018), which constrain model updates to prevent rapid changes and preserve previously learned 120 knowledge; (2) memory-replay-based methods (Chaudhry et al., 2019; Buzzega et al., 2020; Pham 121 et al., 2021; Caccia et al., 2022; Arani et al., 2022; Yang et al., 2023; Wang et al., 2024b), which store a small subset of examples from prior tasks in a memory buffer and replay them during training; (3) 122 architecture-based methods (Rusu et al., 2016; Li et al., 2019; Konishi et al., 2023; Thapa & Li, 123 2024), which either dynamically expand the model's capacity or isolate specific model weights to 124 prevent interference between tasks; and (4) gradient-projection-based methods (Saha et al., 2021; 125 Xiao et al., 2024), which project the current gradient update into the subspace defined by previous 126 tasks to reduce forgetting. 127

- 128 129
- 2.2 CONTINUAL LEARNING FOR LLM

130 Continual learning for LLM can be categorized into: (1) subspace-based method (Wang et al., 131 2023a); (2) prompt-based method (Qin & Joty, 2022; Razdaibiedina et al., 2023); (3) attention-based 132 method (Zhao et al., 2024); (4) architecture-based method (Wang et al., 2023b; 2024a). Orthogo-133 nal to these existing CL methods for LLMs, which rely on statistical correlations, these approaches 134 (1) tend to memorize spurious features during new task learning, and (2) are highly vulnerable to data distribution shifts, exacerbating forgetting, particularly in limited data settings. In contrast, 135 our method introduces a causality-based approach, which mitigates the impact of data distribution 136 shifts on model parameters and lessens the memorization of spurious features during model training. 137 This offers a novel strategy for improving continual learning performance in LLMs, ensuring better 138 retention of old knowledge and generalization on new task. 139

140 141

142

3 CAUSAL MACHINE LEARNING

Causal inference (Pearl et al., 2016; Pearl & Mackenzie, 2018) focuses on uncovering cause-effect 143 relationships among variables, providing a robust foundation for understanding and modeling data 144 beyond simple correlations. The principles of causality have been increasingly integrated into ma-145 chine learning (Parascandolo et al., 2018; Besserve et al., 2020), enabling models to incorporate 146 causal reasoning for improved generalization and robustness. Causal machine learning methods 147 have been applied to various domains, including few-shot learning (Yue et al., 2020), imitation 148 learning (De Haan et al., 2019), domain adaptation (Gong et al., 2016; Magliacane et al., 2018; 149 Kong et al., 2022) and representation learning (Schölkopf et al., 2021). These methods primarily 150 leverage hard interventions, where variables are explicitly manipulated or fixed to isolate causal 151 effects. While effective in many contexts, such approaches are not directly suitable for the con-152 tinual learning framework. In continual learning, the ability to adapt to new tasks while retaining knowledge from prior tasks is critical, but hard interventions often result in static model parameters. 153 This eliminates catastrophic forgetting but compromises the model's capacity to learn new tasks 154 effectively. 155

In contrast, our proposed method introduces soft interventions, which adjust the dependency of model parameters on input distributions without freezing them entirely. This approach strikes a balance between mitigating forgetting and maintaining adaptability to new tasks. By dynamically refining the influence of input distributions during learning, our method ensures continual improvement across tasks while preserving past knowledge. To the best of our knowledge, this is the first application of causal principles via soft interventions in the context of continual learning, addressing a unique set of challenges not tackled by prior causal ML methods.



Figure 1: Causal Diagram of CL in LLMs: the nodes represent key random variables: X denotes 174 the input data distribution, Y denotes the ground truth data label, Y^* denotes the model outputs 175 and Θ represents the model parameters. The arrows \rightarrow indicate causal relationships between these 176 variables. Specifically, the causal link $X \to \Theta \leftarrow Y$ highlights two major challenges in CL: (i) 177 Forgetting of Previously Learned Knowledge: As the data distribution (X, Y) shifts over time, it 178 directly influences the model parameters Θ , causing them to change and resulting in the loss of 179 previously acquired knowledge. (ii) Memorization of Spurious Features in New Tasks: Due to the 180 strong dependence of Θ on X, the model tends to memorize spurious correlations in the data, lead-181 ing to suboptimal performance on new tasks. (a) Traditional CL Methods: These methods do not 182 employ any intervention to address the causal link between X and Θ . As a result, they suffer from 183 the above challenges, including catastrophic forgetting and spurious feature memorization. (b) Hard Interventions in Traditional Causal Inference: Traditional causal inference methods propose hard 184 interventions on the link $X \to \Theta$, effectively removing the influence of X on Θ . While this pre-185 vents the model from forgetting old knowledge, it also hampers the ability to learn new task-specific information, making the model static and unresponsive to new data. (c) Our Proposed Soft Inter-187 vention: We introduce a novel soft intervention approach that *partially* mitigates the influence of X 188 on Θ , rather than completely severing the connection. This technique offers two key advantages: 189 (1) Reduced forgetting: By lessening the impact of data distribution shifts on model parameters, the 190 model retains previously learned knowledge more effectively while still incorporating new informa-191 tion. (2) Mitigated memorization of spurious features: The reduced reliance of Θ on X makes the 192 model less prone to overfitting spurious correlations, leading to better performance on new tasks. 193

194 195

196

4 DATA AUGMENTATION AND PARAMETER PERTURBATION

197 Traditional methods such as dropout, MixOut (Lee et al., 2020), and data augmentation primarily 198 focus on improving generalization during model training by introducing noise to the input data or model parameters to prevent overfitting. However, these approaches overlook the dependency 199 between model parameters and input data, which can result in suboptimal solutions. In contrast, 200 our proposed method specifically targets mitigating the dependency of model parameters on input 201 distributions, which is a crucial challenge in continual learning where input distributions change 202 over time. This focus on reducing distribution dependency is not a primary design goal of existing 203 methods like dropout or MixOut. 204

205 206

207

5 Method

In this section, we begin by outlining the problem setup and the limitations of existing CL approaches for LLMs in Section 5.1. Next, we introduce our proposed causality-guided CL approach in Section 5.2.

211 212

213

5.1 PROBLEM SETUP AND LIMITATIONS OF EXISTING WORKS

Problem Setup and Formulation Given a pre-trained LLM with parameters θ_0 , the goal of CL is to train the model sequentially on a series of N tasks, each with its own training dataset $\mathcal{D}_1^{tr}, \mathcal{D}_2^{tr}, \dots, \mathcal{D}_N^{tr}$. After training on the final task, the objective is to obtain a model f_{θ} parameter-

ized by θ that not only performs well on the current task but also retains high performance across the test sets of all previously learned tasks, i.e., $\mathcal{D}_1^{te}, \mathcal{D}_2^{te}, \ldots, \mathcal{D}_N^{te}$.

To mitigate forgetting during CL in LLMs, existing approaches aim to minimize the following loss function to update LLMs:

$$\mathcal{L} = \mathbb{E}_{(\boldsymbol{x}, y) \sim \mathcal{D}_{\star}^{tr}}(\boldsymbol{x}, y, \boldsymbol{\theta}) + \lambda \mathcal{L}_{old}(\boldsymbol{\theta})$$
(1)

where $\mathbb{E}_{(\boldsymbol{x},y)\sim \mathcal{D}_t^{tr}}(\boldsymbol{x}, y, \boldsymbol{\theta})$ represents the loss on the training data of the current new task t, i.e., \mathcal{D}_t^{tr} . $\mathcal{L}_{old}(\boldsymbol{\theta})$ is the loss term designed to preserve previously learned knowledge and mitigate forgetting. This term could involve a regularization penalty on the model parameters or a memory replay loss. The hyperparameter λ controls the balance between optimizing performance on the new task and preserving knowledge from old tasks. To better illustrate the limitations of CL optimization described in Eq. 1, we construct a causal diagram (Pearl, 2009) as explained below.

229 **Understanding CL in LLMs through a Causal Diagram** Figure 1 illustrates the causal relation-230 ships involved in CL for LLMs using a causal diagram. In this diagram, each node represents a random variable: X denotes the input data, Y^* represents the model's prediction, Y represents the 231 data ground truth label, and Θ symbolizes the model parameters. The arrows between nodes capture 232 the causal dependencies between these variables. For instance, the arrows $X \to Y^* \leftarrow \Theta$ suggest 233 that both the input data (X) and the model parameters (Θ) influence the predictions (Y^{*}). Simi-234 larly, the arrow $X \to \Theta$ indicates that shifts in the data distribution (X) directly affect the model 235 parameters (Θ). These causal relationships highlight two key limitations in existing CL methods: 236 (1) Changes in the data distribution (X) lead to shifts in the model parameters (Θ), which causes 237 the model to forget previously learned knowledge as it adapts to new data. This phenomenon is 238 known as catastrophic forgetting. (2) The model tends to memorize spurious features of the new 239 task due to its heavy reliance on the input data (X), rather than identifying true causal relationships. 240 Consequently, the performance on new tasks can degrade, especially when faced with unseen data. 241 Figure 1a shows how existing CL approaches, which rely on optimizing Eq. 1, are fundamentally statistic-correlation-based. In these approaches, the model parameters Θ are heavily influenced by 242 the observed data statistics in X, making them vulnerable to shifts and spurious patterns in the data. 243

244 245

265

221 222

223

224

225

226

227

228

5.2 CAUSALITY-GUIDED CONTINUAL LEARNING FOR LLMS

To address the problem of existing approaches that heavily rely on the statistical correlation between input data X and model parameters Θ (i.e., $X \to \Theta$), we propose a novel method grounded in causal inference. Our approach mitigates the causal dependency between the input variable X and the model parameters Θ , enhancing the model's generalization.

251 **Traditional Causal Inference (Hard Intervention)** The do-operation $(do(\Theta = \theta))$: The do-252 operator, introduced by Judea Pearl (Pearl, 2009), is a key concept in causal inference, providing 253 a formal framework to distinguish causation from mere correlation by modeling the effects of in-254 terventions in a system. The do-operation represents an intervention where a variable, in this case, 255 Θ (model parameters), is forcibly set to a specific value θ , disregarding its natural causes. This intervention allows us to estimate the direct causal effect of Θ on another variable, such as Y 256 (model outputs), which is expressed as $P(y|\mathsf{do}(\Theta = \theta))$. Unlike standard conditional probabil-257 ities like $P(y|\Theta = \theta)$, which capture the observed relationships between Θ and Y in the data, 258 $P(y|\mathsf{do}(\Theta = \theta))$ models the outcome when Θ is actively set to θ , effectively severing any natural 259 influences on Θ . However, as illustrated in Figure 1b, implementing a hard intervention that com-260 pletely severs the causal link between the input data X and model parameters Θ would result in the 261 model being immune to forgetting because Θ no longer depends on X. While this prevents catas-262 trophic forgetting, it also hinders the model's ability to learn new tasks, as the model parameters 263 remain fixed and unresponsive to new input data. This limitation emphasizes the trade-off inherent 264 in traditional causal interventions within CL settings.

Our Method (Soft Intervention) To overcome the limitations of hard interventions in traditional causal inference, we propose a novel approach: a partial (soft) intervention on the causal relationship between the data input X and the model parameters Θ . Unlike hard interventions that completely disconnect Θ from X, soft interventions allow for a controlled, partial dependence, effectively balancing the model's need to retain previously learned knowledge while still adapting to new information. By softly intervening on Θ , we can reduce the influence of data shifts on model parameters, thereby mitigating catastrophic forgetting of previously learned tasks. At the same time, this approach retains enough flexibility to allow the model to learn new tasks effectively. The goal is to limit the degree to which Θ is affected by changes in X, thereby reducing the impact of spurious correlations without entirely severing the ability of the model to adjust and learn. As illustrated in Figure 1c, our soft intervention approach modifies the influence of input data on model parameters, striking a balance that retains learned knowledge while enabling the integration of new information. This targeted, controlled intervention can be mathematically expressed as follows:

$$P(y|do(soft(\Theta))) = \sum_{\boldsymbol{x} \in \mathcal{X}} P(\boldsymbol{x}) \sum_{\boldsymbol{\theta} \sim \boldsymbol{\Theta}} P(y|\boldsymbol{x}, \boldsymbol{\theta}) P'(\boldsymbol{\theta}|\boldsymbol{x})$$
(2)

where $P'(\theta|\mathbf{x})$ denotes the modified parameter distribution conditioned on input \mathbf{x} , reflecting the partial (soft) intervention on Θ . We set it to be $P'(\theta|\mathbf{x}) \approx P(\theta|\mathbf{x}) + \mathcal{N}(\mathbf{0}, \sigma^2)$, where σ are learnable standard deviation parameters. This modified posterior distribution $P'(\theta|\mathbf{x})$, diverges from the original $P(\theta|\mathbf{x})$, enabling Θ to partially depend on X. As a result, it reduces the strong dependency of Θ on X. This equation is derived using backdoor adjustment from causal inference. The backdoor path is $\Theta \leftarrow X \rightarrow Y^*$, and the variable X is not a descendant of Θ . Therefore, X satisfies the backdoor criterion with respect to the causal effect of Θ on Y^* . Consequently, we can apply the backdoor adjustment to account for the confounding variable and estimate the causal effect.

To reduce the reliance of the model parameters Θ on data inputs X and enhance the relation between Θ and data label Y in a soft way, we propose maximizing the following learning objective:

$$\mathcal{L}_{causal} = I(\boldsymbol{\Theta}, Y) - \alpha I(\boldsymbol{\Theta}, X) + \lambda \mathcal{L}_{old}(\boldsymbol{\theta})$$
(3)

where $I(\Theta, Y)$ denotes the mutual information between Θ and Y, defined as $I(\Theta, Y) = \int P(y, \theta) \log \frac{P(y, \theta)}{P(y)P(\theta)} dy d\theta$, which quantifies the dependence or shared information between the variables Y and Θ . Maximizing $I(\Theta, Y)$ would enhance the prediction of label y with parameters θ . $I(\Theta, X)$ denotes the mutual information between Θ and X. Minimizing $I(\Theta, X)$ reduces the dependence of Θ on X. $\alpha > 0$ is a weighting constant that balances the two mutual information terms. However, calculating mutual information is computationally intractable due to the difficulty of modeling the joint distribution $P(y, \theta)$, and marginal distribution $P(\theta)$ in high-dimensional spaces. Therefore, inspired by the variational inference Alemi et al. (2017), we propose maximizing the following variational objective:

$$\mathcal{L}_{causal} \approx \frac{1}{M} \sum_{i=1}^{i=M} \mathbb{E}_{\boldsymbol{\epsilon} \sim P(\boldsymbol{\epsilon})}[\log P(y|\mathsf{do}(\mathsf{soft}(\boldsymbol{\Theta})), \boldsymbol{\epsilon})] - \alpha \mathbb{KL}(P'(\boldsymbol{\theta}|\boldsymbol{x}_i), Q(\boldsymbol{\theta})) + \lambda \mathcal{L}_{old}(\boldsymbol{\theta}) \quad (4)$$

where M denotes the number of training data points. $P'(\boldsymbol{\theta}|\boldsymbol{x}_i)$ denotes the modified model parame-ter posterior distribution given the input x_i as Eq. 2, and $Q(\theta)$ is the model prior distribution, which is set to be standard normal distribution, i.e., $Q(\theta) = \mathcal{N}(0, I)$. We put the detailed derivations of Eq. 4 in Appendix A. In addition, to simulate samples from the distribution $P(\epsilon)$, i.e., generate noisy versions of the input text, we can use the following transformations: (1) WordNet synonym replacement, which randomly replaces words with their synonyms; (2) Word deletion, which ran-domly removes words from the sentence; (3) Word order swaps, which randomly swap the positions of words in the sentence; and (4) Random synonym insertion, which inserts a synonym of a random word at a random position. These techniques collectively transform the input text to build noisy data from $P(\epsilon)$. Furthermore, it is important to note that the last term in Eq. 4, $\lambda \mathcal{L}_{old}(\theta)$, represents the loss term introduced in previous works to mitigate forgetting of old knowledge. The hyperparameter λ is thus not part of our method. Consequently, α is the only hyperparameter specific to our approach. We refer to our proposed approach as Causality-Guided Continual Learning (CGCL). The detailed steps of our algorithm are provided in Algorithm 1.

324 Algorithm 1 Causality-Guided Continual Learning for LLMs. 325 1: **REQUIRE:** pre-trained LLM parameters θ_0 , learning rate η . 326 2: for n = 1 to N do (number of CL tasks) 327 3: for k = 1 to K do (number of CL steps) 328 4: calculate the data intervention by Eq. 2. 5: calculate the causal learning objective by Eq. 4. 6: update the LLM parameters by $\boldsymbol{\theta}_{k+1}^n = \boldsymbol{\theta}_k^n - \eta \nabla \mathcal{L}_{causal}(\boldsymbol{\theta}_k^n)$ 330 7: end for 331 8: end for 332 333 334 EXPERIMENT 6 335 336 In this Section, we first provide the experiment setup in Section 6.1. Then, we present the experiment 337 results in Section 6.2. Next, we present more detailed analysis and ablation study in Section 6.3. 338 339 6.1 Setup 340 341 **Datasets** Following (Razdaibiedina et al., 2023; Wang et al., 2023a), we use the benchmark that includes a variety of NLP tasks, each accompanied by expert-crafted instructions, to provide a more 342 practical evaluation framework for CL in LLMs. We use two benchmark datasets to evaluate the 343 performance of CL methods in LLMs. 344 345 • standard CL benchmark: The four text classification tasks (Zhang et al., 2015) are shuf-346 fled into three distinct sequences, forming order 1, 2, and 3, for use in the standard CL 347 benchmark. 348 • long sequence CL benchmark: This benchmark consists of a total of 15 tasks, which present 349 additional challenges to existing CL approaches. Specifically, it includes five tasks from the 350 CL benchmark, four tasks from the GLUE benchmark (MNLI, QQP, RTE, SST2) (Wang 351 et al., 2019b), five tasks from the SuperGLUE benchmark (WiC, CB, COPA, MultiRC, 352 BoolQ) (Wang et al., 2019a), and the IMDB movie reviews dataset (Maas et al., 2011), 353 creating long sequence CL benchmark orders 4, 5, and 6. The task sequence order can be 354 found in Appendix B. 355 356 **Baselines** We compare to the following SOTA baselines: LFPT5 (Qin & Joty, 2022) is a replay-357 based method that continuously trains a soft prompt designed to both solve current tasks and generate pseudo-labeled samples from previously learned domains. These generated samples are then utilized in experience replay to reinforce past knowledge. EPI (Wang et al., 2023b) employs a parameter 359 isolation strategy to reduce forgetting by allocating a small set of task-specific parameters for each 360 task, which are learned alongside a shared pre-trained model. ProgPrompt (Razdaibiedina et al., 361 2023) incrementally learns a new soft prompt for each new task and sequentially append it to the 362 prompts learned from previous tasks. O-LoRA (Wang et al., 2023a) which builds on LoRA (Hu et al., 2022) and learns tasks in distinct low-rank vector subspace that are maintained orthogonal to 364 each other, effectively minimizing interference between tasks. SAPT (Zhao et al., 2024) aligns the 365 parameter-efficient tuning (PET) learning and selection through a shared attention-based learning 366 and selection module. In addition, we also compare to the data augmentation (DA) techniques in 367 (Wei & Zou, 2019) as an additional strong baseline. 368 369 **Pre-trained Models** Following the setting in (Wang et al., 2023a; Zhao et al., 2024), we use the pre-trained T5-Large (Raffel et al., 2020) and LLaMA-2-7B (Touvron et al., 2023).

Evaluation Metrics We define $a_{i,j}$ as the evaluation performance on the *j*-th task after training on the *i*-th task. Following (Wang et al., 2023a), we utilize the average performance across the CL task sequence to assess the performance of CL on LLM. The overall performance across all tasks after completing the training on the final task is defined as: $A_{\mathcal{T}} = \frac{1}{\mathcal{T}} \sum_{t=1}^{\mathcal{T}} a_{\mathcal{T},t}$.

376

370 371

Implementation Details Our implementation is based on the codebase of O-LORA (Wang et al., 2023a). All the experiment results are averaged over 3 runs and are performed on 4×A6000 NVIDIA

		Standard CL Benchmark				Long Sequence Benchmark			
	Order-1	Order-2	Order-3	Avg	Order-4	Order-5	Order-6	Avg	
ProgPrompt	55.39	56.67	51.38	54.48	46.28	48.91	36.75	43.98	
LFPT5	52.20	54.31	49.52	52.01	40.46	46.25	32.83	39.85	
EPI	43.62	47.29	42.09	44.33	31.43	40.37	27.51	33.10	
SAPT-LoRA	60.51	61.77	58.23	60.17	47.34	55.89	40.93	48.05	
O-LoRA	59.42	59.46	55.68	58.19	48.20	52.06	39.78	46.68	
O-LoRA + DA	61.83	62.61	59.97	61.47	51.57	54.72	43.64	49.98	
O-LoRA + CGCL	73.66±2.03	$75.38{\pm}1.76$	72.75±1.94	73.93±1.89	64.81±2.53	67.60±2.67	60.87±2.36	64.43±2.5	

Table 1: The overall results on two continual learning benchmarks with T5-Large model with sam-ple size of 100.

Table 2: The overall results on two continual learning benchmarks with LLaMA-2-7B model with sample size of 100.

	Standard CL Benchmark				Long Sequence Benchmark			
	Order-1	Order-2	Order-3	Avg	Order-4	Order-5	Order-6	Avg
ProgPrompt	62.03	45.69	61.40	56.37	52.35	51.82	52.31	52.16
LFPT5	57.17	42.51	53.76	51.15	46.08	47.34	50.75	48.06
EPI	46.32	32.49	52.71	43.84	43.24	42.53	43.67	43.15
SAPT-LoRA	69.37	51.95	67.83	63.05	59.63	59.21	61.32	60.05
O-LoRA	68.24	50.63	65.68	61.52	57.45	56.79	59.28	57.84
O-LoRA + DA	69.05	52.76	67.72	63.18	59.91	59.68	62.86	60.82
O-LoRA + CGCL	67.91±1.93	69.94±1.82	$72.17{\pm}1.58$	$70.01{\pm}1.76$	66.71±2.06	65.35±1.97	69.99±1.83	67.35±1.9

GPU using DeepSpeed repository. We set the word perturbation rate in each sentence to 10% in order to generate noisy samples from $P(\epsilon)$, and $\alpha = 0.003$. LoRA configuration: r = 8. The learning rate is set to be 1e-4. More implementation details can be found in Appendix C.

6.2 RESULTS

We present the CL results with a sample size of 100 per task for the T5-Large model in Table 1 and the LLaMA-2-7B model in Table 2. Our method significantly exceeds the SOTA CL performance, achieving improvements of over 12% with T5-Large and more than 7% with LLaMA-2-7B on the standard CL benchmark. On the long-sequence CL benchmark, our approach demonstrates even greater gains, improving SOTA performance by over 15% with T5-Large and more than 6% with LLaMA-2-7B. These results emphasize the effectiveness of our causality-guided approach in data-efficient learning scenarios.

Effect of Sample Size for Each Task To evaluate the impact of sample size on task performance, we conduct experiments using a sample size of 500 for each task. This larger sample size allows us to explore how varying the amount of training data influences the effectiveness of our CGCL approach. The CL results for the T5-large model with a sample size of 500 are presented in Table 3, while the results for the LLaMA-2-7B model are detailed in Table 4. Our findings indicate that even with this increased sample size, our method continues to deliver substantial improvements in overall CL performance. This suggests that our approach is robust and effective, demonstrating its ability to enhance learning outcomes regardless of the amount of available training data. The consistent performance gains across both models further validate the advantages of incorporating a causality-guided perspective in CL tasks.

Table 3: The overall results on two continual learning benchmarks with T5-Large model with sample size of 500.

		Standard CL Benchmark				Long Sequence Benchmark			
	Order-1	Order-2	Order-3	Avg	Order-4	Order-5	Order-6	Avg	
ProgPrompt	70.31	71.18	72.05	71.18	52.20	64.31	62.16	59.56	
LFPT5	62.04	61.37	60.90	61.43	47.95	46.26	43.35	45.85	
EPI	64.29	66.73	65.64	65.55	49.32	61.43	58.62	56.46	
SAPT-LoRA	75.32	75.86	76.37	75.85	60.28	70.24	69.83	66.78	
O-LoRA	74.15	75.29	75.71	75.05	59.83	69.38	69.20	66.14	
O-LoRA + DA	75.81	76.52	76.93	76.42	62.56	70.91	71.22	68.23	
O-LoRA + CGCL	75.39±1.65	$78.16{\pm}1.72$	$78.42{\pm}1.43$	77.32±1.59	69.30±1.90	72.39±1.87	$76.88{\pm}1.16$	$72.86{\pm}1.68$	

	Standard CL Benchmark				Long Sequence Benchmark			
	Order-1	Order-2	Order-3	Avg	Order-4	Order-5	Order-6	Avg
ProgPrompt	62.56	70.25	60.58	64.46	57.03	56.67	61.82	58.51
LFPT5	51.43	60.18	54.71	55.44	55.97	50.60	52.29	52.95
EPI	57.95	53.52	59.41	56.96	53.35	52.38	54.67	53.47
SAPT-LoRA	68.70	75.67	70.32	71.56	64.51	63.42	66.35	64.76
O-LoRA	66.51	73.36	68.23	69.37	62.64	61.06	65.42	63.04
O-LoRA + DA	68.62	75.79	71.15	71.85	63.89	62.90	66.78	64.52
O-LoRA + CGCL	74.81±1.08	$78.03{\pm}0.95$	79.61±0.79	$77.48{\pm}0.93$	68.36±1.02	66.54±1.29	69.46±1.30	$68.12{\pm}1.21$

Table 4: The overall results on two continual learning benchmarks with LLaMA-2-7B model with sample size of 500.

6.3 ANALYSIS

Individual Task Performance Evaluation To assess the performance of individual (both old and new) tasks, we conduct a comparative analysis of each task under a data limit of 100 examples, utilizing both the T5-large model and LLaMA-2-7B within the context of task order 6. The performance results for the T5-large model are presented in Figure 2, while the corresponding results for LLaMA-2-7B are depicted in Figure 3. Our findings reveal that, with the implementation of our CGCL approach, the performance of each task shows a significant improvement across most datasets. This enhancement can be attributed to the effectiveness of our causality-guided methodology, which diminishes the dependency of model parameters on input data. Specifically, this reduction in dependency leads to two key benefits: (1) it mitigates the risk of forgetting previously learned task knowledge, and (2) it lessens the likelihood of memorizing spurious features associated with new tasks. Consequently, our proposed approach elevates superior overall performance across both previously learned old and new tasks.



Figure 2: Performance of individual tasks in the CL sequence using the T5-large model, with a data limit of 100 samples per task under the task order 6.

Hyperparameter Sensitivity Analysis To evaluate the sensitivity of the hyperparameter α in Eq. 4, we perform a sensitivity analysis, with the results summarized in Table 5. The findings indicate that our method is not very sensitive to α .





Figure 3: Performance of individual tasks in the CL sequence using the LLaMA-2-7B model, with a data limit of 100 samples per task under the task order 6.

Table 5: Hyperparameter analysis of α on CL performance under sample size of 100 with T5-large

α	Order 4	Order 5	Order 6
0.001	62.73	68.17	62.32
0.003	64.81	67.60	60.87
0.005	63.06	65.28	61.53

514 515 516

517

524

525 526

527

528 529

530

531

532

533

534

507

508 509

7 CONCLUSION

This paper explores continual learning for LLMs in a limited data setting. We introduce a novel causality-guided approach that addresses two key challenges: mitigating forgetting of previously learned tasks and avoiding the memorization of spurious features in new task data. Extensive experiments on large-scale pre-trained models, including T5-large and LLaMA-2, highlight the effectiveness of the proposed method.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and et al. Gpt-4 technical report, 2024. URL https://arxiv.org/abs/2303.08774.
- Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. Deep variational information bottleneck. In *International Conference on Learning Representations*, 2017.
- Elahe Arani, Fahad Sarfraz, and Bahram Zonooz. Learning fast, learning slow: A general continual learning method based on complementary learning system. In *International Conference on Learning Representations*, 2022.
- Michel Besserve, Arash Mehrjou, Rémy Sun, and Bernhard Schölkopf. Counterfactuals uncover the modular structure of deep generative models. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=SJxDDpEKvH.
- 539 Simone Bombari and Marco Mondelli. How spurious features are memorized: Precise analysis for random and NTK features. In *Forty-first International Conference on Machine Learning*, 2024.

540	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
540	Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
542	Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,
545	Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Baniamin Chasa, Jack Clerk, Christenhan Barnen, Sam McCandlich, Alas Badford, Ilya Suteksyan
544	and Dario Amodei I anguage models are few shot learners. In Advances in Neural Information
545	Processing Systems, volume 33, 2020.
547	
548	Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark ex-
549	perience for general continual learning: a strong, simple baseline. Advances in neural information
550	processing systems, 55:15920–15950, 2020.
551	Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky.
552	New insights on reducing abrupt representation change in online continual learning. In Interna-
553	tional Conference on Learning Representations, 2022.
554	Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Aianthan, and Philip HS Torr. Riemannian
555	walk for incremental learning: Understanding forgetting and intransigence. In Proceedings of the
556	European conference on computer vision (ECCV), pp. 532–547, 2018.
557	Arslan Chaudhry Marc' Aurelia Ranzata Marcus Rahrhach and Mahamed Elhoseiny Efficient
558	lifelong learning with a-gem. In International Conference on Learning Representations 2019
559	incloing learning with a geni. In machanolaa conjerence on Learning Representations, 2019.
560	Pim De Haan, Dinesh Jayaraman, and Sergey Levine. Causal confusion in imitation learning. Ad-
561	vances in neural information processing systems, 32, 2019.
562	Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard
563	Schölkopf. Domain adaptation with conditional transferable components. In International con-
564	ference on machine learning, pp. 2839–2848. PMLR, 2016.
565	Edward I Hu, velong shen, Phillin Wallis, Zevuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu, Wang
500	and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In <i>International Con</i> -
562	ference on Learning Representations, 2022.
569	$\mathbf{Z}' = \mathbf{W} \cdot \mathbf{W}'' \cdot \mathbf{W} = \mathbf{W} \cdot \mathbf{W}$
570	Lixuan Ke, 11jia Snao, Haowei Lin, Taisuya Konisni, Gyunak Kim, and Bing Liu. Continual pre-
571	tions. 2023.
572	
573	Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In <i>International Conference</i>
574	on Learning Representations, 2013.
575	James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
576	Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom-
577	ing catastrophic forgetting in neural networks. Proceedings of the national academy of sciences,
578	114(13):3521–3526, 2017.
579	Lingjing Kong, Shaoan Xie, Weiran Yao, Yujia Zheng, Guangyi Chen, Petar Stojanov, Victor Akin-
580	wande, and Kun Zhang. Partial disentanglement for domain adaptation. In International confer-
581	ence on machine learning, pp. 11455–11472. PMLR, 2022.
582	Tatsuva Konishi Mori Kurokawa Chihiro Ono Zixuan Ke Gyuhak Kim and Ring Liu Parameter-
583	level soft-masking for continual learning. In International Conference on Machine Learning, pp.
384 595	17492–17505. PMLR, 2023.
586	Chaolbuoung Lao Kuunghuun Cha and Warma Kang Miyout, Effective resultation to Cost
587	large-scale pretrained language models. In International Conference on Learning Represented
588	tions, 2020.
589	
590	Xilai Li, Yingbo Zhou, Tianfu Wu, Richard Socher, and Caiming Xiong. Learn to grow: A continual
591	structure learning framework for overcoming catastrophic forgetting. In International conference
592	on machine learning, pp. 5925-5954. FWILK, 2019.
593	Zhizhong Li and Derek Hoiem. Learning without forgetting. <i>IEEE transactions on pattern analysis and machine intelligence</i> , 40(12):2935–2947, 2017.

594 595 596	Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th annual meeting of the</i> <i>association for computational linguistics: Human language technologies</i> , pp. 142–150, 2011.
597 598 599 600	Sara Magliacane, Thijs Van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions. <i>Advances in neural information processing systems</i> , 31, 2018.
601 602 603	Michael McCloskey and Neal J Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. In <i>Psychology of learning and motivation</i> , volume 24, pp. 109–165. Elsevier, 1989.
604 605 606	Giambattista Parascandolo, Niki Kilbertus, Mateo Rojas-Carulla, and Bernhard Schölkopf. Learning independent causal mechanisms. In <i>International Conference on Machine Learning</i> , pp. 4036–4044. PMLR, 2018.
608	Judea Pearl. Causality. Cambridge university press, 2009.
609 610 611	Judea Pearl and Dana Mackenzie. The book of why: the new science of cause and effect. Basic books, 2018.
612 613	Judea Pearl, Madelyn Glymour, and Nicholas P Jewell. <i>Causal inference in statistics: A primer</i> . John Wiley & Sons, 2016.
615 616 617	Quang Pham, Chenghao Liu, and Steven HOI. Dualnet: Continual learning, fast and slow. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), <i>Advances in Neural Information Processing Systems</i> , 2021.
618 619	Chengwei Qin and Shafiq Joty. LFPT5: A unified framework for lifelong few-shot language learning based on prompt tuning of t5. In <i>International Conference on Learning Representations</i> , 2022.
620 621 622	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67, 2020.
624 625 626	Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Alma- hairi. Progressive prompts: Continual learning for language models. In <i>The Eleventh Interna-</i> <i>tional Conference on Learning Representations</i> , 2023.
627 628 629	Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. <i>arXiv preprint arXiv:1606.04671</i> , 2016.
630 631 632	Gobinda Saha, Isha Garg, and Kaushik Roy. Gradient projection memory for continual learning. In <i>International Conference on Learning Representations</i> , 2021.
633 634 635	Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. <i>Proceedings of the IEEE</i> , 109(5):612–634, 2021.
636 637 638	Jeevan Thapa and Rui Li. Bayesian adaptation of network depth and width for continual learning. In <i>Forty-first International Conference on Machine Learning</i> , 2024.
639 640 641	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
642 643 644	Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. Superglue: A stickier benchmark for general-purpose language understanding systems. <i>Advances in neural information processing systems</i> , 32, 2019a.
646 647	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In <i>International Conference on Learning Representations</i> , 2019b.

048	Mingyang Wang, Heike Adel, Lukas Lange, Jannik Strötgen, and Hinrich Schütze. Rehearsal-free
649	modular and compositional continual learning for language models. In <i>Proceedings of the 2024</i>
650	Conference of the North American Chapter of the Association for Computational Linguistics:
651	Human Language Technologies (Volume 2: Short Papers), pp. 469–480, 2024a.

- Xiao Wang, Tianze Chen, Qiming Ge, Han Xia, Rong Bao, Rui Zheng, Qi Zhang, Tao Gui, and Xuan-Jing Huang. Orthogonal subspace learning for language model continual learning. In Findings of the Association for Computational Linguistics: EMNLP 2023, pp. 10658–10671, 2023a.
- Zhenyi Wang, Yan Li, Li Shen, and Heng Huang. A unified and general framework for continual learning. In The Twelfth International Conference on Learning Representations, 2024b.
- Zhicheng Wang, Yufang Liu, Tao Ji, Xiaoling Wang, Yuanbin Wu, Congcong Jiang, Ye Chao, Zhen-cong Han, Ling Wang, Xu Shao, et al. Rehearsal-free continual language learning via efficient parameter isolation. In Proceedings of the 61st Annual Meeting of the Association for Computa-tional Linguistics (Volume 1: Long Papers), pp. 10933–10946, 2023b.
- Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Process-ing (EMNLP-IJCNLP), pp. 6382-6388, 2019.
- Mingqing Xiao, Qingyan Meng, Zongpeng Zhang, Di He, and Zhouchen Lin. Hebbian learning based orthogonal projection for continual learning of spiking neural networks. In The Twelfth International Conference on Learning Representations, 2024.
- Enneng Yang, Li Shen, Zhenyi Wang, Tongliang Liu, and Guibing Guo. An efficient dataset conden-sation plugin and its application to continual learning. Advances in Neural Information Processing Systems, 36, 2023.
- Zhongqi Yue, Hanwang Zhang, Qianru Sun, and Xian-Sheng Hua. Interventional few-shot learning. Advances in neural information processing systems, 33:2734–2746, 2020.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In International conference on machine learning, pp. 3987–3995. PMLR, 2017.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text clas-sification. Advances in neural information processing systems, 28, 2015.
- Weixiang Zhao, Shilong Wang, Yulin Hu, Yanyan Zhao, Bing Qin, Xuanyu Zhang, Qing Yang, Dongliang Xu, and Wanxiang Che. Sapt: A shared attention framework for parameter-efficient continual learning of large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 11641–11661, 2024.

Appendix CAUSAL LEARNING OBJECTIVE А $\mathcal{L}_{causal} = I(\boldsymbol{\Theta}, Y) - \alpha I(\boldsymbol{\Theta}, X)$ To minimize Eq. 5, inspired by (Alemi et al., 2017), we obtain the following the variational objective: Given $I(\Theta, Y)$ and $I(\Theta, X)$, we have the following derivation. (1) For $I(\Theta, Y)$, we can derive the following inequality: $I(\boldsymbol{\Theta}, Y) = \int d\boldsymbol{\theta} \, dy \, P(\boldsymbol{x}, \boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, y)}{P(\boldsymbol{\theta})P(y)} = \int d\boldsymbol{\theta} \, dy \, P(y, \boldsymbol{\theta}) \log \frac{P(y|\boldsymbol{\theta})}{P(y)}$ For $\mathbb{KL}(p(Y|\Theta), q(Y|\Theta)) \ge 0$, we can have the following inequality: $\int dy d\theta P(y|\theta) \log P(y|\theta) \geq \int dy d\theta P(y|\theta) \log q(y|\theta)$ Then, we have the following inequality: $I(\boldsymbol{\Theta}, Y) \ge \int d\boldsymbol{\theta} \, dy \, P(y, \boldsymbol{\theta}) \log \frac{q(y|\boldsymbol{\theta})}{P(y)}$ $= \int d\boldsymbol{\theta} \, dy \, P(y, \boldsymbol{\theta}) \log q(y|\boldsymbol{\theta}) - \int \, dy \, P(y) \log P(y)$ Since $P(y) \log P(y)$ does not depend on the model parameters θ , so this term could be ignored. We can also obtain the following equality: $P(y,\boldsymbol{\theta}) = \int d\boldsymbol{x} P(\boldsymbol{x}, y, \boldsymbol{\theta}) = \int d\boldsymbol{x} P(\boldsymbol{x}) P(y|\boldsymbol{x}) P(\boldsymbol{\theta}|\boldsymbol{x})$ By plugging-in Eq. (10) into Eq. (8), we can have the following bound: $I(\boldsymbol{\Theta}, Y) \geq \int d\boldsymbol{x} \, dy \, d\boldsymbol{\theta} \, P(\boldsymbol{x}) P(\boldsymbol{y}|\boldsymbol{x}) P(\boldsymbol{\theta}|\boldsymbol{x}) \log q(\boldsymbol{y}|\boldsymbol{\theta}).$

(2) For $I(\Theta, X)$, we can derive the following inequality:

$$I(\boldsymbol{\Theta}, X) = \int d\boldsymbol{\theta} \, d\boldsymbol{x} \, P(\boldsymbol{x}, \boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}, \boldsymbol{x})}{P(\boldsymbol{\theta})P(\boldsymbol{x})} = \int d\boldsymbol{\theta} \, d\boldsymbol{x} \, P(\boldsymbol{x}, \boldsymbol{\theta}) \log \frac{P(\boldsymbol{\theta}|\boldsymbol{x})}{P(\boldsymbol{\theta})}$$
(12)

$$= \int d\boldsymbol{\theta} \, d\boldsymbol{x} \, P(\boldsymbol{x}, \boldsymbol{\theta}) \log P(\boldsymbol{\theta} | \boldsymbol{x}) - \int d\boldsymbol{\theta} \, P(\boldsymbol{\theta}) \log P(\boldsymbol{\theta}) \tag{13}$$

where

$$P(\boldsymbol{\theta}) = \int d\boldsymbol{x} \ P(\boldsymbol{\theta}|\boldsymbol{x}) \ P(\boldsymbol{x})$$
(14)

(5)

(6)

(7)

(8)

(9)

(10)

(11)

We can use $r(\theta)$ be a variational approximation to $P(\theta)$, according to the inequality of $\mathbb{KL}(P(\boldsymbol{\theta}), r(\boldsymbol{\theta})) > 0$, we can obtain the following inequality:

$$\int d\boldsymbol{\theta} P(\boldsymbol{\theta}) \log P(\boldsymbol{\theta}) \ge \int d\boldsymbol{\theta} P(\boldsymbol{\theta}) \log r(\boldsymbol{\theta})$$
(15)

Based on Eq. (14), Eq. (12) can be rewritten as (16)

754
755
$$I(\boldsymbol{\Theta}, X) \leq \int d\boldsymbol{x} \, d\boldsymbol{\theta} \, P(\boldsymbol{x}) P(\boldsymbol{\theta}|\boldsymbol{x}) \log \frac{P(\boldsymbol{\theta}|\boldsymbol{x})}{Q(\boldsymbol{\theta})}.$$
 (16)

Combining the two inequality (11) and (16), we have that:

$$I(\boldsymbol{\Theta}, Y) - I(\boldsymbol{\Theta}, X) \ge \int d\boldsymbol{x} dy d\boldsymbol{\theta} P(\boldsymbol{x}) P(\boldsymbol{y}|\boldsymbol{x}) P(\boldsymbol{\theta}|\boldsymbol{x}) \log q(\boldsymbol{y}|\boldsymbol{\theta})$$
(17)

$$-\int d\mathbf{x} d\boldsymbol{\theta} P(\mathbf{x}) P(\boldsymbol{\theta}|\mathbf{x}) \log \frac{P(\boldsymbol{\theta}|\mathbf{x})}{Q(\boldsymbol{\theta})} = L$$
(18)

To efficiently calculate Eq. (17), we approximate the P(x, y) = P(x)P(y|x) by empirical data distribution $P(x, y) = \frac{1}{N} \sum_{n=1}^{N} \delta_{x_n}(x) \delta_{y_n}(y)$, δ_{x_n} is the Dirac delta function on x_n and δ_{y_n} is the Dirac delta function on y_n .

Then we can use the reparameterization trick (Kingma & Welling, 2013) to rewrite the $P(\theta|x)d\theta = P(\epsilon)d\epsilon$. This enables the distribution $P(\theta|x)$ to be reparameterized as a function of ϵ . We then calculate the KL divergence between $P(\theta|x)$ and $Q(\theta)$, and combine all together to minimize the following empirical loss function:

$$\mathcal{L}_{causal} \approx \frac{1}{M} \sum_{i=1}^{i=M} \mathbb{E}_{\boldsymbol{\epsilon} \sim P(\boldsymbol{\epsilon})}[\log P(y_i | f(\boldsymbol{x}_i, \boldsymbol{\epsilon}))] - \alpha \mathbb{KL}(P(\boldsymbol{\theta} | \boldsymbol{x}_i), Q(\boldsymbol{\theta}))$$
(19)

Next, we derive the soft intervention formula as below:

=

$$P(y|\text{do}(\text{soft}(\Theta))) = \sum_{\boldsymbol{x} \in X} \sum_{\boldsymbol{\theta} \sim \Theta} P(y|\boldsymbol{x}, \boldsymbol{\theta}) P'(\boldsymbol{\theta}|\boldsymbol{x}) P(\boldsymbol{x}) \text{(by adjustment formula or backdoor adjustment)}$$
(20)

$$= \sum_{\boldsymbol{x} \in \mathcal{X}} P(\boldsymbol{x}) \sum_{\boldsymbol{\theta} \sim \boldsymbol{\Theta}} P(\boldsymbol{y} | \boldsymbol{x}, \boldsymbol{\theta}) P'(\boldsymbol{\theta} | \boldsymbol{x})$$
(21)

To achieve causal learning during CL, we plug-in the Eq. 21 into Eq. 19, we obtain the following loss function for LLMs:

$$\mathcal{L}_{causal} \approx \frac{1}{M} \sum_{i=1}^{i=M} \mathbb{E}_{\boldsymbol{\epsilon} \sim P(\boldsymbol{\epsilon})}[\log P(y|\mathsf{do}(\mathsf{soft}(\boldsymbol{\Theta})), \boldsymbol{\epsilon})] - \alpha \mathbb{KL}(P'(\boldsymbol{\theta}|\boldsymbol{x}_i), Q(\boldsymbol{\theta}))$$
(22)

B BENCHMARK AND TASK SEQUENCE DETAILS

Benchmark	Order	Task Sequence
	1	dbpedia \rightarrow amazon \rightarrow yahoo \rightarrow ag
Standard CL	2	dbpedia \rightarrow amazon \rightarrow ag \rightarrow yahoo
	3	yahoo \rightarrow amazon \rightarrow ag \rightarrow dbpedia
	4	$mnli \rightarrow cb \rightarrow wic \rightarrow copa \rightarrow qqp \rightarrow boolqa \rightarrow rte \rightarrow imdb \rightarrow$
Long soguenee	4	$Yelp \rightarrow amazon \rightarrow sst-2 \rightarrow dbpedia \rightarrow ag \rightarrow multirc \rightarrow yahoo$
Long sequence	5	$multirc \rightarrow boolqa \rightarrow wic \rightarrow mnli \rightarrow cb \rightarrow copa \rightarrow qqp \rightarrow rte$
	5	\rightarrow imdb \rightarrow sst-2 \rightarrow dbpedia \rightarrow ag \rightarrow Yelp \rightarrow amazon \rightarrow yahoo
	6	$Yelp \rightarrow amazon \rightarrow mnli \rightarrow cb \rightarrow copa \rightarrow qqp \rightarrow rte \rightarrow imdb \rightarrow$
		$\text{sst-2} \rightarrow \text{dbpedia} \rightarrow \text{ag} \rightarrow \text{yahoo} \rightarrow \text{multirc} \rightarrow \text{boolqa} \rightarrow \text{wic}$

Table 6: Task Sequence Orders for CL in LLMs Experiments. Orders 1-3 correspond to the traditional task sequences commonly used in standard continual learning benchmarks (Zhang et al., 2015). Orders 4-6 expand on this by introducing longer sequences, each comprising 15 tasks (Razdaibiedina et al., 2023).

C IMPLEMENTATION DETAILS

• For T5-large model, the per device train batch size is set to be 8, per device eval batch size is set to be 64 and gradient accumulation steps is set to be 4.

810	For Llame? 7P model, the per device train betch size is set to be 4, per device evel betch
811	size is set to be 10 and gradient accumulation steps is set to be 4.
812	size is set to be 10 and gradient accumulation steps is set to be 4.
813	
814	
815	
816	
817	
818	
819	
820	
821	
822	
823	
824	
825	
826	
827	
828	
829	
830	
831	
832	
833	
834	
835	
836	
837	
838	
830	
840	
9/1	
9/0	
9/2	
043	
945	
040	
040	
047	
040	
049	
050	
050	
052	
055	
004	
856	
957	
959	
950	
009	
000	
001	
002	
003	