

Continual Learning of Large Language Models: A Comprehensive Survey

HAIZHOU SHI*, ZIHAO XU, HENGYI WANG, WEIYI QIN, WENYUAN WANG[†], and YIBIN WANG[†], Rutgers University, USA

ZIFENG WANG, Google Cloud AI Research, USA

SAYNA EBRAHIMI, Google DeepMind, USA

HAO WANG*, Rutgers University, USA

The challenge of effectively and efficiently adapting statically pre-trained Large Language Models (LLMs) to ever-evolving data distributions remains predominant. When tailored for specific needs, pre-trained LLMs often suffer from significant performance degradation in previous knowledge domains – a phenomenon known as “*catastrophic forgetting*”. While extensively studied in the Continual Learning (CL) community, this problem presents new challenges in the context of LLMs. In this survey, we provide a comprehensive overview and detailed discussion of the current research progress on LLMs within the context of CL. Besides the introduction of the preliminary knowledge, this survey is structured into four main sections: we first describe an overview of continually learning LLMs, consisting of two directions of continuity: *vertical continuity* (or *vertical continual learning*), i.e., continual adaptation from general to specific capabilities, and *horizontal continuity* (or *horizontal continual learning*), i.e., continual adaptation across time and domains (Section 3). Following vertical continuity, we summarize three stages of learning LLMs in the context of modern CL: Continual Pre-Training (CPT), Domain-Adaptive Pre-training (DAP), and Continual Fine-Tuning (CFT) (Section 4). We then provide an overview of evaluation protocols for continual learning with LLMs, along with currently available data sources (Section 5). Finally, we discuss intriguing questions related to continual learning for LLMs (Section 6). This survey sheds light on the relatively understudied domain of continually pre-training, adapting, and fine-tuning large language models, suggesting the necessity for greater attention from the community. Key areas requiring immediate focus include the development of practical and accessible evaluation benchmarks, along with methodologies specifically designed to counter forgetting and enable knowledge transfer within the evolving landscape of LLM learning paradigms. The full list of papers examined in this survey is available at <https://github.com/Wang-ML-Lab/lm-continual-learning-survey>.

CCS Concepts: • **Computing methodologies** → **Lifelong machine learning**; **Natural language processing**; **Neural networks**.

Additional Key Words and Phrases: Large Language Models, Continual Learning.

ACM Reference Format:

Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin, Wenyuan Wang, Yibin Wang, Zifeng Wang, Sayna Ebrahimi, and Hao Wang. 2025. Continual Learning of Large Language Models: A Comprehensive Survey. *ACM Comput. Surv.* 1, 1 (May 2025), 44 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 INTRODUCTION

Recent advances in large language models (LLMs) have demonstrated considerable potential for achieving artificial general intelligence (AGI) [1, 6, 22, 40, 173, 186, 230, 231]. Researchers have observed that complex abilities such as multi-step reasoning, few-shot in-context learning, and instruction following improve as the scale of parameter size increases [159, 250, 252, 253, 277]. The development of LLMs is impactful and revolutionary, prompting machine learning practitioners to reconsider traditional computational paradigms for once-challenging human-level tasks. However, LLMs

*Correspondence to: Haizhou Shi <haizhou.shi@rutgers.edu> and Hao Wang <hw488@cs.rutgers.edu>.

[†]Work done as visiting students at Rutgers Machine Learning Lab.

are typically trained on static, pre-collected datasets encompassing general domains, leading to gradual performance degradation over time [5, 52, 95, 96, 100, 137] and across different content domains [35, 44, 69, 71, 100, 104, 183, 184, 221]. Additionally, a single pre-trained large model cannot meet every user need and requires further fine-tuning [10, 16, 37, 106, 182, 254, 255, 255, 281, 299, 299]. While one potential solution is re-collecting pre-training data and re-training models with additional specific needs, this approach is prohibitively expensive and impractical in real-world scenarios.

To efficiently adapt LLMs to downstream tasks while minimizing performance degradation on previous knowledge domains, researchers employ the methodology of Continual Learning (CL), also known as *lifelong learning* or *incremental learning* [38, 178, 232, 237]. Inspired by the incremental learning pattern observed in human brains [101, 153, 154, 175], CL trains machine learning models sequentially on a series of tasks with the expectation of maintaining performance across all tasks [57, 58, 113, 124]. Throughout training, models have limited or no access to previous data, posing a challenge in retaining past knowledge as optimization constraints from unseen previous data are absent during current-task learning [124, 135, 213]. This challenge, known as *catastrophic forgetting* [155], has been a central focus in continual learning research since its inception. Over the years, researchers have explored various techniques to mitigate forgetting. These include replay-based methods [30, 207, 213], parameter regularization [4, 113, 196], and model architecture expansion [191, 236]. Together, these techniques have significantly advanced the goal of achieving zero forgetting in continual learning across diverse tasks, model architectures, and learning paradigms.

In the context of training and adapting LLMs sequentially, the significance of CL is undergoing semantic shifts of its own as well. To highlight this ongoing shift, in this paper, we provide a comprehensive overview and detailed discussion of the current research progress on continual LLMs. For the general picture of continual LLMs, we for the first time divide it into two directions of continuity that need to be addressed by practitioners (details in Section 3):

- **Vertical continuity (or vertical continual learning)**, which refers to the ongoing adaptation of LLMs as they transition from large-scale general domains to smaller-scale specific domains, involving shifts in learning objectives and entities of execution. For example, healthcare institutions may develop LLMs tailored to the medical domain while retaining their general reasoning and question answering capabilities for users.
- **Horizontal continuity (or horizontal continual learning)**, which refers to continual adaptation across time and domains, often entails multiple training stages and increased vulnerability to forgetting. For example, social media platforms continuously update LLMs to reflect recent trends, ensuring accurate targeting of downstream services like advertising and recommendations without compromised experience for existing users.

Importantly, separating vertical and horizontal CL transcends mere modification of existing paradigms, like domain-incremental learning, which aligns with horizontal continuity. This distinction offers a robust framework for analyzing complex CL paradigms in language models. For instance, Recyclable Tuning preserves both vertical and horizontal continuity simultaneously [183], and future designs might include zigzagging between horizontal and vertical CL.

In Fig. 1, following *vertical continuity*, we delineate three key stages of LLM learning within modern CL: Continual Pre-Training (CPT), Domain-Adaptive Pre-training (DAP), and Continual Fine-Tuning (CFT) (details in Section 4). In CPT, existing research primarily investigates three types of distributional shifts: temporal, content-level, and language-level. Each presents distinct focuses and challenges. In DAP, CL evaluation and techniques are frequently utilized. However, there is a noticeable lack of diversity in these techniques, considering the maturity of the conventional CL community. In CFT, our focus is on the emerging field of learning LLMs, covering topics such as Continual Instruction Tuning (CIT), Continual Model Refinement (CMR), Continual Model Alignment (CMA), and Continual Multimodal LLMs (CMLLMs). Next, we present a compilation of publicly available evaluation protocols and benchmarks (details in

Section 5). We conclude our survey with a discussion covering emergent properties of continual LLMs, changes in the roles of conventional CL types and memory constraints within the context of continual LLMs, and prospective research directions for this subject (details in Section 6).

In summary, this survey provides a comprehensive review of existing continual learning studies for LLMs, which significantly distinguishes itself from existing literature on related topics [17, 105, 237, 261, 276]. Our survey highlights the underexplored research area of continually developing LLMs, especially in the field of CPT and DAP. We emphasize the needs for increased attention from the community, including the development of practical, accessible, and widely acknowledged evaluation benchmarks. Additionally, methodologies need to be tailored to address forgetting in emerging LLM learning paradigms. We hope this survey can provide a systematic and novel perspective of continual learning in the rapidly-changing field of LLMs and can help the continual learning community contribute to the challenging goals of developing LLMs in a more efficient, reliable, and sustainable manner [8, 25, 95, 219, 268].

2 BACKGROUND AND RELATED WORK

2.1 Large Language Models

Primarily built on the transformer architecture, pre-trained language models (PLMs) have established a universal hidden embedding space through extensive pre-training on large-scale unlabeled text corpora [51, 133, 189]. By scaling parameters to billions or even hundreds of billions and training on massive text datasets [84, 102], PLMs not only demonstrate superior language understanding and generation capabilities but also manifest emergent abilities such as in-context learning, instruction following, and multi-step reasoning [159, 250, 252, 253, 277]. These larger models are commonly referred to as Large Language Models (LLMs). For more detailed introduction, please refer to Appendix A.1.

2.1.1 Pre-training of LLMs. There are two popular pre-training paradigms for LLMs. (1) *Decoder-only models* typically employ auto-regressive language modeling (LM) tasks during pre-training, including the GPT family [1, 22, 173, 186], Gemini family [194, 225], and the open-source Llama family [230, 231]. Specifically, given a sequence of tokens $\mathbf{x} = [x_1, x_2, \dots, x_N]$, LM predicts the next token x_t autoregressively based on all preceding tokens $\mathbf{x}_{<t} = [x_1, x_2, \dots, x_{t-1}]$, and trains the entire network by minimizing the negative log-likelihood $-\sum_{t=1}^N \log P(x_t | \mathbf{x}_{<t})$, where $P(x_1 | \mathbf{x}_{<1}) \triangleq P(x_1)$ is the unconditional probability estimation of the first token. (2) *Encoder-only models*, e.g., BERT [51, 133], use masked language modeling (MLM) as a common pre-training objective. In MLM, for the input sequence \mathbf{x} , a subset of input tokens $m(\mathbf{x})$ are masked and replaced with the special [MASK] token. The pre-training goal is to utilize the unmasked parts $\mathbf{x}_{\setminus m(\mathbf{x})}$ to predict the masked portions $m(\mathbf{x})$. In summary, the overarching goal of MLM is to minimize the negative log-likelihood $-\sum_{\hat{\mathbf{x}} \in m(\mathbf{x})} \log P(\hat{\mathbf{x}} | \mathbf{x}_{\setminus m(\mathbf{x})})$.

2.1.2 Adaptation of LLMs. LLMs are primarily trained to generate linguistically coherent text. However, this training may not align with human values, preferences, or practical needs. Furthermore, the pre-training data can be outdated, leading to knowledge cutoffs or inaccuracies. To address these issues, various computational paradigms such as Instruction Tuning (IT) [288], Model Refinement (MR) [47], and Model Alignment (MA) [174, 187] have been proposed. These approaches adapt LLMs to better meet diverse downstream tasks and user requirements.

Numerous studies show that **Instruction Tuning (IT)** can notably improve LLMs' ability to follow textual instructions [98, 174, 203, 250, 288], leveraging the pre-existing knowledge within LLMs to bridge the gap between general and task-specific performance [251]. Recent works like WizardLM [269] and CodeLM [246] further tailor synthetic data to steer LLMs' behavior through IT. Additionally, IT enhances the interaction between humans and LLMs, providing

a more natural interface and aligning LLM outputs more closely with human expectations and preferences [145]. LLMs make mistakes, such as inaccurate translations or outdated information [47]. Directly fine-tuning the model to correct these mistakes may disrupt its performance on previously learned tasks. To overcome these challenges, **Model Refinement (MR)** is proposed to rectify the model’s errors while preserving its performance on other inputs, with only moderate computing resources [47, 74, 76, 92, 163, 164, 215]. **Model Alignment (MA)** ensures AI systems’ actions and outputs align with human values, ethics, and preferences [174, 187]. MA can be broadly categorized into two types: Reinforcement Learning-based (RL-based) and Supervised Learning-based (SL-based). RL-based approaches [174, 205] are trained to make decisions reinforced by human feedback, using a reward system to guide them towards desirable outcomes. In contrast, SL-based approaches [81, 97, 187] directly train models on datasets of human preferences, aligning their output with demonstrated human values.

2.2 Continual Learning

Humans can accumulate knowledge and skills across tasks without significant performance decline on previous tasks [101, 153, 154, 175]. In contrast, machine learning models, which are typically data-centric, often experience performance degradation on old tasks when trained on new ones, a phenomenon known as “*catastrophic forgetting*.” The challenge of adapting models to a sequence of tasks without forgetting, especially when little to no past data can be preserved, is extensively studied in the continual learning community [38, 178, 232, 237]. Formally, the objective of CL is to find a hypothesis that minimizes risk across all tasks/domains. Consider DIL as an example [112, 213], at t -th learning stage, the ideal training objective $\mathcal{L}(h)$ is defined as

$$\mathcal{L}(h) \triangleq \underbrace{\sum_{i=1}^{t-1} \mathcal{L}_{\mathcal{D}_i}(h)}_{\text{past domains}} + \underbrace{\mathcal{L}_{\mathcal{D}_t}(h)}_{\text{current domain}}, \quad (1)$$

where \mathcal{D}_i denotes the data distribution of the i -th continual learning stage. The objectives for past domains are often challenging to measure or optimize due to the memory constraints (Definition A.3). Therefore, the core of designing CL algorithms lies in identifying a proxy learning objective for the first term without violating the memory constraint. **A more detailed introduction to the formal definition of CL and its techniques can be found in Appendix A.2.**

2.2.1 Types of Continual Learning. To lay the groundwork for subsequent discussions (as illustrated in Table 3 and Section 6.2), we follow the conceptual framework proposed by [112, 232, 237]. There are three primary types of continual learning scenarios: (i) Task-Incremental Learning (TIL), where task indices are available to the model during inference [113, 124]; (ii) Domain-Incremental Learning (DIL), where the model learns a sequence of tasks with the same formulation but without task indices during inference [213]; and (iii) Class-Incremental Learning (CIL), where the model learns new classes of data during training [112, 193].

2.2.2 Techniques of Continual Learning. Existing CL techniques can be roughly categorized into five groups [237]: (i) replay-based, (ii) regularization-based, (iii) architecture-based, (iv) optimization-based, and (v) representation-based. Here, we provide a concise yet comprehensive introduction to the first three categories of continual learning techniques, as they are extensively applied in continual LLMs.

Replay-based methods adopt the relaxed memory constraint by keeping a small buffer of observed data and retraining the model on it when learning new tasks. Although replay-based methods may theoretically lead to loose generalization bounds [213], they are valued for their simplicity, stability, and high performance, even with a small

episodic memory [24, 30, 193, 195]. **Regularization-based methods** adopt a regularization term $\lambda \|\theta - \theta_{t-1}\|_{\Sigma}$ that penalizes large deviation from the history model in the parameter space, where $\|v\|_{\Sigma} = v^{\top} \Sigma v$ is the vector norm evaluated on a positive-semi-definite matrix Σ , and λ is the regularization coefficient, a hyper-parameter introduced to balance the past knowledge retention and current knowledge learning. The matrix Σ introduced is to measure the different level of importance of each parameters and their correlations in retaining the past knowledge. In practice, to reduce computational overhead, diagonal matrices are often designed to encode only the importance of each parameter [4, 113, 197]. **Architecture-based methods**, especially expanding the network architecture dynamically to assimilate new knowledge, is considered the most efficient form of CL [248, 249]. This method primarily tackles adaptation challenges and can achieve zero-forgetting when task IDs are available during inference or can be correctly inferred [71, 256]. However, due to the difficulty of task ID inference, architecture expansion is predominantly utilized in TIL but is scarcely explored in DIL or CIL. In conjunction with pre-trained backbone large models like ViT [54], CoLoR [256] trains various low-rank adaptation (LoRA) [86] modules for different tasks. It estimates and stores prototypes for each task and utilizes the natural clustering ability of the pre-trained model during testing to infer task IDs, selecting the corresponding LoRA component for prediction generation. In the domain of continual LLMs, architecture expansion has resurged in popularity following the rise of parameter-efficient fine-tuning (PEFT) [50, 86, 211], a topic we will delve into shortly [96, 100, 118, 177, 240, 257, 272, 273].

2.2.3 Evaluation Metrics of Continual Learning. There are four evaluation protocols primarily designed for continual learning. **Overall Performance (OP)** [106, 286, 291] calculates the average performance up until the current training stage, measuring the overall ability of a model balancing the performance of each task. As noted in [213], OP corresponds to the primary optimization objective of continual learning, and hence receives the most attention. **Forgetting (F)** represents the largest performance drop observed of each task throughout the training process, averaged over all training stages. It quantifies the negative impact of learning new tasks brought to previously acquired knowledge. Ideally, a robust continual learning framework should achieve **Backward Transfer (BWT)**, where learning new tasks enhances performance on prior tasks. BWT is measured by negating the forgetting, and hence a negative forgetting indicates an improvement in performance on earlier tasks. **Forward Transfer (FWT)** measures the generalization ability of the continual learning algorithms to unseen tasks. It is defined as the difference between the current model's performance evaluated on the future tasks and the randomly initialized model. Refer to Appendix B.1 for more details.

3 CONTINUAL LEARNING MEETS LARGE LANGUAGE MODELS: AN OVERVIEW

Large language models (LLMs) are extensive in various dimensions, including the size of model parameters, pre-training datasets, computational resources, project teams, and development cycles [1, 6, 22, 40, 173, 186, 230, 231]. The substantial scale of LLMs presents notable challenges for development teams, particularly in keeping them updated amidst rapid environmental changes [5, 52, 95, 96, 100]. To illustrate, in 2023, the average daily influx of new tweets exceeds 500 million¹, and training on even a subset of this large volume of data is unaffordable. Recyclable Tuning [183] is the first work to explicitly outline the supplier-consumer structure in the modern LLM production pipeline. On the supplier side, the model is continually pre-trained over a sequence of large-scale unlabeled datasets. After every release of the pre-trained model, the consumer utilizes the stronger and more up-to-date upstream model for downstream tasks. Compared to the upstream supplier, downstream users often lack capacity of collecting and storing large-scale data, maintaining large-scale hardware systems, and training LLMs themselves. In this survey, we extend this framework

¹Source: <https://www.omnicoreagency.com/twitter-statistics>

and further present a comprehensive modern production pipeline encompassing various studies on continual LLM pre-training, adaptation, and deployment (Fig. 1). What sets our framework apart from existing studies [261] is the *incorporation of two directions of continuity: Vertical Continuity and Horizontal Continuity*.

3.1 Vertical Continuity (Vertical Continual Learning)

Definition. Vertical continuity (or vertical continual learning) has long been studied, either implicitly or explicitly, in existing literature. Vertical continuity is characterized by a hierarchical structure encompassing data inclusiveness, task scope, and computational resources. Specifically, the training task transitions gradually from general pre-training to downstream tasks, typically undertaken by distinct entities within the production pipeline [68, 71, 183, 197, 268, 272]. Fig. 1 shows a typical pipeline for vertical continuity in LLMs, i.e., “pre-training” → “domain-adaptive training” → “downstream fine-tuning” [42, 48, 68, 72, 73, 91, 121, 146, 148, 197, 257, 258, 272, 303]:

- **Pre-training.** During the *pre-training* stage, a substantial amount of data from diverse domains is required to develop a general-purpose LLM. This phase demands a sizable research and development team dedicated to training and benchmarking the model, along with considerable computational resources.
- **Domain-Adaptive Pre-training.** Subsequently, downstream institutions may opt for *domain-adaptive pre-training* to tailor the model for specific tasks using domain-specific data unavailable to the upstream supplier.
- **Finetuning.** Finally, the LLM undergoes *fine-tuning* on annotated data for downstream tasks before deployment.

Throughout the process, the unlabeled domain-specific dataset is smaller in scale than the upstream pre-training phase but larger than the final downstream task fine-tuning phase. This pattern extends to computational resources, team size, and other factors. It is important to note that vertical continuity can involve more than three stages [91, 129, 172, 199]. In real-world applications, during domain-adaptive pre-training, additional “layers” can be added to accommodate multiple entities, such as various departments with distinct objectives but operating within the same domain.

Vertical Forgetting. We term the performance degradation (in terms of general knowledge) due to vertical continual learning “*vertical forgetting*”. As shown in Fig. 2, for vertical continual learning, the data distribution of upstream tasks partially covers the downstream, meaning the model might start off at a decent initialization for the subsequent stage of training. Two significant challenges must be addressed to prevent vertical forgetting:

- **Task Heterogeneity.** Stemming from the inherent disparity between the formulation of upstream tasks and downstream tasks, *task heterogeneity* can lead to differences in model structures and training schemes, which has long been recognized as a major hurdle [112, 124, 170, 193, 262]. To mitigate this issue, practitioners often employ methodologies such as freezing shared parameters during downstream phases or reformulating downstream tasks to match the structure of pre-training tasks [118, 177, 240, 257, 272, 273].
- **Inaccessible Upstream Data.** This challenge arises primarily from varying levels of confidentiality across entities undertaking vertical continual learning. Data collected and curated under different protocols may not be accessible to some downstream entities. This scenario is even more challenging than the strict memory constraint presented in conventional CL (Definition A.3), as algorithms for latter case rely on access to previous data at specific points for parameter importance measurement [4, 113] or for replay [24, 30, 195, 213]. To address the challenge of *inaccessible upstream data*, existing methods either use public datasets or generate pseudo-examples to create proxy pre-training datasets [182].

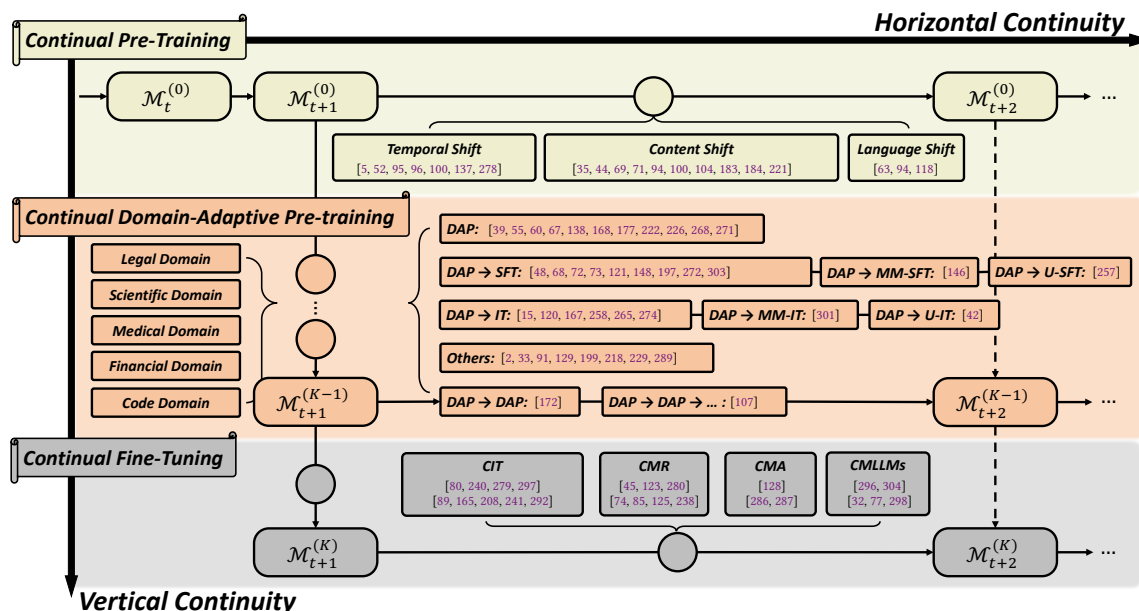


Fig. 1. A high-level overview of the modern pipeline for continually pre-training and fine-tuning LLMs, where two dimensions of continuity are described. **Vertical Continuity (or Vertical Continual Learning):** LLM training can be vertically divided into three stages: (i) Continual Pre-Training (CPT), (ii) Domain-Adaptive Pre-training (DAP), and (iii) Continual Fine-Tuning (CFT). The main focus is the retention of the LLM’s general knowledge (prevention of vertical forgetting). **Horizontal Continuity (or Horizontal Continual Learning):** After the LLMs are deployed, the models are continually updated when a new set of data becomes available. The primary goal is to prevent horizontal forgetting in a long sequence of tasks.

3.2 Horizontal Continuity (Horizontal Continual Learning)

Definition. Horizontal continuity (or horizontal continual learning) refers to continual adaptation across time and domains, a topic extensively explored within the continual learning community. The primary rationale for preserving horizontal continuity lies in the dynamic nature of data distribution over time. To stay updated with these content shifts, an LLM must incrementally learn newly-emerged data. Otherwise, the cost of re-training will become prohibitively expensive and impractical [5, 29, 219, 268]. Empirical evidence has consistently shown that despite their impressive capabilities, LLMs struggle to generalize effectively to future unseen data, particularly in the face of temporal or domain shifts [5, 52, 95, 96]. Additionally, they struggle to retain complete knowledge of past experiences when adapting to new temporal domains, although they do demonstrate a higher level of robustness against catastrophic forgetting [144, 156, 223, 299]. The necessity of employing complex CL algorithms to address challenges in LLMs remains an open question. For instance, during large-scale continual pre-training, large institutions can typically afford the storage costs of retaining all historical data, rendering memory constraints meaningless. Several studies have demonstrated that with full access to historical data, simple sparse replay techniques can effectively mitigate forgetting [62, 181, 208, 223]. In contrast, numerous continual learning studies have showcased superior performance compared to naive solutions, suggesting the importance of continual learning techniques in LLM training [35, 95, 100, 184].

Horizontal Forgetting. We informally define “horizontal forgetting” as the performance degradation on the previous tasks when model is undergoing horizontal continual learning. As illustrated in Fig. 2, horizontal continual learning typically involves training stages of similar scales, with potential distributional overlap among their data. In summary, two main challenges need to be addressed for horizontal continual learning of LLMs:

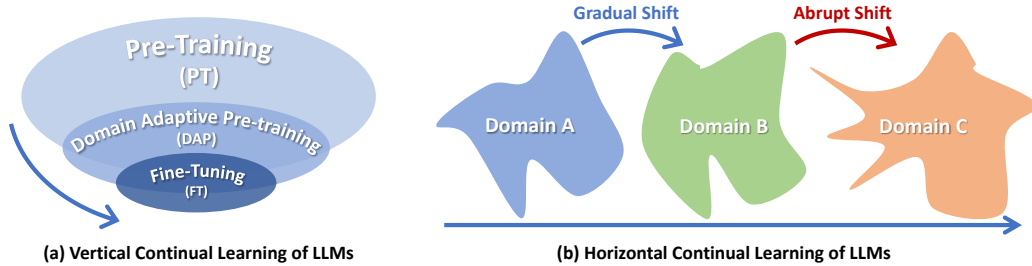


Fig. 2. A diagram showing two different directions of continual learning of LLMs. **(a) Vertical Continual Learning of LLMs:** in this case, the upstream data distribution usually partially covers the subsequent tasks’ data distribution. **(b) Horizontal Continual Learning of LLMs:** No constraints on the data distributions are present on horizontal continual learning. The continual LLMs need to handle the challenge of abrupt distributional shifts and a longer sequence of training.

- **Long Task Sequences.** Horizontal continual learning ideally involves numerous incremental phases, particularly to accommodate temporal shifts in data distribution. A *longer task sequence* entails more update steps of the model, leading to inevitable forgetting of previously learned tasks. To address this challenge, researchers employ established continual learning techniques with stronger constraints, such as continual model ensemble [191].
- **Abrupt Distributional Shift.** In contrast to vertical continuity, where distributional shifts are often predictable, horizontal continual learning does not impose constraints on task properties. Evidence suggests that abrupt changes in task distributions can result in significant horizontal forgetting of the model [204].

4 LEARNING STAGES OF CONTINUAL LARGE LANGUAGE MODELS

Fig. 1 provides an overview of continually learning LLMs. Along the axis of vertical continuity, three main “layers” of modern continual learning emerge. The top layer, Continual Pre-Training (CPT), involves continuous pre-training of LLMs by the supplier on newly-collected data alongside existing data (Section 4.1). The middle layer, Domain-Adaptive Pre-training (DAP), prepares LLMs for domain-specific applications through additional pre-training on domain-specific *unlabeled* data (Section 4.2). The bottom layer, Continual Fine-Tuning (CFT), targets models for final downstream tasks on the consumer side (Section 4.3), where the model needs to be updated after deployment for the specified task.

4.1 Continual Pre-Training (CPT)

4.1.1 CPT: Effectiveness and Efficiency. Before delving into the details of continual pre-training (CPT), it is important to address two fundamental questions: Firstly, regarding *effectiveness*, can CPT enhance performance on downstream tasks beyond that of the initial training on a wide range of data domains? Extensive studies have not only demonstrated the necessity of CPT for improved downstream performance [35, 71, 95, 96, 100, 184], but also shown that when distributional shifts are gradual [95, 278] or somewhat correlated [71], CPT can effectively help model generalize to unseen data. The second question is about *efficiency*: given the large size of an LLM’ parameters and data, both old and new, can we achieve adaptation and knowledge retention in a computationally efficient way? Concerning efficiency, most studies focus on techniques for efficient knowledge retention [95, 96, 100, 118], which significantly overlap with the CL literature addressing catastrophic forgetting [4, 24, 191, 193, 195, 196, 201, 207, 213, 236]. In contrast to prior approaches that fully utilize emergent data, some studies recognize the impracticality of this approach in real production environments. Instead, they concentrate on further improving the efficiency of adaptation. For instance, ELLE [184] employs a function-preserved model expansion to facilitate efficient knowledge growth; [5] and [268] sub-sample training data based on novelty and diversity to enhance training efficiency, achieving superior performance compared

Table 1. **Summary of existing studies on Continual Pre-training of LLMs.** The papers are organized based on their relation to CL: (i) no CL techniques are studied, (ii) CL techniques are studied as solely baselines, and (iii) new CL approaches are proposed. In the table, *Dist. Shift* denotes what type(s) of distributional shifts this particular study considers and is dedicated to solve. In the section of **Continual Learning Tech.**, we mainly categorize three types of continual learning techniques that are studied in the paper: rehearsal (*Rehearsal*), parameter regularization (*Param. Reg.*), and architecture expansion (*Arch. Exp.*). We use “✓”, “✗”, and “♣” to denote “deployed in the proposed method”, “not studied in the paper”, and “studied as a baseline method”, respectively. Note that we do not include naive sequential fine-tuning in this table, as it is universally studied as the important baseline method in all of the papers in the table. The papers with only “♣” [95, 96, 100] means that only existing CL techniques are studied, without proposing new ones, and the papers with only “✗” [63, 69] means that special aspects of fine-tuning are studied, without using CL techniques.

Method	Scenario		Continual Learning Tech.			LLM Arch.	Evaluation	
	Dist. Shift	#Domains	Rehearsal	Param. Reg.	Arch. Exp.		Pre-Training	Downstream
TimeLMs [137]	Temporal	8	✗	✗	✗	RoBERTa	✓	✓
[278]	Content	159	✗	✗	✗	RoBERTa GPT-2	✓	✗
[69]	Content	1	✗	✗	✗	Pythia	✓	✗
[63]	Language	3	✗	✗	✗	GPT	✓	✗
RHO-1 [130]	Other	1	✗	✗	✗	TinyLlama Mistral	✓	✓
[118]	Language	1	✗	P-Freeze♣	Adapter♣ LoRA♣	Llama2	✓	✓
CKL [96]	Temporal	1	Mix-Review♣	P-Freeze♣ RecAdam♣	LoRA♣ K-Adapter♣	T5	✗	✓
LLPT [100]	Temporal	4	ER♣ Logit-KD♣ Rep-KD♣ Contrast-KD♣	oEWC♣	Adapter♣ Layer Exp.♣	RoBERTa	✓	✓
	Content	8	SEED-KD♣					
TemporalWiki [95]	Temporal	5	Mix-Review♣	P-Freeze♣ RecAdam♣	LoRA♣ K-Adapter♣	GPT-2	✓	✓
CPT* [104]	Content	4	DER++♣ KD♣	CPT✓ EWC♣ HAT♣	Adapter♣ DEMIX♣	RoBERTa	✓	✗
ERNIE 2.0 [221]	Content	4	ER✓♣	✗	✗	ERNIE	✗	✓
[5]	Temporal	7	✗	P-Freeze✓	Vocab. Exp.✓	BERT	✗	✓
[44]	Content	5	✗	✗	Vocab. Exp.✓	BERT RoBERTa	✗	✓
DEMIX [71]	Content	8	✗	✗	MoE✓	GPT-3	✓	✓
TempoT5 [52]	Temporal	1	✗	✗	Vocab. Exp.✓ Prompt✓	T5	✗	✓
RecTuning [183]	Content	4	ER✓ KD✓	✗	Adapter✓	RoBERTa	✗	✓
Lifelong-MoE [35]	Content	3	ER✓ KD✓	P-Freeze✓ L2♣	MoE✓	GLaM	✓	✓
ELLE [184]	Content	5	ER✓♣ KD♣	P-Freeze✓	Prompt✓ Layer Exp.✓ Adapter♣	BERT GPT	✓	✓
[94]	Content Language	2	ER✓	✗	✗	GPT-NeoX	✓	✓
CEM [294]	Other	1	ER✓	✗	✗	CuteGPT ChatGLM Qwen-Chat	✗	✓
IR-DRO [36]	Other	1	ER✓	✗	✗	OPT	✗	✓

to full-data training. Though currently underexplored, efficient adaptation in continual pre-training is poised to become significant, given recent findings emphasizing data quality over quantity for LLM generalization [216, 267].

4.1.2 *General Observations on CPT.* Table 1 summarizes the existing studies on continual pre-training (CPT), and here are some key observations we make about CPT.

- **OBS-1: The development of advanced techniques tailored specifically for CPT is at the starting stage and warrants further exploration.** Only about half of the examined papers propose novel techniques for CPT [5, 35, 44, 52, 71, 104, 183, 184, 221], while the remaining half either focus solely on the effects of pure

adaptation without considering CL techniques [63, 69, 137], or conduct empirical studies on the straightforward application of existing CL techniques [95, 96, 100, 118].

- **OBS-2: The diversity of CL techniques incorporated in CPT remains limited.** Most practical implementations of CL techniques for CPT primarily focus on architecture expansion of LLMs [5, 35, 44, 52, 71, 183], with only a few explicitly utilizing replay [35, 183] and parameter regularization [5, 35].
- **OBS-3: There is an apparent gap between the existing studies and the real production environment of CPT.** Except for the recent study [278] which conducts CPT over 159 domains, the longest sequence of pre-training stages explored is 8 [71, 100]. However, this falls short of real-world scenarios where continual pre-training occurs more frequently and persists for months or years. The efficacy of CPT methods in such prolonged scenarios remains uncertain. Additionally, investigating CPT in a task-boundary-free data stream setting is an important avenue for research to be explored in the future as well.

4.1.3 Distributional Shifts in CPT. This survey categorizes distributional shifts of CPT into three main types: (i) *Language Shift*: LLMs sequentially learn different language corpora, e.g., English \rightarrow Chinese [63, 118]. (ii) *Content Shift*: LLMs sequentially learn corpora from different fields, e.g., chemistry \rightarrow biology [35, 44, 69, 71, 100, 183]. (iii) *Temporal Shift*: Distributional shifts occur over time, e.g., news in 2021 \rightarrow news in 2022, with a major focus on timestamp-sensitive knowledge retention and update [5, 52, 95, 96, 100].

Language Shift. [63] focuses on assessing LLMs’ natural ability to learn new languages sequentially. With no explicit CL techniques employed, the study observes consistent positive forward transfer of the knowledge, facilitating new language acquisition regardless of the learning order. Forgetting, on the other hand, emerges as a significant challenge that cannot be mitigated by the increasing size of LLMs. In [118], the degree of forgetting of previously learned language when adapting LLMs to a new language is investigated. Various CL techniques, including parameter freezing, LoRA [86], and (IA)³ [132], are evaluated across multiple dimensions. Preliminary experimental results highlight the non-trivial nature of addressing horizontal forgetting for CPT under the language shift as well.

Content Shift. [278] explores the large-scale CPT over 159 content domains, and shows that CPT on various domains can effectively improve models’ adaptation ability compared to DAP on single domain. Similarly, [69] continues the pre-training phase of Pythia [16] with no complex CL techniques and discovers that learning rate re-warming consistently improves models trained from scratch. Built upon this simple observation, [94] further shows that proper combination of learning rate re-warming and re-decay, and replay of the previous data is sufficient to achieve a comparable performance to full re-training. LLPT [100] establishes a comprehensive training and evaluation protocol for a series of content-level distributional shifts. They assess multiple CL methods and, similar to [63], find consistent forward knowledge transfer, yet horizontal forgetting remains significant. Besides, contrary to the common understanding that experience replay [30] is the most efficient approach to preventing forgetting, the authors find it ineffective in the case of CPT, due to the potential overfitting issue. Recyclable Tuning [183] shows that if the upstream supplier continually pre-trains LLMs, with or without replay, consumer-side efficiency can be boosted by recycling previously learned update components when proper CL techniques are applied.

DEMix [71] incrementally trains and integrates new experts (DEMix layer) for new domains during CPT. To ensure reasonable inference performance during testing when no domain information is available, it proposes a parameter-free probabilistic approach to dynamically estimate a weighted mixture of domains. DEMix’s modularization has been shown to facilitate efficient domain-adaptive pre-training, promote relevant knowledge during inference, and allow for removable components. Lifelong-MoE [35], similar to DEMix [71], incrementally trains domain experts for new domains.

However, Lifelong-MoE differs from DEMix in utilizing a *token-level gating function* to activate multiple experts for intermediate embedding calculation. During training, previous experts' parameters and gating functions remain frozen, and knowledge distillation loss is employed to regulate parameter updates, which thereby makes Lifelong-MoE robust against the issue of horizontal forgetting.

It is noteworthy that some papers draw almost opposite conclusions regarding the significance of CPT for content shifts. For instance, [44] continually pre-trains BERT-based models [51, 133] on five scientific domains and evaluates performance on downstream sentiment analysis. They observe that even the trivial sequential pre-training does not exhibit severe forgetting, prompting reasonable questions about the necessity of CPT.

Temporal Shift. In the context of CPT amid content shifts, Multi-Task Learning (MTL) is often regarded as the upper bound achievable [178, 213, 237]. However, this belief does not fully hold when considering CL under temporal shifts [52, 95, 96], as temporal shifts can introduce conflicting information, posing challenges for LLMs. For instance, the statement “*Lionel Messi plays for team Barcelona*” remains accurate from 2004 to 2021 but becomes false by 2024, as “*Lionel Messi plays for team Inter Miami*” becomes the correct statement.

Hence, as advocated by CKL [96] and TemporalWiki [95], LLMs undergoing continual adaptation to temporal shifts must simultaneously achieve three objectives: (i) retention of old knowledge, (ii) acquisition of new knowledge, and (iii) update of the outdated knowledge. They evaluate the same set of continual learning baseline methods [34, 79, 87, 239], each highlighting distinct aspects of their impact. CKL [96] observes that parameter expansion consistently exhibits robust performance across all experimental conditions. In contrast, replay-based methods struggle to efficiently adapt to new knowledge acquisition and outdated knowledge update, leading to rapid forgetting of newly learned information during training. TemporalWiki [95] constructs a series of temporal corpora and their differential sets from sequential snapshots of Wikipedia, revealing that updating LLMs on these differential sets substantially enhances new knowledge acquisition and updates, requiring significantly less computational resources, and various CL techniques prove effective in mitigating horizontal forgetting during this process. LLPT [100] introduces temporal generalization evaluation for LLMs pre-trained on sequential corpora. Through experiments on a large-scale chronologically-ordered Tweet Stream, the authors demonstrate the superiority of CPT combined with CL techniques to task-specific LMs, in terms of both knowledge acquisition and temporal generalization. Nonetheless, these preliminary experiments do not conclusively determine which specific CL method is more preferable than the others.

Another line of work, Temporal Language Models (TLMs), takes a different approach to address knowledge retention, acquisition, and update under temporal shifts by integrating temporal information into the model [52, 198, 219]. During training, they inject temporal information into training examples as prefixes of prompts, using special tokens [198], explicit year information [52], or syntax-guided structural information [219]. In sequential training experiments conducted by TempoT5 [52], comparison between continually and jointly pre-trained LMs demonstrates that CPT better balances adaptation and forgetting when the replay rate of past data is appropriately set.

Others. CPT as a technique to progressively attain novel knowledge, can be used to refine LLMs' behavior. CEM [294] collects examples where the model's response is incorrect and continually trains the model on these examples, along with a supplemental dataset. RHO-1 [130] proposes Selective Language Modeling (SLM), which employs a reference model to evaluate the perplexity of each token in the training corpus, and continually pre-trains the model on high-perplexity tokens. Similarly, IR-DRO [36] re-trains the model on re-weighted examples from the original pre-training dataset, focusing more on higher-loss sequences.

The significance of addressing temporal shifts through CPT is underscored by several industrial studies. For instance, [5] employs a dynamic vocabulary expansion algorithm and an efficient sub-sampling procedure to conduct CPT on

large-scale emerging tweet data. Conversely, [137] adopts CPT without explicit measures to constrain model updates, releasing a series of BERT-based LMs incrementally trained on new tweet data every three months. Preliminary experimental results demonstrate substantial improvements of continually pre-trained LMs over the base BERT model across downstream tasks. While some studies question the necessity of continually adapting LLMs along the temporal axis for environmental reasons, such as reducing CO₂ emissions [8], the community commonly embraces CPT as a more efficient learning paradigm compared to the traditional “combine-and-retrain” approach.

4.2 Domain-Adaptive Pre-training (DAP)

Background of DAP. Institutions, regardless of size, often possess significant amounts of unlabeled, domain-specific data. This data bridges the gap between general-purpose LLMs trained on diverse corpora and fine-tuned LLMs designed for specific downstream tasks. Leveraging this data as a preparatory stage can facilitate effective adaptation of LLMs to downstream tasks. Such process of “continued/continual/continuous pre-training” [9, 42, 68, 73, 91, 138, 148, 212, 264, 266, 268, 272, 282, 285], “further pre-training” [3, 48, 129, 200, 218], “domain tuning” [197], “knowledge enhancement pre-training” [138], and “knowledge injection training” [258] is unified and termed “*Domain Adaptive Pre-training (DAP)*” [72] for clarity and consistency throughout this survey. In the pioneering work of domain-adaptive pre-training (DAPT) [72], the authors continuously pre-train the language models on a larger domain-specific dataset before fine-tuning them to the downstream tasks, resulting in universally improved performance across various tasks. As the observation above has been validated on multiple domains in parallel, including BioMed, CS, News, and Reviews [72], practitioners commonly accept that employing DAP on additional unlabeled domain-specific data benefits downstream tasks. Consequently, this technique has become widely deployed in many modern LLMs.

Summary of LLMs with DAP. We provide a summary of the existing 41 studies utilizing DAP for LLMs in Table 2. Each entry is characterized by three main features: (i) training process specifications, encompassing the vertical domain for which LLMs are trained, the training pipeline preceding release, and the LLM architecture employed; (ii) adopted continual learning techniques, including rehearsal, parameter regularization, and architecture expansion; and (iii) evaluation metrics for CL, such as backward transfer (forgetting) and forward transfer (adaptation to downstream data).

4.2.1 *General Observation on DAP.* Several key observations emerge regarding the research landscape of DAP (Table 2).

- **OBS-1: DAP predominantly occurs in a single stage.** Continual DAP which involves more than one stage is seldom explored: among all papers listed in Table 2, only one employs two stages of DAP (“PT → DAP → DAP → FT” in Code Llama [199]). It is arguably reasonable to categorize studies that conduct only one stage of DAP and nothing more [9, 39, 67, 138, 168, 177, 218, 226, 268, 271] into CPT rather than DAP. Nevertheless, considering that they aim to adapt a general-purpose LLM to a specific domain, we include them in this section.
- **OBS-2: The notion of interpreting DAP through the lens of CL, whether intentional or not, is widely embraced.** As shown in Table 2, except for the first section (white, 13/41), where papers overlook any potential side effects of DAP leading to vertical forgetting, the remaining sections (all gray, 28/41) either evaluate the potential negative impacts of DAP or proactively employ CL techniques to mitigate the risk of vertical forgetting.
- **OBS-3: Further research of more sophisticated CL techniques for not just DAP, but general vertical continual learning is much needed.** It is supported by the widespread adoption of CL techniques (22/41) for training domain-specific LLMs. However, the diversity of these techniques is limited, with only replay [9, 33, 39, 42, 91, 148, 197, 258, 274, 289] and parameter expansion (LoRA [177, 257, 271, 272]) or Layer/Block expansion [257, 272] utilized. In fact, it appears that individuals may not explicitly recognize that DAP should be

viewed from the perspective of vertical continuity, as they often employ CL techniques unknowingly, e.g., studies deploying replay terming the technique as “data combination” [258] or “data mixing/mixture” [9, 39, 148, 274], without recognizing it as a typical CL solution to vertical continual learning.

4.2.2 Different Domains of DAP. We include work aimed at establishing vertical LLMs across various domains, including legal, medical, financial, scientific, and code. Additionally, we cover other domains such as language and e-commerce.

Legal Domain. In Layer Llama [91], the authors gathered publicly available legal texts from China Courts websites, totaling approximately 10 billion tokens as noted in a GitHub issue. In SaulLM [42], the authors collected the DAP corpus from various jurisdictions in different countries, resulting in a corpus of 30 billion tokens to cover diverse aspects of legal texts. When combined with previously available datasets, the total number of tokens used for legal-domain DAP reaches 94 billion. The substantial volume of DAP data, while offering valuable insights into specific domains, increases the risk of vertical forgetting of the general knowledge due to the large number of update steps involved. To mitigate this issue, SaulLM incorporates general data from Wikipedia, StackExchange, and GitHub into the DAP data, constituting about 2% of the final dataset [42]. Similarly, Lawyer Llama incorporates replaying general-domain data during DAP, but the replay rate is not disclosed [91]. [222] also replays of non-latest business documents during DAP when building a Japanese business-specific LLM.

Medical Domain. Efforts have been made to develop medical specialists by either training an LLM from scratch [66, 143] or fine-tuning publicly-available LLMs to meet specific medical needs [33, 146, 258]. Among these approaches, DAP techniques have been extensively utilized to preserve the communication and instruction-following abilities of a general LLM, preparing it for subsequent medical applications [33, 146, 258]. BioMedGPT [146] is a multi-modal biomedical language model that integrates representations of human language and the language of life (molecules, proteins, cells, genes, etc.). Prior to final multi-modal supervised fine-tuning, the authors initialize the model from Llama2-Chat [231] and conduct DAP using extensive biomedical documents from S2ORC [134], without considering any CL techniques or evaluations. In [68], DAP is performed using Chinese medical encyclopedias and online expert articles, with next-token prediction as the training objective. During DAP, the performance gradually deteriorates on general-domain datasets as the training step increases, but improves on the downstream medical examination tasks [82]. PMC-LLama [258] gathers biomedical papers from S2ORC [134] and medical textbooks for “knowledge injection training.” During this phase, a general language corpus from RedPajama-Data [43] is replayed at a 5% rate within a training batch. However, the paper does not analyze the effectiveness of this operation of mixing in general-domain data for DAP.

To mitigate vertical forgetting, AF Adapter [272] proposes an adapter structure extending the width of Attention layers and FFNs for acquiring domain knowledge and only the adapters are tuned during DAP. Similarly, Hippocrates [2] deploys LoRA during DAP to both have medical-specific knowledge injected and general ability preserved. Me-Llama [265] mixes in about 25% of the general-domain data for DAP on the clinical notes and biomedical articles, which achieves even positive backward transfer on MMLU [82]. HuatuoGPT-II [33] proposes to fuse the DAP into the final SFT, unifying the two stages into one single process. The challenge of such process mainly comes from the data heterogeneity of DAP’s unlabeled corpus. The authors address this challenge by reformulating paragraphs of data into (*instruction, output*) format using existing large language models. They further employ a priority sampling strategy to avoid compromising downstream ability, a pitfall observed in the fixed-rate data mixing strategy [231]. This paper empirically demonstrates the superiority of unified one-stage SFT over two-stage training, questioning the reasonability of the current DAP. On medical-domain data, [197] finds that LMs constrained by CL techniques on source domains

exhibit greater robustness to future domain shifts. Specifically, they identify that parameter regularization techniques like EWC [113], despite slightly higher cost, can facilitate positive forward and backward transfer.

Financial Domain. A gap persists between general-purpose LLMs and existing domain-specific smaller-scale LLMs [7, 259], underscoring the urgent need for more powerful financial-domain experts through the integration of LLMs. Notably, DAP techniques have emerged as crucial tools for tailoring LLMs to the intricacies of the financial domain while mitigating the negative effects of abrupt domain shifts from general to finance [121, 138, 268, 271, 289].

BBT-Fin [138] collects a Chinese financial DAP dataset comprising 80 billion tokens sourced from corporate reports, analyst reports, social media, and financial news. In addition to the conventional masked language modeling (MLM) training objective, BBT-Fin further incorporates triplet masking and span masking techniques during DAP. CFGPT [121] creates CFData, a financial dataset for DAP and SFT, comprising 141 billion tokens. During DAP, CFGPT does not employ CL techniques but utilizes QLoRA [50] for preventing overfitting to downstream data and balancing general response ability and domain-specific ability during SFT. These two methods are typical domain-specific LLMs focusing solely on adaptation to target domains without explicit CL measures or evaluation of vertical forgetting.

In [268], the authors aim to enhance the data efficiency of DAP. When the downstream tasks’ data distribution \mathcal{T} are known, based on the generalization bound [14, 61, 213], the authors propose to sample the subset of DAP data whose distribution \mathcal{D} is similar to the downstream task’s data, i.e., $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}, \mathcal{T})$ is low. When the downstream data distribution is unknown, the authors suggest ensuring *novelty* and *diversity* in the sampled corpus for DAP. This approach significantly enhances DAP efficiency: it utilizes only 10% of the originally collected data yet outperforms models trained on the entire DAP dataset, underscoring the importance of data quality over quantity. WeaverBird [271] introduces an intelligent finance dialogue system, where the encoder is trained on Chinese and English financial documents, alongside expert-annotated financial query-response pairs, using LoRA [87]. Xuanyuan 2.0 [289], akin to HuatuoGPT-II [33], proposes the technique of hybrid-tuning, which fuses the stages of DAP and SFT into one, general-domain data and financial-domain data into one. Notably, the distribution of data in hybrid-tuning is unconventional: financial DAP data comprises only a small portion of 13%. This prompts a pertinent question in line with the investigation on efficient DAP in [268]: Is a large DAP dataset necessary for developing a domain-specific LLM?

Scientific Domain. Vertical scientific LLMs span many subjects [9, 15, 129, 142, 168, 285, 301]. However, among all the studies listed above, only a small fraction of them adopt the technique of DAP. OceanGPT [15] is the first LLM tailored specifically for the ocean domain. It performs DAP on a raw corpus of ocean science literature, prioritizing recent research and historically significant works. K2 [48] pioneers the development of a foundational language model tailored specifically for geoscience. It aggregates geoscience open access literature and Earth science-related Wikipedia pages for DAP. Following this, it undergoes multi-task instruction tuning utilizing LoRA [87] on both a general instruction tuning dataset and the GeoSignal benchmark introduced within the K2 framework. AstroLlama [168] gathers abstracts solely from astronomy papers on arXiv and proceeds pre-training. It observes an improved perplexity on the domain of scholarly astronomy, without providing more quantitative evaluation. MarineGPT [301] is a multi-modal LLM designed specifically for the marine domain. During DAP, MarineGPT incorporates 5 million marine image-text pairs to imbue domain knowledge. This involves training a Q-Former [122] between the frozen visual and text decoder [54, 230].

Another branch of methods proactively integrate in the replay of the general-domain data to mitigate vertical forgetting. GeoGalactica [129] introduces a series of LLMs tailored for geoscience. In the DAP phase, besides the 52-billion-token geoscience corpus, Arxiv papers and Codedata are incorporated, with a mixing ratio of 8:1:1. The authors believe that the inclusion of the Codedata during the model’s pre-training can significantly boost the reasoning ability of the LLMs. Although GeoGalactica pinpoints challenges of DAP, including overfitting, catastrophic forgetting,

Table 2. **Summary of the existing studies that leverage Domain-Adaptive Pre-training of LLMs**, where the papers are organized in four main categories based on whether they (i) adopt the *continual learning techniques* and (ii) perform the evaluation for *backward transfer (forgetting)*. In the column of **Train Proc.** (Training Process), we omit the phase of general Pre-Training. DAP represents Domain-Adaptive Pre-Training; SFT represents Supervised Fine-Tuning; IT represents Instruction Tuning. The prefix G- and D- represent General and Domain-Specific training process [91, 129], and the prefix U- represents them unified [33, 257]. The prefix MM- and LC- represents Multi-Modal and Long-Context training phases [146, 199, 301]. In the column of **Continual Learning Eval.**, we consider two criteria: (i) *Backward Transfer*, i.e., performance degradation on the previous tasks, which is also known as catastrophic forgetting, (ii) *Forward Transfer*, i.e., the performance gained by DAP while transferring the LLMs to the downstream tasks. We use L and Perp. to denote Loss and Perplexity, FT to denote Fine-Tuning, ZS and FS to denote Zero-Shot and Few-Shot Accuracy, HE and LLM to denote the Human Evaluation and LLM Evaluation for generative tasks.

Domain	Method	Train Proc.	LLM Arch.	Continual Learning Tech.			Continual Learning Eval.	
				Rehearsal	Param. Reg.	Arch. Exp.	Backward Transfer	Forward Transfer
Medical	BioMedGPT [146]	DAP → MM-SFT	Llama2	✗	✗	✗	✗	FT
Financial	BBT-Fin [138]	DAP	T5	✗	✗	✗	✗	FT
Financial	CFGPT [121]	DAP → SFT	InternLM	✗	✗	Q-LoRA _(SFT)	✗	HE ¹
Scientific	AstroLlama [168]	DAP	LlaVa	✗	✗	✗	✗	Perp.
Scientific	OceanGPT [15]	DAP → IT	Vicuna Llama2-chat ChatGLM2	✗	✗	LoRA _(IT)	✗	HE
Scientific	K2 [48]	DAP → SFT	Llama	✗	✗	LoRA _(SFT)	✗	Perp. ZS LLM
Scientific	MarineGPT [301]	MM-DAP → MM-IT	Llama	✗	✗	✗	✗	HE
Code	CodeGen [172]	DAP → DAP	CodeGen	✗	✗	✗	✗	Perp. ZS
Code	Comment-Aug [218]	IT → DAP	Llama2 Code Llama InternLM2	✗	✗	✗	✗	ZS
EventTemporal	EcoNet [73] ¹	DAP → FT	BERT RoBERTa	✗	✗	✗	✗	FT
CommonSense	CALM [303]	DAP → FT	T5	✗	✗	✗	✗	FT
Multi-Domain	BLADE [120]	DAP → IT	BLOOMZ	✗	✗	✗	✗	ZS
Scientific	ClimateGPT [229]	DAP → IT → RAG	Llama2	✗	✗	✗	✗	FS Ret.
Medical	[68]	DAP → FT	Llama2	✗	✗	✗	FS FT	FS FT
Financial	[268]	DAP	Pythia	✗	✗	✗	L FS	L FS
Scientific	GeoGalactica [129]	DAP → G-SFT → D-SFT	GAL	✗	✗	✗	ZS	Perp. ZS LLM
Code	StarCoder [226]	DAP	StarCoder	✗	✗	✗	Perp. ZS FS	Perp. ZS FS
Code	DeepSeek-Coder [67]	DAP	DeepSeek-LLM	✗	✗	✗	ZS FS	ZS
Multi-Domain	DAPT [72]	DAP → FT	RoBERTa	✗	✗	✗	Loss	L FT
Financial	WeaverBird [271]	DAP	GLM2	✗	✗	LoRA	✗	HE
Code	IRCoder [177]	DAP	StarCoder DeepSeek-Coder Code Llama	✗	✗	LoRA	✗	ZS
Code	Code Llama [199]	DAP → LC-FT → IT DAP → DAP → LC-FT	Llama2	Replay	✗	✗	✗	Perp. ZS
Legal	SauLLM [42]	DAP → U-IT	Mistral	Replay	✗	✗	✗	Perp. ZS
Medical	PMC-Llama [258]	DAP → IT	Llama	Replay	✗	✗	✗	ZS FT
Scientific	Llama [9]	DAP	Code Llama	Replay	✗	✗	✗	Perp. FS
Multi-Domain	DAS [107]	[DAP] _n	RoBERTa	DER++ [♣]	EWC [♣] HAT [♣] Soft-Masking	Adapter [♣] DEMix [♣]	✗	FT
Medical	Hippocrates [2]	DAP → IT → MA	Llama2 Mistral	✗	✗	LoRA	✗	ZS FS
Language	Sailor [55]	DAP	Qwen1.5	Replay	✗	✗	✗	ZS
Code & Math	Llama Pro [257]	DAP → U-SFT	Llama2	✗	✗	Block Exp. LoRA [♣]	ZS FS	Perp. ZS FS
Medical	AF Adapter [272]	DAP → FT	RoBERTa	✗	✗	Layer Exp. LoRA [♣]	Acc.	L FT
Medical	[197]	DAP → FT	BERT RoBERTa DistilBERT	Replay [♣] GEM [♣]	L2 Reg. [♣] EWC [♣]	✗	L FT	L FT
Medical	HuatuoGPT-II [33]	DAP + U-SFT	Baichuan2	Replay	✗	✗	ZS	ZS HE
Financial	XuanYuan 2.0 [289]	DAP + SFT	BLOOM	Replay	✗	✗	HE	HE
Scientific	PLlama [274]	DAP → IT	GAL	Replay	✗	✗	L	L ZS
E-Commerce	EcomGPT-CT [148]	DAP → SFT	BLOOM	Replay	✗	✗	ZS FS	ZS FS
Legal	Layer Llama [91]	DAP → G-IT → D-IT	Llama	Replay	✗	✗	ZS	ZS
Multi-Domain	AdaptLLM [39]	DAP	Llama	Replay	✗	✗	ZS	ZS FT
Language	Swallow [60]	DAP	Llama2	Replay	✗	✗	FS	FS
Financial	[222]	DAP	Llama2	Replay	✗	✗	Loss ZS	Loss ZS FS RAG
Medical	Me-Llama [265]	DAP → IT	Llama2	Replay	✗	✗	ZS FS	ZS FS FT
Language	Aurora-M [167]	DAP → IT	StarCoder	Replay	✗	✗	ZS	ZS FS HE

maintaining the training stability, and convergence speed, it does not further provide empirical evidence supporting the inclusion of the Codedata, or deploying specific measures to address the challenges proposed above. Llemma [9] focuses on mathematics, initialized from Code Llama [199], and undergoes DAP on a blend of the 55-billion-token mathematical pre-training dataset and general domain data at the ratio of 19:1. In contrast, PLlama [274], designed for plant science, mixes domain-specific and general-domain data at the ratio of 9:1.

Code Domain. The development of LLMs for automatic code filling, debugging, and generation holds significant practical importance [166, 220]. These advancements cover various frameworks, including encoder-only [166], encoder-decoder [242, 245], and decoder-only [67, 172, 227]. There is a growing trend towards decoder-only architectures [220], leveraging models pre-trained on general natural language like Llama [230, 231]. Consequently, there is a shift in the training objective from utilizing code structures to simpler tasks like next token prediction and infilling.

From the perspective of CL, the code domain presents unique advantages and challenges for DAP, compared to other domains. On one hand, its hierarchical structure (*general domain corpus* \rightarrow *multi-language code* \rightarrow *specific programming language*) provides an ideal training pipeline for DAPs [199], offering potential for more efficient training strategies. On the other hand, programming languages adhere to strict grammars, unlike the fuzzy and context-dependent natural language. Consequently, language models should ideally leverage these structures through tailored designs, and adopting the same training objectives as for natural languages may yield sub-optimal results. Therefore, many existing studies omit DAP [147, 242, 245]. In the following section, we will introduce existing code LLMs that employ DAP before the final downstream tasks, discussing both their common attributes and unique characteristics.

Representing a series of notable works that focus solely on adaptation to target domains, CodeGen [172] comprises a suite of LLMs designed for natural language (CodeGen-NL), multi-lingual programming languages (CodeGen-Multi), and mono-lingual programming languages (CodeGen-Mono). These models are trained sequentially, with each subsequent model initialized from the previous one trained on more general-domain data. Comment-Aug [218] addresses the challenge of aligning programming languages with natural languages (PL-NL alignment) by performing DAP on the code augmented with generated additional comments. StarCoder [226] introduces two models: StarCoderBase and StarCoder. StarCoderBase is initially trained on a mixed dataset comprising various programming languages without significant reweighting on the data. Subsequently, StarCoderBase undergoes further fine-tuning on additional 35 billion tokens of Python code, resulting in the development of StarCoder. DeepSeek-Coder-v1.5 [67] originates from DeepSeek-LLM [224] and undergoes pre-training on 2 trillion tokens, comprising 87% source code, 10% English code-related natural language, and 3% Chinese natural language corpus. Initialization from a general-domain LLM results in improved performance across various tasks, including natural language and mathematical reasoning, with minimal performance degradation on coding tasks, which underscores the efficacy of DAP.

As the only work that utilizes the general data replay to mitigate vertical forgetting in the code domain, Code Llama [199] introduces a sophisticated training framework tailored for various coding tasks and model sizes. Initialized from Llama 2 weights, these models undergo DAP on a dataset composed of deduplicated public code, discussions about code, and a subset of natural language data. This mix of natural language data serves as a form of pseudo-replay to maintain the models' proficiency in understanding natural language. Besides replay, architecture expansion has proven effective in acquiring robust coding abilities and preventing vertical forgetting simultaneously. IRCoder [177] utilizes compiler intermediate representations to enhance the multilingual transferability of Code LLMs. By conducting DAP on code grounded in intermediate representations with LoRA [86], IRCoder achieves superior multilingual programming instruction following, enhanced multilingual code understanding, and increased robustness to prompt perturbations. Llama Pro [257] undergoes DAP on a combination of code and math data. It expands the original Llama2 architecture

by dynamically adding multiple identity copies of the transformer blocks. These added blocks initially preserves the original functionality, and will be tuned for DAP. The proposed expansion method is shown to be more resilient against vertical forgetting compared to other parameter-efficient tuning methods like LoRA.

The three aforementioned studies highlight the importance of DAP for code LLMs. However, it is crucial to note that the problem definition and conventional architectures of existing Code LLMs may present challenges of compatibility for DAP deployment, and need to be addressed in the future.

Other Domains. ECONET [73] enhances the model’s ability to reason about event temporal relations through a dedicated DAP phase. Temporal and event indicators are masked out, and a contrastive loss is applied to the recovered masked tokens. Results demonstrate that incorporating this DAP stage significantly improves performance on final tasks compared to direct fine-tuning. Concept-Aware Language Model (CALM) [303] introduces a data-efficient DAP approach for enhancing the concept-centric commonsense reasoning ability of LLMs. It incorporates both generative and discriminative commonsense reasoning tasks specifically tailored for concept-centric reasoning tasks. Consequently, even a small number of data examples for DAP can lead to notable improvements for downstream tasks.

Aurora-M [167] and Swallow [60] adopt the simple replay strategy that mixes in a small portion of general data during DAP for their multi-lingual ability. Furthermore, Sailor [55] studies the optimal strategy of data mixing for DAP, balancing the general knowledge and capacity of different languages. EcomGPT-CT [148] employs a data mixing strategy for DAP which transforms semi-structured E-commerce data into a set of nodes and edges, samples a cluster of nodes, and then extracts and concatenates them into a training example. It combines the general-domain corpus with E-commerce data at a ratio of 2:1, which is significantly lower than the common setting adopted by other works.

Notably, there are some papers studying other effective ways of DAP. AdaptLLM [39] transforms raw corpora into (*raw text, question, answer*) format, creating intrinsic reading comprehension tasks. AdaptLLM demonstrates superior domain-specific knowledge adaptation and minimal vertical forgetting, thereby challenging the data efficiency of conventional DAP. Tag-LLM [212] re-purposes the general-domain LLM into domain-specific one by multi-stage training of domain tags and function tags, without modifying the base LLM’s weights and thereby mitigates forgetting.

4.3 Continual Fine-Tuning (CFT)

Background of Continual Fine-Tuning (CFT). Continual Fine-Tuning (CFT) lies at the bottom layer of the vertical continuity, where models are trained on successive homogeneous tasks drawn from an evolving data distribution. As the service-oriented layer of LLM, it does not require consideration of further adaptation to another downstream tasks, simplifying optimization objectives to a great extent: better adaptation and less forgetting². In the era of LLMs, new computational paradigms in CFT have emerged and attracted significant attention within the research community. These topics include (i) Continual Instruction Tuning (CIT) [292], (ii) Continual Model Refinement (CMR) [74], (iii) Continual Model Alignment (CMA) [128, 287], and (iv) Continual Learning for Multimodal Language Models (CMLLMs) [77, 171]. We summarize existing studies on CFT in Table 3, categorizing studies into sub-categories as listed above. The table includes details on incremental learning types (X-IL), LLM architecture, and employed CL techniques and evaluation metrics. After discussing general observations on CFT in Section 4.3.1, we will delve into each sub-category in detail.

4.3.1 General Observations on CFT. Examining the landscape of continual learning in the context of LLMs, and combined with the results shown in Table 3, we make several key observations about CFT.

²We direct interested readers to additional survey literature on the topic of general CFT [17, 105].

Table 3. **Summary of the existing studies on Continual Fine-Tuning LLMs**, where the papers are organized in five main categories based on what downstream tasks they are designed to tackle, including (i) General Continual Fine-Tuning (CFT); (ii) Continual Instruction Tuning (CIT); (iii) Continual Model Refinement (CMR); (iv) Continual Model Alignment (CMA); (v) Continual Multimodal LLMs (CMLLMs), which is shown in the column of **CFT Type**. The column of **X-IL** shows what continual learning paradigm the study includes [232], where *TIL* represents task-incremental learning, meaning task ID/information is provided during inference; *DIL* represents domain-incremental learning, meaning the tasks are defined in the same format, and no task ID/information is available during inference; *CIL* represents class-incremental learning, meaning the task ID needs to be further inferred when testing.

CFT Type	Method	X-IL	LLM Arch.	Continual Learning Tech.				Continual Learning Eval.		
				Rehearsal	Param. Reg.	Arch. Exp.	Others	Avg. Acc.	Bwd. Trans.	Fwd. Trans.
General	CTR [106]	DIL CIL	BERT	✗	✗	Adapter	✗	✓	✓	✓
	[223]	TIL	BERT	S-Replay	✗	✗	✗	♣	♣	♣
	CIRCLE [281]	DIL	T5	Replay	EWC	Prompt	✗	✓	✓	✓
	ConPET [217]	DIL	Llama	Replay	✗	LoRA	✗	✓	✓	✓
	[10]	DIL CIL	BERT	✗	✗	✗	G-Prompt	✓	✓	✗
	[144]	TIL	DistilBERT ALBERT RoBERTa	ER DER LwF	✗	✗	✗	♣	♣	✗
	SEQ* [299]	TIL CIL	Pythia BERT GPT2	✗	P-Freeze	✗	Tricks for Classifiers	✗	✓	✗
	LFPT5 [182]	DIL	T5	P-Replay	✗	✗	✗	✓	✓	✗
	[254]	DIL	RoBERTa GPT2	Replay	EWC SI RWalk	✗	✗	✓	✓	✗
	LR ADJUST [255]	DIL	XLm-R	✗	✗	✗	LR Scheduling	✓	✓	✓
	C3 [37]	TIL	T5	KD	✗	Prompt Tuning	✗	✓	✓	✗
	CTO [208]	TIL	T0	S-Replay	✗	✗	✗	✓	✓	✓
	RCL [241]	TIL	LLaMA Vicuna Baichuan	Replay	✗	✗	✗	✓	✓	✓
	DynaInst [165]	TIL	BART	Replay	✗	✗	✗	✓	✓	✓
	CITB [292]	TIL	T5	Replay AGEM	L2 EWC	AdapterCL	✗	✓	✓	✓
CIT	SSR [89]	TIL	LLaMA Alpaca	RandSel KMeansSel	✗	✗	✗	✓	✓	✓
	KPIG [80]	DIL TIL	LLaMA Baichuan	DynaInst PCLL DCL	L2 EWC	DARE LM-Cocktail	KPIG	✓	✓	✓
	ConTinTin [279]	TIL	BART	Replay	✗	✗	InstructionSpeak	✓	✓	✓
	O-LoRA [240]	TIL	LLaMA Alpaca	✗	✗	O-LoRA	✗	✓	✓	✓
	SAPT [297]	TIL	T5 LLaMA	✗	✗	✗	SAPT	✓	✓	✓
	InsCL [243]	TIL	LLaMA	Replay	✗	✗	InsCL	✓	✓	✓
	CMR	CMR [125]	DIL	BART	ER MIR MLR	L2 EWC	✗	✗	✓	✓
GRACE [74]	DIL	T5 BERT GPT2	✗	✗	Adapter	✗	✓	✓	✗	
WiKE [85]	DIL	GPT2 GPT-J	✗	✗	Adaptor	✗	✓	✓	✓	
Larimar [45]	DIL	BERT GPT-J	✗	✗	✗	Kanerva Memory	✓	✓	✓	
MELO [280]	DIL	BERT GPT2 T5	✗	✗	LoRA	✗	✓	✓	✓	
CME [123]	DIL	BERT	Replay	✗	✗	Inner-Prod. Reg.	✓	✓	✓	
WISE [238]	DIL	GPT-J Llama2 Mistral	✗	✗	✗	Side Memory	✓	✓	✓	
CMA	COPF [286]	TIL DIL	Llama	Replay	Function Reg.	Prompt	✗	✓	✗	✓
	AMA [128]	DIL	OpenLLaMA Mistral	Replay	L1 L2	LoRA	Adaptive Model Avg.	♣	♣	♣
	CPPO [287]	TIL	GPT2	✗	Weighting	Prompt	✗	✓	✓	✓
CMLLMs	EProj [77]	TIL	InstructBLIP	✗	TSIR	Projector Exp.	✗	✓	✗	✓
	Fwd-Prompt [298]	TIL	InstructBLIP BLIP2	✗	✗	Projector Exp.	✗	✓	✓	✓
	CoIN [32]	TIL	LLaVA	✗	✗	MoE LoRA	✗	✓	✗	✓
	Model Tailor [304]	TIL	InstructBLIP LLaVA	✗	Model Tailor	✗	✗	✓	✓	✓
	RebQ [296]	TIL	VILT	✗	✗	Prompt Tuning	✗	✓	✗	✓

- **OBS-1: There has been a noticeable transition in focus from CIL to TIL and DIL.** It has been a longstanding common sense in the CL community that CIL, as it requires the model to predict the context label and within-context label at the same time [112, 232, 237], is the most challenging CL scenario and hence receives most of the attention from the community. However, among all 35 papers presented in Table 3, only 3 papers study CFT of CIL. The transition of the research focus demonstrates the importance of TIL and DIL in the real-world applications of continual LLMs. More detailed discussion of this transition is included in Section 6.2.
- **OBS-2: In CFT, CL techniques enjoy broader adoption and explicit exploration compared to CPT and DAP.** In Table 3, all 35 papers explicitly deploy the CL techniques, 50% of which develop new techniques that cannot be easily interpreted as trivial combination of existing classic CL techniques, e.g., shared attentive learning framework in SAPT [297], external memory deployed in Larimar [45], and adaptive model averaging method to achieve Pareto-optimal in AMA [128], etc. This underscores the recognition of continual learning as a pivotal component in the development of resilient and adaptive LLMs.

4.3.2 General Continual Fine-Tuning (General CFT). Researchers have long investigated the phenomenon of forgetting resilience in pre-trained LLMs when fine-tuned for downstream tasks [106, 144, 156, 223, 299], despite some discover the opposite [144]. Although the pre-trained weights initially position the model in a flat-loss basin, aiding adaptation to future tasks without severely impacting previous ones [156], zero or near-zero forgetting is only observed at the representation level. This implies that while the model retains its ability to distinguish between task-specific representations, it may still forget specific task details [144, 223, 260, 299]. Therefore, additional measures are necessary when deploying these models in real-world applications [10, 37, 106, 182, 254, 281].

Many studies advance beyond naive sequential fine-tuning, leveraging the inherent anti-forgetting nature of LLMs while avoiding the adoption of overly complex CL techniques [255, 299]. For instance, LR ADJUST [255] proposes a straightforward yet effective method of dynamically adjusting the learning rate to mitigate the overwriting of knowledge from new languages onto old ones. Building on the innate anti-forgetting ability of large language models like Pythia [16], SEQ* [299] introduces several strategies for fine-tuning LLMs on a sequence of downstream classification tasks, such as freezing the LLM and old classifier’s parameters after warm-up, and pre-allocating future classifiers, etc.

Given the minimal forgetting observed at the representation level in CL, some studies aim to tackle the misalignment between the representation space and the decision-making layers by introducing representation-level constraints during CFT. NeiAttn [10] exemplifies this approach by formulating classification tasks as masked language modeling and proposing a neighboring attention mechanism to counteract negative representation drift.

Another line of approaches refines the input/output format and network architectures of pre-trained LLMs to be better suited for CFT. For instance, CTR [106] incorporates two CL-plugin modules, i.e., a task-specific module (TSM) for acquiring task-specific knowledge and a knowledge-sharing module (KSM) for selectively transferring previously learned similar knowledge. CIRCLE [281] manually designs diverse prompt templates for various types of buggy code, unifying them as the cloze task and employs difficulty-based replay to enhance continual program repair. LFPT5 [182] addresses lifelong few-shot language learning by consolidating sequence labeling, text classification, and text generation into a text-to-text generation task. It undergoes prompt tuning on generated pseudo-examples from previous domains when adapting to new tasks. In [291], the authors propose a method for adaptively adding compositional adapters during continual sequence generation tasks. Before training on new domains, a decision stage determines which trained module can be reused. During training, this module also regenerates examples of the past for replay. C3 [37] merges PEFT and in-context learning (ICL) in a teacher-student framework. The teacher model undergoes in-context tuning focused solely on the current domain, while the student model, together with tunable prompts, minimizes the KL-divergence between the output distribution and the ground truth and teacher model simultaneously.

4.3.3 Continual Instruction Tuning (CIT). When the instruction tuning data comes in as a stream, forgetting of the previously learned instructions should be addressed. CT0 [208] represents the inaugural study on Continual Instruction Tuning (CIT) of LLMs, applying the replay method on the base T0 model throughout the process. Many subsequent studies focus on enhancing the replay method used during CIT. For instance, [80] improve replay efficiency by computing Key-Part Information Gain (KPIG) on masked parts to dynamically select replay data, addressing the “half-listening” issue in instruction following. Similarly, SSR [89] uses the LLM to generate synthetic instances for replay, achieving superior or comparable performance to traditional methods at a lower cost.

Other approaches introduce multiple CL techniques during CIT. DynaInst [165] merges parameter regularization with dynamic replay, selectively storing and replaying instances and tasks to enhance outcomes. InstructionSpeak [279] employs negative training and replay instructions to improve both forward transfer and backward transfer. Some

methods incorporate PEFT. Orthogonal Low-Rank Adaptation (O-LoRA) learns new tasks within an orthogonal subspace while preserving LoRA parameters for previous tasks [240] to minimize the interference among different tasks. Shared Attention Framework (SAPT) combines a PET block with a selection module via a Shared Attentive Learning & Selection module, tackling catastrophic forgetting and knowledge transfer concurrently [297]. While regularization-based and architectural-based methods require additional parameter storage and GPU memory, together with replay-based methods they remain for CIT due to the simplicity and effectiveness [243].

4.3.4 Continual Model Refinement (CMR). The concept of model editing was initially explored in [215], which introduced a “reliability-locality-efficiency” principle and proposed a gradient descent editor to address it efficiently. Subsequent research, such as [47] and [163], extended this principle to edit factual knowledge in BERT-based language models and larger models like GPT-J-6B [235] and T5-XXL [189], respectively, using gradient decomposition. These approaches typically update a subset of model parameters to alter the labels of specific inputs. Additionally, memory-based models, as discussed in [164] and [74], incorporate editing through retrieval mechanisms.

Continual Model Refinement (CMR) extends model refinement horizontally, presenting updated sample pairs $(\mathbf{x}_e, y_e, \hat{y}_e)_{N}^{e=1}$ sequentially as a stream. [125] initially introduces this idea, evaluating various CL methods with a dynamic sampling algorithm. Many CMR methods employ a retrieval mechanism. For instance, [74] uses hidden activations of the language model as a “key” to activate updated parameters only when input x_0 resembles updated sample pairs; [280] improves this approach’s efficiency by integrating LoRA [86]; [45] augments the LLM with an external episodic memory, modeling CMR as an ongoing memory refresh. Meanwhile, some methods focus solely on updating a subset of model parameters. For example, [85] addresses the issue of “toxicity buildup and flash” in single-editing methods like ROME [157], adapting it to the CL context with a knowledge-aware layer selection algorithm. WISE [238] addresses the “impossible triangle” of reliability, locality, and generalization in existing lifelong model refinement methods. It introduces a side memory system that enables knowledge sharding and merging, successfully achieving all three objectives simultaneously.

While all these works pioneer research in CMR, the exploration of CMR of LLMs remains open. [75] highlights a potential problem: the location for storing the fact may not coincide with the best place for editing it. This challenges the classical “locate and edit” paradigm used by several existing methods [157, 158], and could become a significant concern for CMR [85]. Other questions, including whether such problem setting fits LLMs and whether more memory/computationally efficient methods of CMR could be developed for LLMs, are yet to be answered.

4.3.5 Continual Model Alignment (CMA). When LLMs undergo the phase of MA, vertical forgetting of previous knowledge usually occurs. In [128], the authors refer to this phenomenon of catastrophic forgetting induced caused by MA as the “Alignment Tax.” Notably, even a single stage of MA can diminish the model’s performance capabilities, as it restricts the model’s responses to a narrower subset of the training distribution.

Continual Model Alignment (CMA) aims to continuously refine LLMs to align with evolving human values, ethics, and data. The static nature of LLM training on historical data sets can lead to discrepancies between the models’ outputs and current factual accuracies, societal norms, and standards, making CMA a crucial process for maintaining their adaptability and alignment with contemporary contexts. Likewise, there are two types of CMA frameworks: RL-based and SL-based. In the realm of RL-based CMA, two significant contributions have been noted. [128] identifies the conflicts between the existing CL techniques and RLHF, and proposes Adaptive Model Averaging (AMA), adaptively finding appropriate ratios for the combination of model layers to gain maximal rewards with minimal tax; Continual Proximal Policy Optimization (CPPO) [287] proposes a weighting strategy for different examples deciding its usage of policy

enhancement or knowledge retention, mitigating the alignment tax over time. For SL-based CMA, Continual Optimal Policy Fitting (COPF) [286] presents a solution adapted from the Direct Policy Optimization (DPO) [188], solving its potential risks of sub-optimal policy fitting and over-optimization in the context of CMA.

4.3.6 Continual Multimodal Large Language Models (CMLLMs). Continually training multi-modal models like CLIP [185] has been long studied [171, 300], while the problem of continually training MLLMs still remains underexplored. Several existing studies have investigated the causes of catastrophic forgetting when continually training MLLMs. [298] performs singular value decomposition on input embeddings, revealing a significant disparity among different input embeddings. This discrepancy causes the model to learn irrelevant information for previously trained tasks, resulting in catastrophic forgetting and negative forward transfer. [284] observes that minority collapse may lead to catastrophic forgetting, when the imbalance ratio between majority and minority classes approaches infinity during fine-tuning. It further identifies hallucination as a contributing factor to performance degradation in MLLMs.

Continual Fine-Tuning MLLMs. In contrast to traditional continual learning methods that involve full-model fine-tuning for new tasks, continual fine-tuning for MLLMs focuses on refining specific layers when adapting to new tasks [32, 77, 284, 298, 304]. Given the strong capabilities of pre-trained models, training specific layers suffices, and can simultaneously reduce computational demands. [296] additionally considers a continual learning scenario, Continual Missing Modality Learning (CMML), where different modalities are emerging throughout the incremental learning stages. All the aforementioned studies collectively indicate that MLLMs still suffer from catastrophic forgetting, which manifests in two ways: along the direction of *vertical continuity*, a performance decline on pre-trained tasks following fine-tuning for downstream tasks; and along the axis of *horizontal continuity*, a performance degrade on previously fine-tuned tasks after fine-tuning for new tasks. [298] also observes negative forward transfer, where the performance of unseen tasks degrades when learning new tasks, indicating a decline in model generalization capability.

While traditional CL methods are applicable, some may not yield optimal results, as evidenced by various experiments [77, 298]. For instance, [77] observes a consistent efficacy of replay-based and model expansion strategies across diverse scenarios of continual fine-tuning MLLMs, but regularization-based methods only perform well on models that have been jointly instruction-tuned on multiple tasks. Other works seek to develop ad-hoc solutions for continual learning MLLMs. [77] proposes EProj to expand the projection layer in MLLMs for each new task and utilizes task-similarity-informed regularization (TIR) to enhance performance. [298] introduces Fwd-Prompt, a prompt tuning method that projects prompt gradient to both the residual space and the pre-trained subspace to minimize the interference between tasks and reuse pre-trained knowledge respectively, fostering positive forward transfer without relying on previous samples. [304] focuses on the forgetting of the pre-trained MLLMs after fine-tuned on specific tasks and proposes model tailor to compensate the selected subset that are critical for enhancing target task performance. [296] presents a novel method named Reconstruct before Query (RebQ), leveraging the multi-modal knowledge from a pre-trained model to reconstruct the absent information for the missing modality. Recently, MoE (Mixture-of-Experts) framework has gained attention which resembles the architecture-based methods in CL. It provides the model with the ability to learn different intentions from distinct experts, e.g., [32] first introduces MoELoRA to fine-tune LLaVA, effectively mitigate the catastrophic forgetting of MLLMs in CoIN and the results demonstrate the effectiveness.

5 EVALUATION PROTOCOLS AND DATASETS

Continual LLMs' Evaluation Protocols. LLanguage Model Analysis (LAMA) is an evaluation framework designed to *probe the world knowledge* embedded in language models [179]. LAMA converts each world fact into a cloze statement,

which is then input into the language models to predict the correct answer. It has been extensively utilized in work on CPT under the temporal shifts [95, 96]. FUAR (Forgotten / (Updated + Acquired) Ratio) is proposed for CPT to address the OP’s drawback of not able to accurately reflect the model’s behavior. A FUAR value of 1 represents an equal trade-off between the knowledge forgetting and knowledge learning, while a FUAR less than 1 suggests high learning efficacy. In TRACE [241], the authors propose a set of “**X-Delta**” metrics for continual instruction tuning, quantifying the forward transfer on specific abilities of LLMs, which is a straightforward extension of FWT. Specifically, the authors construct three sets of evaluation tasks to benchmark the ability of LLMs, including *general ability*, *instruction following*, and *safety*. For more detailed introduction to these evaluation protocols, please refer to Appendix B.2.

Datasets. In this section, we provide a comprehensive review of the datasets available for benchmarking continual LLMs, as illustrated in Table 4. We provide information about these datasets’ types, what distributional shifts and semantic domains they include, and their sources and applications. We intentionally exclude datasets used for domain-adaptive pre-training LLMs in vertical domains such as legal, medical, and financial, unless they are specifically designed for continual domain-adaptive pre-training. Furthermore, we omit datasets used in general continual fine-tuning, as they have already been extensively studied in existing works [17, 105]. For details, please refer to Appendix B.3.

6 DISCUSSION

6.1 Intriguing Properties Emergent in Continual LLMs

Beyond the well-established resilience of pre-trained large language models (LLMs) against catastrophic forgetting compared to downstream-specific models [106, 144, 156, 223, 299], there is a notable lack of exploration into other intriguing properties of LLMs when trained continually. In [275], it is observed that when fine-tuned sequentially and cyclically on a series of documents, large models exhibit a phenomenon known as “*anticipatory recovering*.” This refers to the LLMs’ ability to recover forgotten information on documents even before encountering them again. This suggests that LLMs may possess the capability of sequential memorization, which could pave the way for research into more complex structured learning environments as model parameters scale up.

6.2 Conventional Types of Incremental Learning

As mentioned in Section 2.2.1, three types of incremental learning are prevalent [232]. Among them, class-incremental learning (CIL) has historically attracted significant attention from the community [193, 262]. However, in the context of continually pre-training and adapting large language models (LLMs), we observe a decreased interest in CIL but an increased focus on task-incremental learning (TIL) and domain-incremental learning (DIL). Given that language models are inherently designed for content generation and are pre-trained with the pretext generative task of next-word prediction, it is natural to emphasize the patterns of generative tasks and integrate the traditional CIL paradigm into the broader framework of language modeling, discarding the incremental classification head [26, 210]. However, the declining attention to CIL does not suggest that it is not impactful in the field of continual learning for LLMs. Techniques such as vocabulary expansion [5, 44] and learning routing function in the MoE system [35] can be seen as an extension of expanding the classification head in CIL, and previously validated techniques of CIL can be directly applied.

The importance of DIL is self-evident, given the shared task definition and input-output format in continual pre-training (CPT) and domain-adaptive pre-training (DAP). On the other hand, TIL attracts significant interest as it plays a crucial role in instruction tuning, where instructions can be seen as natural-language-encoded task indices [80, 89, 165, 208, 240, 243, 279, 297]. It is worth noting that the boundary between TIL and DIL becomes

Table 4. **Summary of the existing benchmarks publicly available for Continual Learning LLMs.** In the column of **Name**, we use the superscript ^{***} to denote the lack of the dataset name and the name shown is that of the original paper. In this table, we deliberately omit the datasets used for domain-adaptive pre-training the vertical LLMs, as their main focus of development is not on continual learning. We also omit the datasets used for general continual fine-tuning, as they are extensively discussed in other existing surveys [17, 105].

Name	Type	Shift	Domain	#Stages	Scale	Sources	Applications	Comment
*TimeLMs [137]	CPT	Temporal	Social Media	8	#Examples: 123.86M	Tweets	[137]	code
CC-RecentNews [96]	CPT	Temporal	News	1	#Tokens: ~168M	Web	[96]	code
TWiki [95]	CPT	Temporal	General Knowledge	5	#Tokens: 4.7B	Wikipedia	[95]	code
*DAPT [72]	CPT DAP	Content	Multi-Domain	4	Size: 160GB	BioMed [134], CS [134], News [283], Reviews [78]	[72] [183] [184]	code
*CPT [104]	CPT	Content	Multi-Domain	4	#Examples: 3.12M	Yelp [270], SzORC [134], AG-News [290]	[104]	code
*DEMIX [71]	CPT	Content	Multi-Domain	8	#Tokens: 73.8B	1B [31], CS [134], Legal [27], Med [134] WebText [64], RealNews [283], Reddit [13], Reviews [169]	[71]	code
*DAS [107]	CPT DAP	Content	Multi-Domain	6	Size: 4.16GB	Yelp [270], Reviews [169], Papers [134], PubMed	[107]	code
SuperNI [244]	CIT	Content	Multi-Domain	16	#Tasks: 1616 #Examples: ~5M	GitHub	[243, 292]	code
CITB [292]	CIT	Content	Multi-Domain	19	#Tasks: 38	SuperNI [244]	[292]	code
CoIN [32]	CIT	Content	Multi-Domain	8	#Examples: ~1.14M	ReCOCO [103], ReCOCO+ [151], ReCOCOg [151] ImageNet [49], VQA-v2 [65], ScienceQA [140] TextVQA [214], GQA [93], VizWiz [70], OCR-VQA [160]	[32]	code
TRACE [241]	CIT	Content	Multi-Domain	8	#Examples: 56,000	ScienceQA [140], FOMC [209], MeetingBank [88] C-STANCE [293], 20Minutes [108], CodeXGLUE [141], NumGLUE [162]	[241]	code
NATURAL-INSTRUCTION [161]	CIT	Content	Multi-Domain	6	#Examples: 193k	CosmosQA [90], DROP [56], Essential-Terms [110] MCTACO [302], MultiRC [109], QASC [111] Quoref [46], ROPES [127], Winogrande [202]	[161]	code
IMDB [149]	CMA	Content	Social Media	1	Size: 217.35 MB	IMDB	[286]	code
HH-RLHF [11]	CMA	Content	General Knowledge	1	Size: 28.1 MB	Human Feedback	[286]	code
Reddit TL-DR [234]	CMA	Content	Social Media	2	Size: 19.6 GB	Reddit	[286, 287]	code
Common Sense QA [128] Reading Comprehension [128] Translation [128]	CMA	Content	Multi-Domain	6	#Examples: ~ 41.16M	ARC Easy and Challenge [41], Race [115], PIQA [18] SQuAD [190], DROP [56] WMT 2014 French to English [19]	[128]	see sources
FEVER [228]	CMR	Content	General Knowledge	1	#Examples: 420k	Wikipedia	[47, 76]	code
VitaminC [206]	CMR	Content	General Knowledge	1	#Examples: 450k	Wikipedia	[164]	code
zsRE [117]	CMR	Content	General Knowledge	1	#Examples: 120M	Wikireading [83]	[45, 74-76, 157, 158]	-
T-rer [59]	CMR	Content	General Knowledge	1	#Examples: 11M	Dpmedia abstracts [23]	[53, 119]	code
NQ [114]	CMR	Content	General Knowledge	1	#Examples: 320k	Google queries, Wikipedia	[74]	code
CounterFact [157]	CMR	Content	General Knowledge	1	#Examples: 22k	zsRE [117]	[45, 85, 157, 280]	code
SCOTUS [28]	CMR	Temporal	Law	1	#Examples: 9.2k	Supreme Court Database	[74]	code

somewhat blurred in continual instruction tuning. Language models demonstrate the capability to infer domain information for unseen instructions, suggesting a convergence of TIL and DIL in certain contexts.

6.3 Roles of Memory in Continual LLMs

Previous continual learning research, drawing inspiration from human learning patterns, primarily emphasizes the storage efficiency of past data. However, this focus may no longer hold true in the context of continual LLMs. In the direction of relaxing memory constraints, institutions with access to training data may opt to retain full access without restricting memory size, given that the cost of memory storage is more than affordable. In such scenarios, as highlighted in [233], the challenge shifts from storage efficiency to computational efficiency. To achieve continual learning goals, models must efficiently adapt to new data (efficient adaptation) and select key experiences for replay (efficient replay) [99, 268]. Therefore, it is essential to reassess the existing memory constraint and prioritize optimizing computational efficiency for continual learning of LLMs by restricting the number of updates and FLOPs [180, 247].

On the other end of the spectrum, studies with tightened memory constraints remain vital in modern continual learning of LLMs. As shown in Fig. 1, upstream suppliers of LLMs typically do not provide training data with the released model weights. Consequently, consumers must adapt these models to downstream data without access to the actual replay data. Various rehearsal-free continual strategies are applied in this scenario, such as collecting data examples from alternate sources [9, 42, 199, 258], leveraging the generative capabilities of LLMs to produce pseudo-examples for replay [182], and implementing regularization techniques in the parameter space [107, 197]. Continual learning under the strict memory constraint is also driven by data privacy concerns, where preserving data on the server side is prohibited. In these scenarios, researchers must rely on online continual learning methods [150, 181], where data examples are only utilized for training as they arrive in a stream, and numerous efforts are already underway to develop LLMs capable of operating under these constraints [20].

6.4 Prospective Directions

Theories of Continual LLMs. It is widely recognized that the continual learning community tends to prioritize empirical research over theoretical exploration. Nevertheless, there are efforts to establish theoretical foundations for CL. In [237], the authors utilize second-order Taylor expansions around optimal parameters to derive an inter-task generalization error bound based on the maximum eigenvalue and l_2 -norm of parameter differences. Another line of approaches leverages task/domain discrepancies to construct a multi-task generalization bound. For instance, Unified Domain Incremental Learning (UDIL) in [213] proposes upper bounds for intra-domain and cross-domain distillation losses, unifying various replay-based DIL techniques under a single adaptive generalization bound. However, applying these existing theories directly to continual LLMs can be imprudent, given their pre-trained, large-scale nature. Consequently, there is a notable gap in research focusing on continually learning LLMs with robust theoretical guarantees and understanding the forgetting behaviors of LLMs from a theoretical perspective.

Efficient Replay for Knowledge Retention for Continual LLMs. While the storage budget can theoretically be infinite (Section 6.3), replaying past experiences without specific design can lead to inefficient updates in current domain learning, resulting in slow convergence. Beyond sparse replay solutions that control data mixture ratios [129, 199, 274], there is ongoing exploration of efficient replay for continual LLMs. For example, KPIG [80] enhances replay efficiency by calculating Key-Part Information Gain (KPIG) on masked segments, enabling the dynamic selection of replay data. [99] introduces a forgetting forecasting mechanism based on output changes during adaptation, later used for selective

replay in continual model refinement (CMR). More sophisticated and accurate data mixing strategies and efficient replay sample selection mechanisms are needed and hence we mark it as a significant research focus in the future.

Continual LLMs with Controllable Memory. The long-term memory inherent in the whole set of parameters of LLMs often lacks interpretability and explicit manipulability, which is crucial in certain application areas such as machine unlearning [21], where the continually pre-trained models need to constantly roll back to a previous version predating the inclusion of the revoked data and retrain the model from that point onward. This example illustrates the benefits of equipping LLMs with an external, controllable memory. As part of continual model refinement (CMR), memory systems for continual learning have been explored in several studies. Larimar [45] suggests integrating the Kanerva Machine [263] as an episodic memory for multi-fact model editing. This memory system supports basic operations like *writing, reading, and generating*, as well as advanced operations such as *sequential writing and forgetting*. It enables one-shot knowledge updates without costly retraining or fine-tuning. Other memory systems like Hopfield Networks [192] hold promise for future investigation as well.

Continual LLMs with Custom Preferences. In service-oriented contexts, users often require different trade-offs between domain expertise, ethics, values, or tones of expression. Efficiently building customized LLMs for individual users and offering flexible adjustment options is a challenging task. Early attempts in this direction include Imprecise Bayesian Continual Learning (IBCL), which, under certain assumptions, guarantees the generation of Pareto-optimal models based on user preferences by combining two model posteriors in the parameter space [139]. While empirical validation is limited in scale, this approach paves the way for future research in this area.

7 CONCLUSION

In this work, we offer a comprehensive survey on continual LLMs, summarizing recent advancements in their training and deployment from a continual learning standpoint. We categorize the problems and tasks based on their positions within our proposed broader framework of modern stratified continual learning of LLMs. While there is a widespread and growing interest in this area across the community, we also note several missing cornerstones, including algorithmic diversity and a fundamental understanding of large models' behaviors such as knowledge forgetting, transfer, and acquisition. With a holistic yet detailed approach, we aim for this survey to inspire more practitioners to explore continual learning techniques, ultimately contributing to the development of robust and self-evolving AI systems.

REFERENCES

- [1] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [2] E. C. Acikgoz, O. B. İnce, R. Bench, A. A. Boz, İ. Kesen, A. Erdem, and E. Erdem. Hippocrates: An open-source framework for advancing large language models in healthcare. *arXiv preprint arXiv:2404.16621*, 2024.
- [3] M. Agarwal, Y. Shen, B. Wang, Y. Kim, and J. Chen. Structured code representations enable data-efficient adaptation of code language models, 2024.
- [4] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuytelaars. Memory aware synapses: Learning what (not) to forget. In *Proceedings of the European conference on computer vision (ECCV)*, pages 139–154, 2018.
- [5] S. Amba Hombaiah, T. Chen, M. Zhang, M. Bendersky, and M. Najork. Dynamic language models for continuously evolving content. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2514–2524, 2021.
- [6] R. Anil, A. M. Dai, O. Firat, M. Johnson, D. Lepikhin, A. Passos, S. Shakeri, E. Taropa, P. Bailey, Z. Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.
- [7] D. Araci. Finbert: Financial sentiment analysis with pre-trained language models, 2019.
- [8] G. Attanasio, D. Nozza, F. Bianchi, and D. Hovy. Is it worth the (environmental) cost? limited evidence for temporal adaptation via continuous training, 2023.
- [9] Z. Azerbayev, H. Schoelkopf, K. Paster, M. D. Santos, S. McAleer, A. Q. Jiang, J. Deng, S. Biderman, and S. Welleck. Llemma: An open language model for mathematics. *CoRR*, abs/2310.10631, 2023.

- [10] X. Bai, J. Shang, Y. Sun, and N. Balasubramanian. Enhancing continual learning with global prototypes: Counteracting negative representation drift, 2023.
- [11] Y. Bai, A. Jones, K. Ndousse, A. Askell, A. Chen, N. DasSarma, D. Drain, S. Fort, D. Ganguli, T. Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.
- [12] S. Banerjee and A. Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.
- [13] J. Baumgartner, S. Zannettou, B. Keegan, M. Squire, and J. Blackburn. The pushshift reddit dataset, 2020.
- [14] S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning from different domains. *Machine learning*, 79:151–175, 2010.
- [15] Z. Bi, N. Zhang, Y. Xue, Y. Ou, D. Ji, G. Zheng, and H. Chen. Oceanpnt: A large language model for ocean science tasks. *CoRR*, abs/2310.02031, 2023.
- [16] S. Biderman, H. Schoelkopf, Q. G. Anthony, H. Bradley, K. O’Brien, E. Hallahan, M. A. Khan, S. Purohit, U. S. Prashanth, E. Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR, 2023.
- [17] M. Biesialska, K. Biesialska, and M. R. Costa-jussà. Continual lifelong learning in natural language processing: A survey. In D. Scott, N. Bel, and C. Zong, editors, *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6523–6541, Barcelona, Spain (Online), Dec. 2020. International Committee on Computational Linguistics.
- [18] Y. Bisk, R. Zellers, J. Gao, Y. Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [19] O. Bojar, C. Buck, C. Federmann, B. Haddow, P. Koehn, J. Leveling, C. Monz, P. Pecina, M. Post, H. Saint-Amand, et al. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*, pages 12–58, 2014.
- [20] J. Bornschein, Y. Li, and A. Rannen-Triki. Transformers for supervised online continual learning, 2024.
- [21] L. Bourtole, V. Chandrasekaran, C. A. Choquette-Choo, H. Jia, A. Travers, B. Zhang, D. Lie, and N. Papernot. Machine unlearning, 2020.
- [22] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [23] M. Brümmer, M. Dojchinovski, and S. Hellmann. Dbpedia abstracts: A large-scale, open, multilingual nlp training corpus. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 3339–3343, 2016.
- [24] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara. Dark experience for general continual learning: a strong, simple baseline. *Advances in neural information processing systems*, 33:15920–15930, 2020.
- [25] H. Cao, Z. Liu, X. Lu, Y. Yao, and Y. Li. Instructmol: Multi-modal integration for building a versatile and reliable molecular assistant in drug discovery. *CoRR*, abs/2311.16208, 2023.
- [26] X. Cao, H. Lu, L. Huang, X. Liu, and M.-M. Cheng. Generative multi-modal models are good class incremental learners. *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2024.
- [27] Caselaw Access Project. Caselaw access project, 2018.
- [28] I. Chalkidis, T. Pasini, S. Zhang, L. Tomada, S. F. Schwemer, and A. Søgaard. Fairlex: A multilingual benchmark for evaluating fairness in legal text processing. *arXiv preprint arXiv:2203.07228*, 2022.
- [29] A. Chaudhry, M. Ranzato, M. Rohrbach, and M. Elhoseiny. Efficient lifelong learning with a-gem. In *ICLR*, 2019.
- [30] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- [31] C. Chelba, T. Mikolov, M. Schuster, Q. Ge, T. Brants, P. Koehn, and T. Robinson. One billion word benchmark for measuring progress in statistical language modeling, 2014.
- [32] C. Chen, J. Zhu, X. Luo, H. Shen, L. Gao, and J. Song. Coin: A benchmark of continual instruction tuning for multimodal large language model, 2024.
- [33] J. Chen, X. Wang, A. Gao, F. Jiang, S. Chen, H. Zhang, D. Song, W. Xie, C. Kong, J. Li, X. Wan, H. Li, and B. Wang. Huatuogpt-ii, one-stage training for medical adaption of llms. *CoRR*, abs/2311.09774, 2023.
- [34] S. Chen, Y. Hou, Y. Cui, W. Che, T. Liu, and X. Yu. Recall and learn: Fine-tuning deep pretrained language models with less forgetting. In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7870–7881, Online, Nov. 2020. Association for Computational Linguistics.
- [35] W. Chen, Y. Zhou, N. Du, Y. Huang, J. Laudon, Z. Chen, and C. Cui. Lifelong language pretraining with distribution-specialized experts. In *International Conference on Machine Learning*, pages 5383–5395. PMLR, 2023.
- [36] X. Chen, Z. Wang, D. Sow, J. Yang, T. Chen, Y. Liang, M. Zhou, and Z. Wang. Take the bull by the horns: Hard sample-reweighted continual training improves llm generalization. *arXiv preprint arXiv:2402.14270*, 2024.
- [37] Y. Chen, S. Zhang, G. Qi, and X. Guo. Parameterizing context: Unleashing the power of parameter-efficient fine-tuning and in-context tuning for continual table semantic parsing. *Advances in Neural Information Processing Systems*, 36, 2024.
- [38] Z. Chen and B. Liu. *Lifelong machine learning*, volume 1. Springer.
- [39] D. Cheng, S. Huang, and F. Wei. Adapting large language models via reading comprehension, 2024.
- [40] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.

- [41] P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [42] P. Colombo, T. P. Pires, M. Boudiaf, D. Culver, R. Melo, C. Corro, A. F. T. Martins, F. Esposito, V. L. Raposo, S. Morgado, and M. Desa. Saullm-7b: A pioneering large language model for law, 2024.
- [43] T. Computer. Redpajama: an open dataset for training large language models, 2023.
- [44] A. Cossu, T. Tuytelaars, A. Carta, L. Passaro, V. Lomonaco, and D. Bacciu. Continual pre-training mitigates forgetting in language and vision, 2022.
- [45] P. Das, S. Chaudhury, E. Nelson, I. Melnyk, S. Swaminathan, S. Dai, A. Lozano, G. Kollias, V. Chenthamarakshan, S. Dan, et al. Larimar: Large language models with episodic memory control. *arXiv preprint arXiv:2403.11901*, 2024.
- [46] P. Dasigi, N. F. Liu, A. Marasović, N. A. Smith, and M. Gardner. Quoref: A reading comprehension dataset with questions requiring coreferential reasoning. In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5925–5932, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [47] N. De Cao, W. Aziz, and I. Titov. Editing factual knowledge in language models. *arXiv preprint arXiv:2104.08164*, 2021.
- [48] C. Deng, T. Zhang, Z. He, Y. Xu, Q. Chen, Y. Shi, L. Fu, W. Zhang, X. Wang, C. Zhou, Z. Lin, and J. He. K2: A foundation language model for geoscience knowledge understanding and utilization, 2023.
- [49] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR09*, 2009.
- [50] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [51] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [52] B. Dhingra, J. R. Cole, J. M. Eisenschlos, D. Gillick, J. Eisenstein, and W. W. Cohen. Time-aware language models as temporal knowledge bases. *Transactions of the Association for Computational Linguistics*, 10:257–273, 2022.
- [53] Q. Dong, D. Dai, Y. Song, J. Xu, Z. Sui, and L. Li. Calibrating factual knowledge in pretrained language models. *arXiv preprint arXiv:2210.03329*, 2022.
- [54] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [55] L. Dou, Q. Liu, G. Zeng, J. Guo, J. Zhou, W. Lu, and M. Lin. Sailor: Open language models for south-east asia. *arXiv preprint arXiv:2404.03608*, 2024.
- [56] D. Dua, Y. Wang, P. Dasigi, G. Stanovsky, S. Singh, and M. Gardner. Drop: A reading comprehension benchmark requiring discrete reasoning over paragraphs. *arXiv preprint arXiv:1903.00161*, 2019.
- [57] S. Ebrahimi, M. Elhoseiny, T. Darrell, and M. Rohrbach. Uncertainty-guided continual learning with bayesian neural networks. *arXiv preprint arXiv:1906.02425*, 2019.
- [58] S. Ebrahimi, F. Meier, R. Calandra, T. Darrell, and M. Rohrbach. Adversarial continual learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 386–402. Springer, 2020.
- [59] H. Elsahar, P. Vougiouklis, A. Remaci, C. Gravier, J. Hare, F. Laforest, and E. Simperl. T-rex: A large scale alignment of natural language with knowledge base triples. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- [60] K. Fujii, T. Nakamura, M. Loem, H. Iida, M. Ohi, K. Hattori, H. Shota, S. Mizuki, R. Yokota, and N. Okazaki. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities. *arXiv preprint arXiv:2404.17790*, 2024.
- [61] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, and V. Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.
- [62] S. Garg, M. Farajtabar, H. Pouransari, R. Vemulapalli, S. Mehta, O. Tuzel, V. Shankar, and F. Faghri. Tic-clip: Continual training of clip models. In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.
- [63] E. Gogoulou, T. Lesort, M. Boman, and J. Nivre. Continual learning under language shift, 2024.
- [64] A. Gokaslan and V. Cohen. Openwebtext corpus, 2019.
- [65] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering, 2017.
- [66] Y. Gu, R. Tinn, H. Cheng, M. Lucas, N. Usuyama, X. Liu, T. Naumann, J. Gao, and H. Poon. Domain-specific language model pretraining for biomedical natural language processing. *ACM Transactions on Computing for Healthcare*, 3(1):1–23, Oct. 2021.
- [67] D. Guo, Q. Zhu, D. Yang, Z. Xie, K. Dong, W. Zhang, G. Chen, X. Bi, Y. Wu, Y. K. Li, F. Luo, Y. Xiong, and W. Liang. Deepseek-coder: When the large language model meets programming – the rise of code intelligence, 2024.
- [68] Z. Guo and Y. Hua. Continuous training and fine-tuning for domain-specific language models in medical question answering, 2023.
- [69] K. Gupta, B. Thérien, A. Ibrahim, M. L. Richter, Q. Anthony, E. Belilovsky, I. Rish, and T. Lesort. Continual pre-training of large language models: How to (re)warm your model?, 2023.
- [70] D. Gurari, Q. Li, A. J. Stangl, A. Guo, C. Lin, K. Grauman, J. Luo, and J. P. Bigham. Vizviz grand challenge: Answering visual questions from blind people, 2018.
- [71] S. Gururangan, M. Lewis, A. Holtzman, N. A. Smith, and L. Zettlemoyer. DEMix layers: Disentangling domains for modular language modeling. In M. Carpuat, M.-C. de Marneffe, and I. V. Meza Ruiz, editors, *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5557–5576, Seattle, United States, July 2022. Association for Computational Linguistics.

- [72] S. Gururangan, A. Marasović, S. Swayamdipta, K. Lo, I. Beltagy, D. Downey, and N. A. Smith. Don't stop pretraining: Adapt language models to domains and tasks. In D. Jurafsky, J. Chai, N. Schlueter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8342–8360, Online, July 2020. Association for Computational Linguistics.
- [73] R. Han, X. Ren, and N. Peng. ECONET: Effective continual pretraining of language models for event temporal reasoning. In M.-F. Moens, X. Huang, L. Specia, and S. W.-t. Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5367–5380, Online and Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics.
- [74] T. Hartvigsen, S. Sankaranarayanan, H. Palangi, Y. Kim, and M. Ghassemi. Aging with grace: Lifelong model editing with discrete key-value adaptors. In *Advances in Neural Information Processing Systems*, 2023.
- [75] P. Hase, M. Bansal, B. Kim, and A. Ghandeharioun. Does localization inform editing? surprising differences in causality-based localization vs. *Knowledge Editing in Language Models*, 2023.
- [76] P. Hase, M. Diab, A. Celikyilmaz, X. Li, Z. Kozareva, V. Stoyanov, M. Bansal, and S. Iyer. Do language models have beliefs? methods for detecting, updating, and visualizing model beliefs. *arXiv preprint arXiv:2111.13654*, 2021.
- [77] J. He, H. Guo, M. Tang, and J. Wang. Continual instruction tuning for large multimodal models, 2023.
- [78] R. He and J. McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *Proceedings of the 25th International Conference on World Wide Web*, WWW '16, page 507–517, Republic and Canton of Geneva, CHE, 2016. International World Wide Web Conferences Steering Committee.
- [79] T. He, J. Liu, K. Cho, M. Ott, B. Liu, J. Glass, and F. Peng. Analyzing the forgetting problem in pretrain-finetuning of open-domain dialogue response models. In P. Merlo, J. Tiedemann, and R. Tsarfaty, editors, *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1121–1133, Online, Apr. 2021. Association for Computational Linguistics.
- [80] Y. He, X. Huang, M. Tang, L. Meng, X. Li, W. Lin, W. Zhang, and Y. Gao. Don't half-listen: Capturing key-part information in continual instruction tuning, 2024.
- [81] D. Hendrycks, C. Burns, S. Basart, A. Critch, J. Li, D. Song, and J. Steinhardt. Aligning ai with shared human values, 2023.
- [82] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring massive multitask language understanding. *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [83] D. Hewlett, A. Lacoste, L. Jones, I. Polosukhin, A. Fandrianto, J. Han, M. Kelcey, and D. Berthelot. Wikireading: A novel large-scale language understanding task over wikipedia. *arXiv preprint arXiv:1608.03542*, 2016.
- [84] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [85] C. Hu, P. Cao, Y. Chen, K. Liu, and J. Zhao. Wilke: Wise-layer knowledge editor for lifelong knowledge editing, 2024.
- [86] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [87] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2022.
- [88] Y. Hu, T. Ganter, H. Deilamsalehy, F. Deroncourt, H. Foroosh, and F. Liu. Meetingbank: A benchmark dataset for meeting summarization, 2023.
- [89] J. Huang, L. Cui, A. Wang, C. Yang, X. Liao, L. Song, J. Yao, and J. Su. Mitigating catastrophic forgetting in large language models with self-synthesized rehearsal, 2024.
- [90] L. Huang, R. L. Bras, C. Bhagavatula, and Y. Choi. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning, 2019.
- [91] Q. Huang, M. Tao, Z. An, C. Zhang, C. Jiang, Z. Chen, Z. Wu, and Y. Feng. Lawyer llama technical report. *arXiv preprint arXiv:2305.15062*, 2023.
- [92] Z. Huang, Y. Shen, X. Zhang, J. Zhou, W. Rong, and Z. Xiong. Transformer-patcher: One mistake worth one neuron. *arXiv preprint arXiv:2301.09785*, 2023.
- [93] D. A. Hudson and C. D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering, 2019.
- [94] A. Ibrahim, B. Thérien, K. Gupta, M. L. Richter, Q. Anthony, T. Lesort, E. Belilovsky, and I. Rish. Simple and scalable strategies to continually pre-train large language models. *arXiv preprint arXiv:2403.08763*, 2024.
- [95] J. Jang, S. Ye, C. Lee, S. Yang, J. Shin, J. Han, G. Kim, and M. Seo. Temporalwiki: A lifelong benchmark for training and evaluating ever-evolving language models. 2022.
- [96] J. Jang, S. Ye, S. Yang, J. Shin, J. Han, G. Kim, S. J. Choi, and M. Seo. Towards continual knowledge learning of language models. In *ICLR*, 2022.
- [97] J. Ji, T. Qiu, B. Chen, B. Zhang, H. Lou, K. Wang, Y. Duan, Z. He, J. Zhou, Z. Zhang, F. Zeng, K. Y. Ng, J. Dai, X. Pan, A. O'Gara, Y. Lei, H. Xu, B. Tse, J. Fu, S. McAleer, Y. Yang, Y. Wang, S.-C. Zhu, Y. Guo, and W. Gao. Ai alignment: A comprehensive survey, 2024.
- [98] Z. Jiang, Z. Sun, W. Shi, P. Rodriguez, C. Zhou, G. Neubig, X. V. Lin, W. tau Yih, and S. Iyer. Instruction-tuned language models are better knowledge learners, 2024.
- [99] X. Jin and X. Ren. What will my model forget? forecasting forgotten examples in language model refinement, 2024.
- [100] X. Jin, D. Zhang, H. Zhu, W. Xiao, S.-W. Li, X. Wei, A. Arnold, and X. Ren. Lifelong pretraining: Continually adapting language models to emerging corpora. In A. Fan, S. Ilic, T. Wolf, and M. Gallé, editors, *Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models*, pages 1–16, virtual+Dublin, May 2022. Association for Computational Linguistics.
- [101] E. R. Kandel, J. H. Schwartz, T. M. Jessell, S. Siegelbaum, A. J. Hudspeth, S. Mack, et al. *Principles of neural science*, volume 4. McGraw-hill New York, 2000.

- [102] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [103] S. Kazemzadeh, V. Ordonez, M. Matten, and T. Berg. ReferItGame: Referring to objects in photographs of natural scenes. In A. Moschitti, B. Pang, and W. Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 787–798, Doha, Qatar, Oct. 2014. Association for Computational Linguistics.
- [104] Z. Ke, H. Lin, Y. Shao, H. Xu, L. Shu, and B. Liu. Continual training of language models for few-shot learning. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10205–10216, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics.
- [105] Z. Ke and B. Liu. Continual learning of natural language processing tasks: A survey, 2023.
- [106] Z. Ke, B. Liu, N. Ma, H. Xu, and S. Lei. Achieving forgetting prevention and knowledge transfer in continual learning. In *NeurIPS*, 2021.
- [107] Z. Ke, Y. Shao, H. Lin, T. Konishi, G. Kim, and B. Liu. Continual pre-training of language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- [108] T. Kew, M. Kostrzewa, and S. Ebling. 20 minuten: A multi-task news summarisation dataset for German. In H. Ghorbel, M. Sokhn, M. Cieliebak, M. Hürlimann, E. de Salis, and J. Guerne, editors, *Proceedings of the 8th edition of the Swiss Text Analytics Conference*, pages 1–13, Neuchatel, Switzerland, June 2023. Association for Computational Linguistics.
- [109] D. Khashabi, S. Chaturvedi, M. Roth, S. Upadhyay, and D. Roth. Looking beyond the surface: A challenge set for reading comprehension over multiple sentences. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 252–262, 2018.
- [110] D. Khashabi, T. Khot, A. Sabharwal, and D. Roth. Learning what is essential in questions. In R. Levy and L. Specia, editors, *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL 2017)*, pages 80–89, Vancouver, Canada, Aug. 2017. Association for Computational Linguistics.
- [111] T. Khot, P. Clark, M. Guerquin, P. Jansen, and A. Sabharwal. Qasc: A dataset for question answering via sentence composition, 2020.
- [112] G. Kim, C. Xiao, T. Konishi, Z. Ke, and B. Liu. A theoretical study on solving continual learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 5065–5079. Curran Associates, Inc., 2022.
- [113] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [114] T. Kwiakowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [115] G. Lai, Q. Xie, H. Liu, Y. Yang, and E. Hovy. Race: Large-scale reading comprehension dataset from examinations. *arXiv preprint arXiv:1704.04683*, 2017.
- [116] A. Lazaridou, A. Kuncoro, E. Gribovskaya, D. Agrawal, A. Liska, T. Terzi, M. Gimenez, C. de Masson d’Autume, T. Kocisky, S. Ruder, et al. Mind the gap: Assessing temporal generalization in neural language models. *Advances in Neural Information Processing Systems*, 34:29348–29363, 2021.
- [117] O. Levy, M. Seo, E. Choi, and L. Zettlemoyer. Zero-shot relation extraction via reading comprehension. *arXiv preprint arXiv:1706.04115*, 2017.
- [118] C.-A. Li and H.-Y. Lee. Examining forgetting in continual pre-training of aligned large language models, 2024.
- [119] D. Li, A. S. Rawat, M. Zaheer, X. Wang, M. Lukasik, A. Veit, F. Yu, and S. Kumar. Large language models with controllable working memory. *arXiv preprint arXiv:2211.05110*, 2022.
- [120] H. Li, Q. Ai, J. Chen, Q. Dong, Z. Wu, Y. Liu, C. Chen, and Q. Tian. Blade: Enhancing black-box large language models with small domain-specific models. *arXiv preprint arXiv:2403.18365*, 2024.
- [121] J. Li, Y. Bian, G. Wang, Y. Lei, D. Cheng, Z. Ding, and C. Jiang. Cfgpt: Chinese financial assistant with large language model, 2023.
- [122] J. Li, D. Li, S. Savarese, and S. Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
- [123] L. Li and X. Qiu. CONTINUAL MODEL EVOLVEMENT WITH INNER-PRODUCT RESTRICTION, 2023.
- [124] Z. Li and D. Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [125] B. Y. Lin, S. Wang, X. Lin, R. Jia, L. Xiao, X. Ren, and S. Yih. On continual model refinement in out-of-distribution data streams. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3128–3139, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [126] C.-Y. Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004.
- [127] K. Lin, O. Tafjord, P. Clark, and M. Gardner. Reasoning over paragraph effects in situations, 2019.
- [128] Y. Lin, H. Lin, W. Xiong, S. Diao, J. Liu, J. Zhang, R. Pan, H. Wang, W. Hu, H. Zhang, H. Dong, R. Pi, H. Zhao, N. Jiang, H. Ji, Y. Yao, and T. Zhang. Mitigating the alignment tax of rlhf, 2024.
- [129] Z. Lin, C. Deng, L. Zhou, T. Zhang, Y. Xu, Y. Xu, Z. He, Y. Shi, B. Dai, Y. Song, B. Zeng, Q. Chen, T. Shi, T. Huang, Y. Xu, S. Wang, L. Fu, W. Zhang, J. He, C. Ma, Y. Zhu, X. Wang, and C. Zhou. Geogalactica: A scientific large language model in geoscience, 2023.
- [130] Z. Lin, Z. Gou, Y. Gong, X. Liu, Y. Shen, R. Xu, C. Lin, Y. Yang, J. Jiao, N. Duan, et al. Rho-1: Not all tokens are what you need. *arXiv preprint arXiv:2404.07965*, 2024.
- [131] H. Liu, C. Li, Q. Wu, and Y. J. Lee. Visual instruction tuning, 2023.
- [132] H. Liu, D. Tam, M. Muqeeth, J. Mohta, T. Huang, M. Bansal, and C. A. Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context learning. *Advances in Neural Information Processing Systems*, 35:1950–1965, 2022.

- [133] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [134] K. Lo, L. L. Wang, M. Neumann, R. Kinney, and D. Weld. S2ORC: The semantic scholar open research corpus. In D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4969–4983, Online, July 2020. Association for Computational Linguistics.
- [135] V. Lomonaco, D. Maltoni, and L. Pellegrini. Rehearsal-free continual learning over small non-i.i.d. batches, 2020.
- [136] D. Lopez-Paz and M. Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30, 2017.
- [137] D. Loureiro, F. Barbieri, L. Neves, L. Espinosa Anke, and J. Camacho-collados. TimeLMs: Diachronic language models from Twitter. In V. Basile, Z. Kozareva, and S. Stajner, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 251–260, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [138] D. Lu, H. Wu, J. Liang, Y. Xu, Q. He, Y. Geng, M. Han, Y. Xin, and Y. Xiao. Bbt-fin: Comprehensive construction of chinese financial domain pre-trained language model, corpus and benchmark. *CoRR*, abs/2302.09432, 2023.
- [139] P. Lu, M. Caprio, E. Eaton, and I. Lee. Ibcl: Zero-shot model generation for task trade-offs in continual learning, 2023.
- [140] P. Lu, S. Mishra, T. Xia, L. Qiu, K.-W. Chang, S.-C. Zhu, O. Tafjord, P. Clark, and A. Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering, 2022.
- [141] S. Lu, D. Guo, S. Ren, J. Huang, A. Svyatkovskiy, A. Blanco, C. Clement, D. Drain, D. Jiang, D. Tang, G. Li, L. Zhou, L. Shou, L. Zhou, M. Tufano, M. Gong, M. Zhou, N. Duan, N. Sundaresan, S. K. Deng, S. Fu, and S. Liu. Codexglue: A machine learning benchmark dataset for code understanding and generation, 2021.
- [142] H. Luo, Q. Sun, C. Xu, P. Zhao, J. Lou, C. Tao, X. Geng, Q. Lin, S. Chen, and D. Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. *arXiv preprint arXiv:2308.09583*, 2023.
- [143] R. Luo, L. Sun, Y. Xia, T. Qin, S. Zhang, H. Poon, and T.-Y. Liu. Biogpt: generative pre-trained transformer for biomedical text generation and mining. *Briefings in Bioinformatics*, 23(6), Sept. 2022.
- [144] Y. Luo, Z. Yang, X. Bai, F. Meng, J. Zhou, and Y. Zhang. Investigating forgetting in pre-trained representations through continual learning, 2023.
- [145] Y. Luo, Z. Yang, F. Meng, Y. Li, J. Zhou, and Y. Zhang. An empirical study of catastrophic forgetting in large language models during continual fine-tuning, 2023.
- [146] Y. Luo, J. Zhang, S. Fan, K. Yang, Y. Wu, M. Qiao, and Z. Nie. Biomedgpt: Open multimodal generative pre-trained transformer for biomedicine. *arXiv preprint arXiv:2308.09442*, 2023.
- [147] Z. Luo, C. Xu, P. Zhao, Q. Sun, X. Geng, W. Hu, C. Tao, J. Ma, Q. Lin, and D. Jiang. Wizardcoder: Empowering code large language models with evol-instruct, 2023.
- [148] S. Ma, S. Huang, S. Huang, X. Wang, Y. Li, H.-T. Zheng, P. Xie, F. Huang, and Y. Jiang. Ecomgpt-ct: Continual pre-training of e-commerce large language models with semi-structured data, 2023.
- [149] A. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th annual meeting of the association for computational linguistics: Human language technologies*, pages 142–150, 2011.
- [150] Z. Mai, R. Li, J. Jeong, D. Quispe, H. Kim, and S. Sanner. Online continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51, 2022.
- [151] J. Mao, J. Huang, A. Toshev, O. Camburu, A. Yuille, and K. Murphy. Generation and comprehension of unambiguous object descriptions, 2016.
- [152] V. Mazzia, A. Pedrani, A. Caciolai, K. Rottmann, and D. Bernardi. A survey on knowledge editing of neural networks. *arXiv preprint arXiv:2310.19704*, 2023.
- [153] D. McCaffary. Towards continual task learning in artificial neural networks: current approaches and insights from neuroscience. *arXiv preprint arXiv:2112.14146*, 2021.
- [154] J. L. McClelland, B. L. McNaughton, and R. C. O’Reilly. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, 102(3):419, 1995.
- [155] M. McCloskey and N. J. Cohen. Catastrophic interference in connectionist networks: The sequential learning problem. volume 24 of *Psychology of Learning and Motivation*, pages 109–165. Academic Press, 1989.
- [156] S. V. Mehta, D. Patil, S. Chandar, and E. Strubell. An empirical investigation of the role of pre-training in lifelong learning. *Journal of Machine Learning Research*, 24(214):1–50, 2023.
- [157] K. Meng, D. Bau, A. Andonian, and Y. Belinkov. Locating and editing factual associations in gpt. *Advances in Neural Information Processing Systems*, 35:17359–17372, 2022.
- [158] K. Meng, A. S. Sharma, A. Andonian, Y. Belinkov, and D. Bau. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*, 2022.
- [159] S. Min, X. Lyu, A. Holtzman, M. Artetxe, M. Lewis, H. Hajishirzi, and L. Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- [160] A. Mishra, S. Shekhar, A. K. Singh, and A. Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*, 2019.
- [161] S. Mishra, D. Khashabi, C. Baral, and H. Hajishirzi. Natural instructions: Benchmarking generalization to new tasks from natural language instructions. *arXiv preprint arXiv:2104.08773*, 2021.
- [162] S. Mishra, A. Mitra, N. Varshney, B. Sachdeva, P. Clark, C. Baral, and A. Kalyan. Numglue: A suite of fundamental yet challenging mathematical reasoning tasks, 2022.

- [163] E. Mitchell, C. Lin, A. Bosselut, C. Finn, and C. D. Manning. Fast model editing at scale. *arXiv preprint arXiv:2110.11309*, 2021.
- [164] E. Mitchell, C. Lin, A. Bosselut, C. D. Manning, and C. Finn. Memory-based model editing at scale. In *International Conference on Machine Learning*, pages 15817–15831. PMLR, 2022.
- [165] J. Mok, J. Do, S. Lee, T. Taghavi, S. Yu, and S. Yoon. Large-scale lifelong learning of in-context instructions and how to tackle it. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12573–12589, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [166] A. Moradi Dakhel, V. Majdinasab, A. Nikanjam, F. Khomh, M. C. Desmarais, and Z. M. J. Jiang. Github copilot ai pair programmer: Asset or liability? *Journal of Systems and Software*, 203:111734, 2023.
- [167] T. Nakamura, M. Mishra, S. Tedeschi, Y. Chai, J. T. Stillerman, F. Friedrich, P. Yadav, T. Laud, V. M. Chien, T. Y. Zhuo, et al. Aurora-m: The first open source multilingual language model red-teamed according to the us executive order. *arXiv preprint arXiv:2404.00399*, 2024.
- [168] T. D. Nguyen, Y. Ting, I. Ciuca, C. O’Neill, Z. Sun, M. Jablonska, S. Kruk, E. Perkowski, J. W. Miller, J. Li, J. Peek, K. Iyer, T. Rózanski, P. Khetarpal, S. Zaman, D. Brodrick, S. J. R. Méndez, T. Bui, A. Goodman, A. Accomazzi, J. P. Naiman, J. Cranney, K. Schawinski, and UniverseTBD. Astrollama: Towards specialized foundation models in astronomy. *CoRR*, abs/2309.06126, 2023.
- [169] J. Ni, J. Li, and J. McAuley. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 188–197, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [170] Z. Ni, H. Shi, S. Tang, L. Wei, Q. Tian, and Y. Zhuang. Revisiting catastrophic forgetting in class incremental learning. *arXiv preprint arXiv:2107.12308*, 2021.
- [171] Z. Ni, L. Wei, S. Tang, Y. Zhuang, and Q. Tian. Continual vision-language representation learning with off-diagonal information. In *Proceedings of the 40th International Conference on Machine Learning*, pages 26129–26149, 2023.
- [172] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong. Codegen: An open large language model for code with multi-turn program synthesis. *ICLR*, 2023.
- [173] OpenAI. Introducing chatgpt. [online]. available: <https://openai.com/blog/chatgpt>. 2022.
- [174] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe. Training language models to follow instructions with human feedback, 2022.
- [175] C. Pallier, S. Dehaene, J.-B. Poline, D. LeBihan, A.-M. Argenti, E. Dupoux, and J. Mehler. Brain imaging of language plasticity in adopted adults: Can a second language replace the first? *Cerebral cortex*, 13(2):155–161, 2003.
- [176] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.
- [177] I. Paul, J. Luo, G. Glavaš, and I. Gurevych. Ircoder: Intermediate representations make language models robust multilingual code generators, 2024.
- [178] A. Pentina. *Theoretical foundations of multi-task lifelong learning*. PhD thesis, 2016.
- [179] F. Petroni, T. Rocktäschel, S. Riedel, P. Lewis, A. Bakhtin, Y. Wu, and A. Miller. Language models as knowledge bases? In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2463–2473, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [180] A. Prabhu, H. A. Al Kader Hammoud, P. K. Dokania, P. H. Torr, S.-N. Lim, B. Ghanem, and A. Bibi. Computationally budgeted continual learning: What does matter? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3698–3707, 2023.
- [181] A. Prabhu, Z. Cai, P. Dokania, P. Torr, V. Koltun, and O. Sener. Online continual learning without the storage constraint, 2023.
- [182] C. Qin and S. Joty. Lfpt5: A unified framework for lifelong few-shot language learning based on prompt tuning of t5. In *International Conference on Learning Representations*, 2021.
- [183] Y. Qin, C. Qian, X. Han, Y. Lin, H. Wang, R. Xie, Z. Liu, M. Sun, and J. Zhou. Recyclable tuning for continual pre-training. *arXiv preprint arXiv:2305.08702*, 2023.
- [184] Y. Qin, J. Zhang, Y. Lin, Z. Liu, P. Li, M. Sun, and J. Zhou. ELLE: Efficient lifelong pre-training for emerging data. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2789–2810, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [185] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [186] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [187] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [188] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [189] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.
- [190] P. Rajpurkar, R. Jia, and P. Liang. Know what you don’t know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*, 2018.
- [191] R. Ramesh and P. Chaudhari. Model zoo: A growing” brain” that learns continually. *arXiv preprint arXiv:2106.03027*, 2021.

- [192] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M. Holzleitner, M. Pavlović, G. K. Sandve, V. Greiff, D. Kreil, M. Kopp, G. Klambauer, J. Brandstetter, and S. Hochreiter. Hopfield networks is all you need, 2021.
- [193] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- [194] M. Reid, N. Savinov, D. Teplyashin, D. Lepikhin, T. Lillicrap, J.-b. Alayrac, R. Soricut, A. Lazaridou, O. Firat, J. Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [195] M. Riemer, I. Cases, R. Ajemian, M. Liu, I. Rish, Y. Tu, and G. Tesauro. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.
- [196] H. Ritter, A. Botev, and D. Barber. Online structured laplace approximations for overcoming catastrophic forgetting. *Advances in Neural Information Processing Systems*, 31, 2018.
- [197] S. Rongali, A. Jagannatha, B. P. S. Rawat, and H. Yu. Continual domain-tuning for pretrained language models, 2021.
- [198] G. D. Rosin, I. Guy, and K. Radinsky. Time masking for temporal language models. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining, WSDM '22*, page 833–841, New York, NY, USA, 2022. Association for Computing Machinery.
- [199] B. Rozière, J. Gehring, F. Gloeckle, S. Sootla, I. Gat, X. E. Tan, Y. Adi, J. Liu, R. Sauvestre, T. Remez, J. Rapin, A. Kozhevnikov, I. Evtimov, J. Bitton, M. Bhatt, C. C. Ferrer, A. Grattafiori, W. Xiong, A. Défossez, J. Copet, F. Azhar, H. Touvron, L. Martin, N. Usunier, T. Scialom, and G. Synnaeve. Code llama: Open foundation models for code, 2024.
- [200] A. N. Rubungo, C. Arnold, B. P. Rand, and A. B. Dieng. Llm-prop: Predicting physical and electronic properties of crystalline solids from their text descriptions. *CoRR*, abs/2310.14029, 2023.
- [201] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [202] K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019.
- [203] V. Sanh, A. Webson, C. Raffel, S. H. Bach, L. Sutawika, Z. Alyafeai, A. Chaffin, A. Stiegler, T. L. Scao, A. Raja, M. Dey, M. S. Bari, C. Xu, U. Thakker, S. S. Sharma, E. Szczechla, T. Kim, G. Chhablani, N. Nayak, D. Datta, J. Chang, M. T.-J. Jiang, H. Wang, M. Manica, S. Shen, Z. X. Yong, H. Pandey, R. Bawden, T. Wang, T. Neeraj, J. Rozen, A. Sharma, A. Santilli, T. Fevry, J. A. Fries, R. Teehan, T. Bers, S. Biderman, L. Gao, T. Wolf, and A. M. Rush. Multitask prompted training enables zero-shot task generalization, 2022.
- [204] F. Sarfraz, E. Arani, and B. Zonooz. Error sensitivity modulation based experience replay: Mitigating abrupt representation drift in continual learning. *arXiv preprint arXiv:2302.11344*, 2023.
- [205] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms, 2017.
- [206] T. Schuster, A. Fisch, and R. Barzilay. Get your vitamin c! robust fact verification with contrastive evidence. *arXiv preprint arXiv:2103.08541*, 2021.
- [207] J. Schwarz, W. Czarneci, J. Luketina, A. Grabska-Barwinska, Y. W. Teh, R. Pascanu, and R. Hadsell. Progress & compress: A scalable framework for continual learning. In *International conference on machine learning*, pages 4528–4537. PMLR, 2018.
- [208] T. Scialom, T. Chakrabarty, and S. Muresan. Fine-tuned language models are continual learners. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 6107–6122, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics.
- [209] A. Shah, S. Paturi, and S. Chava. Trillion dollar words: A new financial dataset, task & market analysis, 2023.
- [210] Y. Shao, Y. Guo, D. Zhao, and B. Liu. Class-incremental learning based on label generation. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1263–1276, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [211] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017.
- [212] J. Shen, N. Tenenholtz, J. B. Hall, D. Alvarez-Melis, and N. Fusi. Tag-llm: Repurposing general-purpose llms for specialized domains. *arXiv preprint arXiv:2402.05140*, 2024.
- [213] H. Shi and H. Wang. A unified approach to domain incremental learning with memory: Theory and algorithm. *Advances in Neural Information Processing Systems*, 36, 2024.
- [214] A. Singh, V. Natarajan, M. Shah, Y. Jiang, X. Chen, D. Batra, D. Parikh, and M. Rohrbach. Towards vqa models that can read, 2019.
- [215] A. Sinitin, V. Plokhotnyuk, D. Pyrkin, S. Popov, and A. Babenko. Editable neural networks. *arXiv preprint arXiv:2004.00345*, 2020.
- [216] L. Soldaini, R. Kinney, A. Bhagia, D. Schwenk, D. Atkinson, R. Authur, B. Bogin, K. Chandu, J. Dumas, Y. Elazar, V. Hofmann, A. H. Jha, S. Kumar, L. Lucy, X. Lyu, N. Lambert, I. Magnusson, J. Morrison, N. Muennighoff, A. Naik, C. Nam, M. E. Peters, A. Ravichander, K. Richardson, Z. Shen, E. Strubell, N. Subramani, O. Tafjord, P. Walsh, L. Zettlemoyer, N. A. Smith, H. Hajishirzi, I. Beltagy, D. Groeneveld, J. Dodge, and K. Lo. Dolma: an open corpus of three trillion tokens for language model pretraining research, 2024.
- [217] C. Song, X. Han, Z. Zeng, K. Li, C. Chen, Z. Liu, M. Sun, and T. Yang. Conpet: Continual parameter-efficient tuning for large language models, 2023.
- [218] D. Song, H. Guo, Y. Zhou, S. Xing, Y. Wang, Z. Song, W. Zhang, Q. Guo, H. Yan, X. Qiu, and D. Lin. Code needs comments: Enhancing code llms with comment augmentation, 2024.
- [219] Z. Su, J. Li, Z. Zhang, Z. Zhou, and M. Zhang. Efficient continue training of temporal language model with structural information. In H. Bouamor, J. Pino, and K. Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6315–6329, Singapore, Dec. 2023. Association for Computational Linguistics.

- [220] Q. Sun, Z. Chen, F. Xu, K. Cheng, C. Ma, Z. Yin, J. Wang, C. Han, R. Zhu, S. Yuan, Q. Guo, X. Qiu, P. Yin, X. Li, F. Yuan, L. Kong, X. Li, and Z. Wu. A survey of neural code intelligence: Paradigms, advances and beyond, 2024.
- [221] Y. Sun, S. Wang, Y. Li, S. Feng, H. Tian, H. Wu, and H. Wang. Ernie 2.0: A continual pre-training framework for language understanding. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8968–8975, Apr. 2020.
- [222] K. Takahashi, T. Omi, K. Arima, and T. Ishigaki. Pretraining and updating language-and domain-specific large language model: A case study in japanese business domain. *arXiv preprint arXiv:2404.08262*, 2024.
- [223] M. Tao, Y. Feng, and D. Zhao. Can bert refrain from forgetting on sequential tasks? a probing study. In *The Eleventh International Conference on Learning Representations*, 2022.
- [224] D.-A. Team. Deepseek llm: Scaling open-source language models with longtermism, 2024.
- [225] G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [226] S. Team. Starcoder: may the source be with you!, 2023.
- [227] S. Team. Starcoder 2 and the stack v2: The next generation, 2024.
- [228] J. Thorne, A. Vlachos, C. Christodoulopoulos, and A. Mittal. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*, 2018.
- [229] D. Thulke, Y. Gao, P. Pelsler, R. Brune, R. Jalota, F. Fok, M. Ramos, I. van Wyk, A. Nasir, H. Goldstein, et al. Climategpt: Towards ai synthesizing interdisciplinary research on climate change. *arXiv preprint arXiv:2401.09646*, 2024.
- [230] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [231] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [232] G. M. Van de Ven, T. Tuytelaars, and A. S. Tolias. Three types of incremental learning. *Nature Machine Intelligence*, 4(12):1185–1197, 2022.
- [233] E. Verwimp, R. Aljundi, S. Ben-David, M. Bethge, A. Cossu, A. Gepperth, T. L. Hayes, E. Hüllermeier, C. Kanan, D. Kudithipudi, C. H. Lampert, M. Mundt, R. Pascanu, A. Popescu, A. S. Tolias, J. van de Weijer, B. Liu, V. Lomonaco, T. Tuytelaars, and G. M. van de Ven. Continual learning: Applications and the road forward, 2024.
- [234] M. Völske, M. Potthast, S. Syed, and B. Stein. Tl; dr: Mining reddit to learn automatic summarization. In *Proceedings of the Workshop on New Frontiers in Summarization*, pages 59–63, 2017.
- [235] B. Wang and A. Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [236] L. Wang, X. Zhang, Q. Li, J. Zhu, and Y. Zhong. Coscl: Cooperation of small continual learners is stronger than a big one. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVI*, pages 254–271. Springer, 2022.
- [237] L. Wang, X. Zhang, H. Su, and J. Zhu. A comprehensive survey of continual learning: Theory, method and application. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–20, 2024.
- [238] P. Wang, Z. Li, N. Zhang, Z. Xu, Y. Yao, Y. Jiang, P. Xie, F. Huang, and H. Chen. Wise: Rethinking the knowledge memory for lifelong model editing of large language models. *arXiv preprint arXiv:2405.14768*, 2024.
- [239] R. Wang, D. Tang, N. Duan, Z. Wei, X. Huang, J. Ji, G. Cao, D. Jiang, and M. Zhou. K-Adapter: Infusing Knowledge into Pre-Trained Models with Adapters. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1405–1418, Online, Aug. 2021. Association for Computational Linguistics.
- [240] X. Wang, T. Chen, Q. Ge, H. Xia, R. Bao, R. Zheng, Q. Zhang, T. Gui, and X. Huang. Orthogonal subspace learning for language model continual learning. In H. Bouamor, J. Pino, and K. Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10658–10671, Singapore, Dec. 2023. Association for Computational Linguistics.
- [241] X. Wang, Y. Zhang, T. Chen, S. Gao, S. Jin, X. Yang, Z. Xi, R. Zheng, Y. Zou, T. Gui, Q. Zhang, and X. Huang. Trace: A comprehensive benchmark for continual learning in large language models, 2023.
- [242] Y. Wang, H. Le, A. D. Gotmare, N. D. Q. Bui, J. Li, and S. C. H. Hoi. Codet5+: Open code large language models for code understanding and generation, 2023.
- [243] Y. Wang, Y. Liu, C. Shi, H. Li, C. Chen, H. Lu, and Y. Yang. Insl: A data-efficient continual learning paradigm for fine-tuning large language models with instructions, 2024.
- [244] Y. Wang, S. Mishra, P. Alipoormolabashi, Y. Kordi, A. Mirzaei, A. Arunkumar, A. Ashok, A. S. Dhanasekaran, A. Naik, D. Stap, E. Pathak, G. Karamanolakis, H. G. Lai, I. Purohit, I. Mondal, J. Anderson, K. Kuznia, K. Doshi, M. Patel, K. K. Pal, M. Moradshahi, M. Parmar, M. Purohit, N. Varshney, P. R. Kaza, P. Verma, R. S. Puri, R. Karia, S. K. Sampat, S. Doshi, S. Mishra, S. Reddy, S. Patro, T. Dixit, X. Shen, C. Baral, Y. Choi, N. A. Smith, H. Hajishirzi, and D. Khashabi. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks, 2022.
- [245] Y. Wang, W. Wang, S. Joty, and S. C. Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *EMNLP*, 2021.
- [246] Z. Wang, C.-L. Li, V. Perot, L. T. Le, J. Miao, Z. Zhang, C.-Y. Lee, and T. Pfister. Codeclm: Aligning language models with tailored synthetic data. *arXiv preprint arXiv:2404.05875*, 2024.

- [247] Z. Wang, Z. Zhan, Y. Gong, G. Yuan, W. Niu, T. Jian, B. Ren, S. Ioannidis, Y. Wang, and J. Dy. Sparcl: Sparse continual learning on the edge. *Advances in Neural Information Processing Systems*, 35:20366–20380, 2022.
- [248] Z. Wang, Z. Zhang, S. Ebrahimi, R. Sun, H. Zhang, C.-Y. Lee, X. Ren, G. Su, V. Perot, J. Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. *European Conference on Computer Vision*, 2022.
- [249] Z. Wang, Z. Zhang, C.-Y. Lee, H. Zhang, R. Sun, X. Ren, G. Su, V. Perot, J. Dy, and T. Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022.
- [250] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*, 2021.
- [251] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le. Finetuned language models are zero-shot learners, 2022.
- [252] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [253] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [254] M. Weyssow, X. Zhou, K. Kim, D. Lo, and H. Sahraoui. On the usage of continual learning for out-of-distribution generalization in pre-trained language models of code. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2023*, page 1470–1482, New York, NY, USA, 2023. Association for Computing Machinery.
- [255] G. Winata, L. Xie, K. Radhakrishnan, S. Wu, X. Jin, P. Cheng, M. Kulkarni, and D. Preotiuc-Pietro. Overcoming catastrophic forgetting in massively multilingual continual learning. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 768–777, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [256] M. Wistuba, P. T. Sivaprasad, L. Balles, and G. Zappella. Continual learning with low rank adaptation. In *NeurIPS 2023 Workshop on Distribution Shifts (DistShifts)*, 2023.
- [257] C. Wu, Y. Gan, Y. Ge, Z. Lu, J. Wang, Y. Feng, P. Luo, and Y. Shan. Llama pro: Progressive llama with block expansion, 2024.
- [258] C. Wu, W. Lin, X. Zhang, Y. Zhang, Y. Wang, and W. Xie. Pmc-llama: Towards building open-source language models for medicine. *arXiv preprint arXiv:2305.10415*, 6, 2023.
- [259] S. Wu, O. Irsoy, S. Lu, V. Dabrovolski, M. Dredze, S. Gehrmann, P. Kambadur, D. S. Rosenberg, and G. Mann. Bloomberggpt: A large language model for finance. *CoRR*, abs/2303.17564, 2023.
- [260] T. Wu, M. Caccia, Z. Li, Y.-F. Li, G. Qi, and G. Haffari. Pretrained language model in continual learning: A comparative study. In *International conference on learning representations*, 2021.
- [261] T. Wu, L. Luo, Y.-F. Li, S. Pan, T.-T. Vu, and G. Haffari. Continual learning for large language models: A survey, 2024.
- [262] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu. Large scale incremental learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 374–382, 2019.
- [263] Y. Wu, G. Wayne, A. Graves, and T. Lillicrap. The kanerva machine: A generative distributed memory. *arXiv preprint arXiv:1804.01756*, 2018.
- [264] J. Xie, Y. Liang, J. Liu, Y. Xiao, B. Wu, and S. Ni. Quert: Continual pre-training of language model for query understanding in travel domain search. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, page 5282–5291, New York, NY, USA, 2023. Association for Computing Machinery.
- [265] Q. Xie, Q. Chen, A. Chen, C. Peng, Y. Hu, F. Lin, X. Peng, J. Huang, J. Zhang, V. Keloth, et al. Me llama: Foundation large language models for medical applications. *arXiv preprint arXiv:2402.12749*, 2024.
- [266] Q. Xie, W. Han, X. Zhang, Y. Lai, M. Peng, A. Lopez-Lira, and J. Huang. PIXIU: A large language model, instruction data and evaluation benchmark for finance. *CoRR*, abs/2306.05443, 2023.
- [267] S. M. Xie, S. Santurkar, T. Ma, and P. S. Liang. Data selection for language models via importance resampling. *Advances in Neural Information Processing Systems*, 36, 2024.
- [268] Y. Xie, K. Aggarwal, and A. Ahmad. Efficient continual pre-training for building domain specific large language models, 2023.
- [269] C. Xu, Q. Sun, K. Zheng, X. Geng, P. Zhao, J. Feng, C. Tao, and D. Jiang. Wizardlm: Empowering large language models to follow complex instructions. *arXiv preprint arXiv:2304.12244*, 2023.
- [270] H. Xu, B. Liu, L. Shu, and P. S. Yu. Bert post-training for review reading comprehension and aspect-based sentiment analysis, 2019.
- [271] S. Xue, F. Zhou, Y. Xu, H. Zhao, S. Xie, Q. Dai, C. Jiang, J. Zhang, J. Zhou, D. Xiu, and H. Mei. Weaverbird: Empowering financial decision-making with large language model, knowledge base, and search engine. *CoRR*, abs/2308.05361, 2023.
- [272] Y. Yan, K. Xue, X. Shi, Q. Ye, J. Liu, and T. Ruan. Af adapter: Continual pretraining for building chinese biomedical language model. In *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 953–957, Los Alamitos, CA, USA, dec 2023. IEEE Computer Society.
- [273] S. Yang, M. A. Ali, C.-L. Wang, L. Hu, and D. Wang. Moral: Moe augmented lora for llms’ lifelong learning, 2024.
- [274] X. Yang, J. Gao, W. Xue, and E. Alexandersson. Pllama: An open-source large language model for plant science. *CoRR*, abs/2401.01600, 2024.
- [275] Y. Yang, M. Jones, M. C. Mozer, and M. Ren. Reawakening knowledge: Anticipatory recovery from catastrophic interference via structured training, 2024.
- [276] Y. Yang, J. Zhou, X. Ding, T. Huai, S. Liu, Q. Chen, L. He, and Y. Xie. Recent advances of foundation language models-based continual learning: A survey. *arXiv preprint arXiv:2405.18653*, 2024.

- [277] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [278] Ç. Yildiz, N. K. Ravichandran, P. Punia, M. Bethge, and B. Ermis. Investigating continual pretraining in large language models: Insights and implications. *arXiv preprint arXiv:2402.17400*, 2024.
- [279] W. Yin, J. Li, and C. Xiong. ConTinTin: Continual learning from task instructions. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3062–3072, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [280] L. Yu, Q. Chen, J. Zhou, and L. He. Melo: Enhancing model editing with neuron-indexed dynamic lora. *arXiv preprint arXiv:2312.11795*, 2023.
- [281] W. Yuan, Q. Zhang, T. He, C. Fang, N. Q. V. Hung, X. Hao, and H. Yin. Circle: continual repair across programming languages. In *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis, ISSTA 2022*, page 678–690, New York, NY, USA, 2022. Association for Computing Machinery.
- [282] X. Yue, X. Qu, G. Zhang, Y. Fu, W. Huang, H. Sun, Y. Su, and W. Chen. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*, 2023.
- [283] R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, and Y. Choi. Defending against neural fake news. *Advances in neural information processing systems*, 32, 2019.
- [284] Y. Zhai, S. Tong, X. Li, M. Cai, Q. Qu, Y. J. Lee, and Y. Ma. Investigating the catastrophic forgetting in multimodal large language models, 2023.
- [285] D. Zhang, X. Hu, S. Zhoubian, Z. Du, K. Yang, Z. Wang, Y. Yue, Y. Dong, and J. Tang. Sciglm: Training scientific language models with self-reflective instruction annotation and tuning. *CoRR*, abs/2401.07950, 2024.
- [286] H. Zhang, L. Gui, Y. Zhai, H. Wang, Y. Lei, and R. Xu. Copf: Continual learning human preference through optimal policy fitting. *arXiv preprint arXiv:2310.15694*, 2023.
- [287] H. Zhang, Y. Lei, L. Gui, M. Yang, Y. He, H. Wang, and R. Xu. Cppo: Continual learning for reinforcement learning with human feedback.
- [288] S. Zhang, L. Dong, X. Li, S. Zhang, X. Sun, S. Wang, J. Li, R. Hu, T. Zhang, F. Wu, and G. Wang. Instruction tuning for large language models: A survey, 2024.
- [289] X. Zhang and Q. Yang. Xuanyuan 2.0: A large chinese financial chat model with hundreds of billions parameters. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM '23*, page 4435–4439, New York, NY, USA, 2023. Association for Computing Machinery.
- [290] X. Zhang, J. Zhao, and Y. LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [291] Y. Zhang, X. Wang, and D. Yang. Continual sequence generation with adaptive compositional modules. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3653–3667, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [292] Z. Zhang, M. Fang, L. Chen, and M.-R. Namazi-Rad. CITB: A benchmark for continual instruction tuning. In H. Bouamor, J. Pino, and K. Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9443–9455, Singapore, Dec. 2023. Association for Computational Linguistics.
- [293] C. Zhao, Y. Li, and C. Caragea. C-STANCE: A large dataset for Chinese zero-shot stance detection. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13369–13385, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [294] H. Zhao, H. Han, J. Shi, C. Du, J. Liang, and Y. Xiao. Large language model can continue evolving from mistakes. *arXiv preprint arXiv:2404.08707*, 2024.
- [295] H. Zhao, H. Wang, Y. Fu, F. Wu, and X. Li. Memory-efficient class-incremental learning for image classification. *IEEE Transactions on Neural Networks and Learning Systems*, 33(10):5966–5977, 2022.
- [296] S. Zhao, X. Zou, T. Yu, and H. Xu. Reconstruct before query: Continual missing modality learning with decomposed prompt collaboration, 2024.
- [297] W. Zhao, S. Wang, Y. Hu, Y. Zhao, B. Qin, X. Zhang, Q. Yang, D. Xu, and W. Che. Sapt: A shared attention framework for parameter-efficient continual learning of large language models, 2024.
- [298] J. Zheng, Q. Ma, Z. Liu, B. Wu, and H. Feng. Beyond anti-forgetting: Multimodal continual instruction tuning with positive forward transfer, 2024.
- [299] J. Zheng, S. Qiu, and Q. Ma. Learn or recall? revisiting incremental learning with pre-trained language models, 2023.
- [300] Z. Zheng, M. Ma, K. Wang, Z. Qin, X. Yue, and Y. You. Preventing zero-shot transfer degradation in continual learning of vision-language models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19125–19136, 2023.
- [301] Z. Zheng, J. Zhang, T. Vu, S. Diao, Y. H. W. Tim, and S. Yeung. Marinegpt: Unlocking secrets of ocean to the public. *CoRR*, abs/2310.13596, 2023.
- [302] B. Zhou, D. Khashabi, Q. Ning, and D. Roth. "going on a vacation" takes longer than "going for a walk": A study of temporal commonsense understanding, 2019.
- [303] W. Zhou, D.-H. Lee, R. K. Selvam, S. Lee, B. Y. Lin, and X. Ren. Pre-training text-to-text transformers for concept-centric common sense. 2021.
- [304] D. Zhu, Z. Sun, Z. Li, T. Shen, K. Yan, S. Ding, K. Kuang, and C. Wu. Model tailor: Mitigating catastrophic forgetting in multi-modal large language models, 2024.

SUPPLEMENTARY MATERIAL

A PRELIMINARIES

In this section, we provide an overview of the fundamental concepts of large language models (LLMs) and continual learning (CL). We begin by introducing the notation used in this paper. Subsequently, we discuss the pre-training and downstream adaptation of LLMs, as well as mainstream LLM families (Appendix A.1), followed by an introduction to basic continual learning techniques studied by the community (Appendix A.2).

Notation. We denote scalars with lowercase letters, vectors with lowercase boldface letters, and matrices with uppercase boldface letters. The l_2 -norm of vectors and the Frobenius norm of a matrix are represented by $\|\cdot\|_2$. For a vector $\boldsymbol{v} = [v_1, v_2, \dots, v_n]^\top$, $\|\boldsymbol{v}\|_2 = (\sum_{i=1}^n v_i^2)^{1/2}$; for a matrix $\boldsymbol{A} \in \mathbb{R}^{m \times n}$, $\|\boldsymbol{A}\|_2 = (\sum_{ij} A_{ij}^2)^{1/2}$. We use $\epsilon_{\mathcal{D}}$, $\mathcal{L}_{\mathcal{D}}$ to denote the error function, and loss function that is deployed for training, respectively, where the subscript is used to denote the error/loss measured by taking the expectation on the data distribution \mathcal{D} . We further use $\widehat{\mathcal{L}}_S$ to represent the empirical evaluation of the loss function \mathcal{L} over the set of examples S . Probability and expectation are denoted by P and \mathbb{E} , respectively. We use $[m]$ to denote the set of positive integers up to m , $\{1, \dots, m\}$.

A.1 Large Language Models

In the past two decades, neural language modeling has emerged as the dominant field of deep learning, marked by significant and rapid advancements. Primarily built on the transformer architecture, pre-trained language models (PLMs) like BERT have established a universal hidden embedding space through extensive pre-training on large-scale unlabeled text corpora. Following the pre-training and fine-tuning paradigms, PLMs exhibit promising performance across various natural language processing tasks after being fine-tuned upon small amounts of task-specific data [51, 133, 189]. Research on scaling laws indicates that increasing model size enhances the capacity of language models [84, 102]. By scaling parameters to billions or even hundreds of billions and training on massive text datasets, PLMs not only demonstrate superior language understanding and generation capabilities but also manifest emergent abilities such as in-context learning, instruction following, and multi-step reasoning, which are absent in small-scale language models like BERT [159, 250, 252, 253, 277]. These larger models are commonly referred to as Large Language Models (LLMs).

A.1.1 Pre-Training of LLMs. Pre-training is essential for language models to acquire broad language representations. Decoder-only models typically employ probability language modeling (LM) tasks during pre-training. LM, in this context, specifically refers to auto-regressive LM. Given a sequence of tokens $\boldsymbol{x} = [x_1, x_2, \dots, x_N]$, LM predicts the next token x_t autoregressively based on all preceding tokens $\boldsymbol{x}_{<t} = [x_1, x_2, \dots, x_{t-1}]$, and trains the entire network by minimizing the negative log-likelihood:

$$\mathcal{L}_{\text{LM}}(\boldsymbol{x}) \triangleq - \sum_{t=1}^N \log P(x_t | \boldsymbol{x}_{<t}), \quad (2)$$

where $P(x_1 | \boldsymbol{x}_{<1}) \triangleq P(x_1)$ is the unconditional probability estimation of the first token. The three most popular families of decoder-only models are GPT, PaLM, and LLaMA. The GPT family, developed by OpenAI, includes models such as GPT-2 [186], GPT-3 [22], ChatGPT [173], and GPT-4 [1]. Notably, GPT-3 was the first LLM to exhibit emergent abilities not found in smaller PLMs. Another notable family, Gemini, developed by Google, is comparable to the GPT family [194, 225]. While both GPT and Gemini families are closed-source, LLaMA, released by Meta, is currently the most popular open-source family of LLMs [230, 231]. The weights of these models are made available to the research community under non-commercial licenses.

Masked language modeling (MLM) task serves as a common pre-training objective for encoder-only models like BERT [51, 133]. In MLM, for the input sequence \mathbf{x} , a subset of input tokens $m(\mathbf{x})$ are masked and replaced with the special [MASK] token. The pre-training goal is to utilize the unmasked parts $\mathbf{x}_{\setminus m(\mathbf{x})}$ to predict the masked portions $m(\mathbf{x})$. In summary, the overarching goal of MLM is to minimize the negative log-likelihood:

$$\mathcal{L}_{\text{MLM}}(\mathbf{x}) \triangleq - \sum_{\hat{\mathbf{x}} \in m(\mathbf{x})} \log P(\hat{\mathbf{x}} | \mathbf{x}_{\setminus m(\mathbf{x})}). \quad (3)$$

Some encoder-decoder architecture models, such as T5 [189], also utilize Sequence-to-Sequence MLM task as the pre-training objective. They take masked sentences as encoder inputs and utilize the decoder to sequentially predict the masked tokens.

A.1.2 Adaptation of LLMs. LLMs are primarily trained to generate linguistically coherent text. However, this training may not align with human values, preferences, or practical needs. Furthermore, the pre-training data can be outdated, leading to knowledge cutoffs or inaccuracies. To address these issues, various computational paradigms such as Instruction Tuning (IT) [288], Model Refinement (MR) [47], and Model Alignment (MA) [174, 187] have been proposed. These approaches adapt LLMs to better meet diverse downstream tasks and user requirements.

DEFINITION A.1 (INSTRUCTION TUNING, IT). Let $h(\mathbf{x})$ be a language model that takes as input data \mathbf{x} , typically consisting of natural language instructions or queries. Instruction Tuning (IT) is a specialized training approach designed to enhance the model’s ability to accurately and effectively respond to specific instructions. The objective of IT is to refine h by adjusting its parameters using a designated set of training examples $\mathcal{I} = \{(\mathbf{x}_i, \hat{\mathbf{y}}_i)\}_{i=1}^N$ drawn from the IT data distribution $\mathcal{D}_{\mathcal{I}}$, where $\hat{\mathbf{y}}_i$ represents the desired output for \mathbf{x} . This set is curated to target specific tasks or functionalities that require improved performance. Formally, IT seeks to find an optimal refined hypothesis h^* that satisfies:

$$h^* \triangleq \arg \min_{h'} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{\mathcal{I}}} [-\log P(\hat{\mathbf{y}} | \mathbf{x}, h')] \approx \arg \min_{h'} \sum_{i=1}^N -\log P(\hat{\mathbf{y}}_i | \mathbf{x}_i, h'). \quad (4)$$

REMARK. The task of Model Alignment (MA) is usually formulated in the same problem definition as IT, with an alignment dataset of size M as $\mathcal{A} = \{(\mathbf{x}_a, \mathbf{y}_a, \hat{\mathbf{y}}_a)\}_{a=1}^M$, where \mathbf{y}_a represents the model’s original decision for input \mathbf{x}_a , and $\hat{\mathbf{y}}_a$ denotes the aligned decision that adheres to specified ethical guidelines or desired outcomes.

DEFINITION A.2 (MODEL REFINEMENT, MR). Suppose we have a model $h(\mathbf{x})$ taking data \mathbf{x} (e.g., natural language queries) as inputs. Consider a size- N editing set $\mathcal{E} = \{(\mathbf{x}_e, \mathbf{y}_e, \hat{\mathbf{y}}_e)\}_{e=1}^N$, where $\hat{\mathbf{y}}_e$ denotes the true label of \mathbf{x}_e , but the model incorrectly outputs \mathbf{y}_e for \mathbf{x}_e . Model Refinement (MR) aims to efficiently update the model from h to h' such that it correctly predicts the editing set \mathcal{E} , while preserving the original outputs outside \mathcal{E} . Formally, we aim to find h' satisfying

$$h'(\mathbf{x}_0) = \begin{cases} \hat{\mathbf{y}}_0 & \text{if } (\mathbf{x}_0, \hat{\mathbf{y}}_0) \in \mathcal{E}, \\ h(\mathbf{x}_0) & \text{o.w.} \end{cases} \quad (5)$$

A.2 Continual Learning

Humans gradually accumulate knowledge and skills across tasks without significant performance decline on previous tasks [101, 153, 154, 175]. In contrast, machine learning models are usually data-centric, minimizing the training loss on the subsequent tasks will cause the model fail on the old ones, which phenomenon is phrased as “catastrophic forgetting”. Addressing this challenge is a focal point in continual learning research. The problem of efficiently adapting models to

a sequence of tasks without forgetting is extensively studied in the continual learning community [38, 178, 232, 237]. These studies are typically conducted under the following memory constraint of CL.

DEFINITION A.3 (MEMORY CONSTRAINT OF CONTINUAL LEARNING). Suppose T sets of observations $\{S_t \sim \mathcal{T}_t\}_{t=1}^T$ come in as a sequence, where $\{\mathcal{T}_t\}_{t=1}^T$ denotes the T task distributions. At the learning stage $t > 1$, the sets of observations $\{S_i\}_{i=1}^{t-1}$ are not accessible (**strong**) or only partially accessible (**relaxed**).

REMARK. In early stages of CL, works mostly focused on the strong memory constraint [4, 113, 124, 135]; as the research field progresses, more focus was put on relaxing the memory constraint to a small buffer for replay [24, 30, 193, 213]; some modern CL works completely discard the memory constraint but put focus on the computational budget [181, 233].

A.2.1 Three Types of Continual Learning. There are three outstanding types of continual learning scenarios: task-incremental learning (TIL), domain-incremental learning (DIL), and class-incremental learning (CIL). To establish a groundwork for subsequent discussions (as illustrated in Table 3 and Section 6.2), we adhere to the conceptual framework proposed by [112, 232, 237] and offer formal definitions for these three continual learning scenarios.

DEFINITION A.4 (TASK-INCREMENTAL LEARNING, TIL). Suppose T task distributions $\{\mathcal{T}_t\}_{t=1}^T$ come in as a sequence, where \mathcal{T}_t denotes the joint distribution over the t -th task's input space and the label space $(\mathcal{X}_t, \mathcal{Y}_t)$. Denote $\mathcal{X} \triangleq \bigcup_{t=1}^T \mathcal{X}_t$ and $\mathcal{Y} \triangleq \bigcup_{t=1}^T \mathcal{Y}_t$ as the union of the input and label spaces, respectively. Under the memory constraint defined in Definition A.3, Task-Incremental Learning (TIL) aims to find the optimal hypothesis $h^* : \mathcal{X} \times [T] \rightarrow \mathcal{Y}$ that satisfies:

$$h^* = \arg \min_h \sum_{t=1}^T \mathbb{E}_{(x,y) \sim \mathcal{T}_t} [\mathbb{1}_{h(x,t) \neq y}]. \quad (6)$$

DEFINITION A.5 (DOMAIN-INCREMENTAL LEARNING, DIL). Suppose T domain distributions $\{\mathcal{D}_t\}_{t=1}^T$ come in as a sequence, where \mathcal{D}_t denotes the t -th joint distribution over the shared input space and label space $(\mathcal{X}, \mathcal{Y})$. Under the memory constraint defined in Definition A.3, Domain-Incremental Learning (DIL) aims to find the optimal hypothesis $h^* : \mathcal{X} \rightarrow \mathcal{Y}$ that satisfies:

$$h^* = \arg \min_h \sum_{t=1}^T \mathbb{E}_{(x,y) \sim \mathcal{D}_t} [\mathbb{1}_{h(x) \neq y}]. \quad (7)$$

DEFINITION A.6 (CLASS-INCREMENTAL LEARNING, CIL). Suppose T task distributions $\{\mathcal{T}_t\}_{t=1}^T$ come in as a sequence, where \mathcal{T}_t denotes the joint distribution over the t -th task's input space and the label space $(\mathcal{X}_t, \mathcal{Y}_t)$. Denote $\mathcal{X} \triangleq \bigcup_{t=1}^T \mathcal{X}_t$ and $\mathcal{Y} \triangleq \bigcup_{t=1}^T \mathcal{Y}_t$ as the union of the input and label spaces, respectively. Under the memory constraint defined in Definition A.3, Class-Incremental Learning (CIL) aims to find the optimal hypothesis $h^* : \mathcal{X} \rightarrow [T] \times \mathcal{Y}$ that satisfies:

$$h^* = \arg \min_h \sum_{t=1}^T \mathbb{E}_{(x,y) \sim \mathcal{T}_t} [\mathbb{1}_{h(x) \neq (t,y)}]. \quad (8)$$

REMARK. In TIL, it is common to have a shared input space $\mathcal{X} = \mathcal{X}_t, \forall t \in [T]$, but the space of the label distribution \mathcal{Y}_t can be distinct ($\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset, \forall i \neq j$), partially shared ($\mathcal{Y}_i \cap \mathcal{Y}_j \neq \emptyset, \exists i \neq j$), or shared across different tasks ($\mathcal{Y} = \mathcal{Y}_t, \forall t \in [T]$). In DIL, the tasks are defined in the same format, i.e., same input space \mathcal{X} and same output space \mathcal{Y} . During the inference, no task IDs are provided for the hypothesis, which means the continual learning model needs to capture the pattern between the domain-invariant features and the labels. DIL is commonly perceived as more difficult than TIL. CIL is commonly viewed as the most challenging continual learning scenario, as the model needs to infer the label and the task ID at the same time. Another possible formulation of CIL is to represent it as DIL but the output label spaces are disjoint, $\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset, \forall i \neq j$.

A.2.2 Techniques of Continual Learning. The objective of CL is to find a hypothesis that minimizes risk across all tasks/domains. Consider DIL as an example [213], at t -th learning stage, the ideal training objective $\mathcal{L}(h)$ is defined as

$$\mathcal{L}(h) \triangleq \underbrace{\sum_{i=1}^{t-1} \mathcal{L}_{\mathcal{D}_i}(h)}_{\text{past domains}} + \underbrace{\mathcal{L}_{\mathcal{D}_t}(h)}_{\text{current domain}} . \quad (9)$$

The objectives for past domains are often challenging to measure or optimize due to the memory constraints (Definition A.3). Therefore, the core of designing CL algorithms lies in identifying a proxy learning objective for the first term without violating the memory constraint. Existing CL techniques can be roughly categorized into 5 groups: (i) replay-based, (ii) regularization-based, (iii) architecture-based, (iv) optimization-based, and (v) representation-based [237]. Here, we provide a concise yet comprehensive introduction to the first three categories of continual learning techniques, as they find extensive application in continually learning large language models.

Replay-Based Methods. Replay-based methods adopt the relaxed memory constraint by keeping a small buffer of observed data $\{M_i\}_{i=1}^{t-1}$ for each task \mathcal{T}_i . Formally, they seek to optimize the following empirical training objective:

$$\widehat{\mathcal{L}}_{\text{replay}}(h) \triangleq \underbrace{\sum_{i=1}^{t-1} \widehat{\mathcal{L}}_{M_i}(h)}_{\text{proxy for past domains}} + \underbrace{\widehat{\mathcal{L}}_{S_t}(h)}_{\text{current domain}} , \quad (10)$$

where $\widehat{\mathcal{L}}_S$ denotes the empirical loss term evaluated on the set of examples S . Often regarded as a simplistic solution to CL, replay-based methods may theoretically lead to loose generalization bounds [213]. Despite this, they are valued for their simplicity, stability, and high performance, even with a small episodic memory [30, 195]. For instance, DER++ [24] demonstrates consistent performance enhancement by replaying a small set of past examples along with their logits (known as dark experience replay). ESM-ER [204] introduces error sensitivity modulation (ESM) to mitigate abrupt representational drift caused by high-error new examples. A significant focus in replay-based CL is enhancing sample efficiency for buffer maintenance. For instance, [193] prioritizes exemplar selection based on herding to accurately model class mean throughout class-incremental learning. [295] propose storing low-fidelity examples to achieve memory-efficient exemplar set maintenance.

Regularization-Based Methods. Suppose $h_{\theta_{t-1}}$ is the hypothesis yielded after the $t-1$ -th stage of training, parameterized by θ_{t-1} . Regularization-based methods utilize a regularization term as a proxy for past domain losses, determined by the distance in the parameter space.

$$\widehat{\mathcal{L}}_{\text{reg}}(h_\theta) \triangleq \underbrace{\lambda \cdot \|\theta - \theta_{t-1}\|_\Sigma}_{\text{proxy for past domains}} + \underbrace{\widehat{\mathcal{L}}_{S_t}(h_\theta)}_{\text{current domain}} , \quad (11)$$

where $\|v\|_\Sigma = v^\top \Sigma v$ is the vector norm evaluated on a positive-semi-definite matrix Σ , and λ is the regularization coefficient, a hyper-parameter introduced to balance the past knowledge retention and current knowledge learning. The matrix Σ introduced is to measure the different level of importance of each parameters and their correlations in retaining the past knowledge. In practice, to reduce computational overhead, diagonal matrices are often designed to encode only the importance of each parameter. For example, Elastic Weight Consolidation (EWC) [113] adopts a Bayesian perspective, using diagonal values from the Fisher Information Matrix (FIM) as an approximation for the Hessian matrix of parameters. This forms a sequential Maximize A Posteriori (MAP) optimization for continual learning.

Memory Aware Synapses (MAS) [4] computes parameter importance in an online and unsupervised manner, defining importance by accumulated absolute gradient during training. It is also worth noting that when $\Sigma = I$ degenerates to an identity matrix, the regularization term simplifies to a basic l_2 -penalty term, equally penalizing each parameter, which can be surprisingly effective in some cases of continual LLMs [197].

Architecture-Based Methods. Expanding the network architecture dynamically to assimilate new knowledge is deemed the most efficient form of CL [248, 249]. This method primarily tackles adaptation challenges and can achieve zero-forgetting when task IDs are available during inference or can be correctly inferred [71, 256]. However, due to the difficulty of task ID inference, architecture expansion is predominantly utilized in TIL but is scarcely explored in DIL or CIL. Progressive Neural Networks (PNN) [201] proposes learning laterally connected neurons as new tasks arise, ensuring non-forgetting and enabling transfer of previously learned neurons for future tasks. In conjunction with pre-trained backbone large models like ViT [54], CoLoR [256] trains various low-rank adaptation (LoRA) [86] modules for different tasks. It estimates and stores prototypes for each task and utilizes the natural clustering ability of the pre-trained model during testing to infer task IDs, selecting the corresponding LoRA component for prediction generation. In the domain of continual LLMs, architecture expansion has resurged in popularity following the rise of parameter-efficient fine-tuning (PEFT) [50, 86, 211], a topic we will delve into shortly [96, 100, 118, 177, 240, 257, 272, 273].

B EVALUATION PROTOCOLS AND DATASETS

In Appendix B.1, we review common continual learning evaluation metrics and provide formal definitions. In Appendix B.2, we introduce metrics designed specifically for continual LLMs. Finally, in Appendix B.3, we outline the datasets available for each discussed topic.

B.1 Evaluation Metrics of Continual Learning

In the realm of conventional continual learning, where task streams take the form of classification, many metrics rely on the concept of Accuracy Matrix [136, 213]. Extending this notion to the context of continually learning LLMs, we introduce the **Performance Matrix** $P \in \mathbb{R}^{T \times T}$, where T represents the total number of training stages. Each entry of P corresponds to a performance metric evaluated on the models, such as perplexity on pre-training data [35, 69, 100], zero-shot/few-shot evaluation metrics on downstream data without fine-tuning [9, 42, 48, 172, 199, 258], fine-tuned accuracies on downstream tasks [5, 35, 96, 183], and probing accuracies from fine-tuning add-on components evaluated on downstream tasks [144, 223, 299]. In P , $P_{i,j}$ denotes the model’s performance after training on task i and evaluating on task j . With this Performance Matrix definition, we introduce the primary evaluation protocols widely adopted.

Overall Performance (OP). The Overall Performance (OP) [106, 286, 291] is a natural extension of the concept of Average Accuracy [136, 213]. The OP measured up until training stage t is the average performance of the model trained right after the stage t . Denote it as OP_t and we have:

$$OP_t \triangleq \frac{1}{t} \sum_{i=1}^t P_{t,i}. \quad (12)$$

As noted in [213], the OP corresponds to the primary optimization objective defined in Definition A.4, A.5, and A.6. In much of the continual learning literature, once all T tasks are completed, the final OP (OP_T) is reported, with the subscript T often omitted for brevity. In some works, OP is weighted by the importance of tasks $\widetilde{OP} \triangleq \frac{1}{T} \sum_{i=1}^T w_i P_{t,i}$, where $w_i = N_i / \sum_{j=1}^T N_j$ represents the ratio of data. In some literature, \widetilde{OP} is referred to as “example accuracy” [37], “whole accuracy” [217], or “edit success rate” in CMR [74].

Forgetting (F). Define F_t as the forgetting up to task t , which represents the largest performance drop observed throughout the training process, averaged over t training stages:

$$F_t \triangleq \frac{1}{t-1} \sum_{j=1}^{t-1} \left[\max_{l \in [t-1]} \{P_{l,j} - P_{t,j}\} \right]. \quad (13)$$

Typically, researchers report the average forgetting $F = F_T$ at the end of the entire training process. Forgetting quantifies the impact of learning new tasks on previously acquired knowledge. Ideally, a robust continual learning framework should achieve **Backward Transfer (BWT)**, where learning new tasks enhances performance on prior tasks. This enhancement is typically measured by negating the forgetting, thus indicating an improvement in performance on earlier tasks. The concepts of Forgetting and Backward Transfer underpin various evaluation metrics, such as knowledge retention [100], performance on unchanged knowledge [95], average increased perplexity (AP⁺) [184], and test and edit retention rate in CMR [74].

Forward Transfer (FWT). Forward Transfer measures the generalization ability of the continual learning algorithms. Formally, forward transfer FWT_t up to training stage t is defined as

$$\text{FWT}_t \triangleq \frac{1}{t-1} \sum_{i=2}^t P_{i-1,i} - b_i, \quad (14)$$

where b_i is the baseline performance of the model evaluated on task i before undergoing continual learning. Strictly speaking, the definition of b_i is not the same as defined in the previous work [136, 213], where it is used to denote the performance of a random initialization of the model. Additionally, we extend the notation of forward transfer in the vertical direction to represent the performance improvement on downstream tasks resulting from domain-adaptive pre-training (see Table 2). Forward Transfer is alternatively referred to as temporal generalization [100] or knowledge transfer [116] in some literature. In this section, we introduce the evaluation protocols and datasets for continual LLMs.

B.2 Continual LLMs' Evaluation Protocols

Language Model Analysis (LAMA). Language Model Analysis (LAMA) is an evaluation framework designed to *probe the world knowledge* embedded in language models [179]. It converts each world fact into a cloze statement, which is then inputted into the language models to predict the correct answer. LAMA has been extended for continual pre-training, particularly for those under the temporal shifts [95, 96]. In CKL, three LAMA benchmarks are constructed for different dimensions: InvariantLAMA assesses knowledge retention on time-invariant facts, UpdatedLAMA focuses on knowledge update, and NewLAMA evaluates knowledge acquisition [96].

Forgotten / (Updated + Acquired) Ratio (FUAR). As the performance of a pre-trained LLM is decomposed into a fine-grained set in CKL [96], OP becomes a too general metric and cannot accurately reflect the balance and trade-offs of the model's behavior. To address this issue, CKL proposes a joint evaluation metric FUAR (Forgotten / (Updated + Acquired) Ratio) for continual pre-training. A FUAR value of 1 represents an equal trade-off between the knowledge forgetting and knowledge learning: for each piece of updated or acquired knowledge, one piece of time-invariant knowledge is forgotten on average. A FUAR less than 1 suggests high learning efficacy, where more than one piece of knowledge is acquired at the expense of forgetting one piece of time-invariant knowledge.

X-Delta. In TRACE [241], the authors propose a set of "X-Delta" metrics for continual instruction tuning, quantifying the forward transfer on specific abilities of LLMs. Let's denote a set of M datasets $\{X_1, X_2, \dots, X_M\}$ for task X . The baseline performances of the pre-trained LLM evaluated on these tasks are denoted as $\{b_1^X, \dots, b_M^X\}$. The model

undergoes continuous fine-tuning on a different set of tasks, distinct from those used for evaluation. Throughout the sequential training process, the performance of the model after learning task t on evaluation tasks X_i is $R_{t,i}^X$. The X-Delta ΔR_t^X after learning task t is defined as:

$$\Delta R_t^X \triangleq \frac{1}{M} \sum_{m=1}^M (R_{t,i}^X - b_i^X). \quad (15)$$

In the public TRACE benchmark, the authors construct three sets of evaluation tasks to benchmark the ability of LLMs, including *general ability*, *instruction following*, and *safety* [241].

NLG Score. In continual model alignment, three prominent metrics used to evaluate different aspects of Natural language generation (NLG) are BLEU-4 [176], METEOR [12], and ROUGE-L [126]. BLEU-4 [176], designed for machine translation (MT), evaluates the precision of n-grams between the machine-generated and reference texts, focusing especially on four-word sequences to gauge fluency and adequacy. METEOR [12] also targets MT but aims to improve correlation with human judgment by considering synonyms and stemming, thus providing a more nuanced assessment of translation quality. On the other hand, ROUGE-L [126] is commonly applied in summarization tasks, assessing the longest common subsequence between the generated summary and a set of reference summaries, effectively measuring the recall of essential content. Each metric has its strengths and is tailored to specific kinds of language processing tasks, reflecting different dimensions of text generation quality.

B.3 Datasets

In this section, we provide a comprehensive review of the datasets available for benchmarking continual LLMs, as illustrated in Table 4. We intentionally exclude datasets used for domain-adaptive pre-training LLMs in vertical domains such as legal, medical, and financial, unless they are specifically designed for continual domain-adaptive pre-training. Furthermore, we omit datasets used in general continual fine-tuning, as they have already been extensively studied in existing works [17, 105].

Datasets for Continual Pre-Training (CPT) and Domain Adaptive Pre-Training (DAP). Current research lacks a widely recognized benchmark for evaluating continual pre-training LLMs under temporal shifts. TimeLMs utilizes a series of Twitter corpora collected until 2022, sequentially pre-training RoBERTa models quarterly [137]. CC-RecentNews, adopted as unlabeled pre-training data for LMs in CKL [96], consists of recent news and serves as a single-stage dataset. Additionally, CKL introduces InvariantLAMA, NewLAMA, and UpdatedLAMA to assess the principles of continual knowledge learning. TWiki, a dataset derived from the articles of Wikipedia between August and December 2021, is curated and cleaned in TemporalWiki [95]. This dataset facilitates the exploration of incremental learning by providing the Diffsets between neighboring snapshots. For works that study the content-level distributional shifts in CPT and DAP, researchers often resort to a similar set of publicly available datasets [134, 169, 270] to construct their own test beds for continual learning algorithms. The *DAPT dataset, developed by [72], comprises four domains: BioMed and Computer Science from S2ORC [134], News from [283], and Reviews from [78]. In *DAPT’s original study, each domain undergoes individual domain adaptive pre-training stages to demonstrate the universality of DAP’s effectiveness. Subsequent works, such as ELLE [184] and Recyclable Tuning [183], follow suit by employing these domains for multi-stage CPT. DEMix [71] presents another large-scale dataset, featuring eight semantic domains with over 73.8 billion tokens. Alongside the training set, it includes eight additional datasets for validating the generalization ability of LLMs. On a smaller scale, *CPT [104] and *DAS [107] datasets consist of four and eight domains, respectively,

with approximately 3.12 million examples and a size of 4.16GB each. These datasets are constructed similarly to the aforementioned ones.

Datasets for Continual Instruction Tuning. Measuring the effectiveness of CIT is crucial, particularly because traditional evaluation metrics may not be suitable for LLMs: many of them are overly simplistic and fail to comprehensively assess the model’s ability to learn continually. New benchmarks and metrics are required to evaluate both the retention of old knowledge and the integration of new instructions. TRACE [241] stands as a continual learning benchmark designed specifically for LLMs, encompassing diverse tasks such as multilingual capabilities, code generation, and mathematical reasoning. CITB [292] represents another benchmark for CIT, incorporating both learning and evaluation protocols. It in addition demonstrates that replay generally yields the best performance across all methods. CoIN [32] extends the benchmark to MLLMs, incorporating a balanced and diverse set of instructions from vision-language datasets.

Datasets for Continual Model Refinement. Most datasets for continual model refinement can be categorized into two types [152]: fact checking and question answering. For fact checking, models are asked to verify the truthfulness of certain claims, typically modeled as a classification task. Key datasets include FEVER [228] (used by [47, 76]) and VitaminC [206] (used by [164]), both sourced from Wikipedia. For question answering, models are tasked with providing specific answers instead of choices. Zero-shot Relation Extraction (zsRE) [117] is the most widely employed dataset for this purpose [45, 74–76, 157, 158], alongside Natural Questions (NQ) [114] and T-rex [59]. [157] adapted zsRE with additional counterfactuals to create the more challenging CounterFact dataset, used by [45, 85, 280]. Beyond these two categories, SCOTUS [28] is also utilized [74] in the assessment of continual model refinement through a document classification task for U.S. Supreme Court cases into 11 topics.

Datasets for Continual Model Alignment. In the domain of reinforcement learning with human feedback (RLHF), several datasets are commonly employed across different studies to evaluate the adaptation and effectiveness of models under varying scenarios and continuous learning conditions. The IMDB [149] and HH-RLHF [11] dataset, as introduced in [286] within their study on continual learning through optimal policy fitting, leverages data gathered from interactive RL scenarios to model human preferences dynamically. Similarly, the Reddit TL;DR dataset [234] used by [286, 287] is focused on text summarization, providing a robust platform for testing the longevity and adaptability of learning algorithms under evolving conditions. Lastly, Common Sense QA [18, 41, 115], Reading Comprehension [56, 190], and Translation [19], which are utilized in [128] are selected to assess the challenges of aligning RL agents with human expectations without incurring significant performance penalties. Each of these datasets is pivotal in advancing the understanding of continual learning and the interplay between human feedback and machine learning adaptation.

Datasets for Continual Multimodal Large Language Models. Following LLaVA [131], many MLLMs adopt the pattern of instruction tuning to enable assessing alignment with human intention and knowledge preservation for reasoning. Thus, traditional tasks like image classification can be transformed to VQA tasks to evaluate the ability of MLLMs, which are otherwise challenging to assess using conventional methods. Several benchmarks have been proposed to evaluate the CL method for MLLMs. MCIT [77] proposes the first continual instruction tuning benchmarks, Benchmark1 and Benchmark2. The difference between benchmark1 and benchmark2 is that benchmark2 includes Multi-task Joint Instruction Tuning, which aims to explore whether multi-task joint instruction tuning improves the model’s continual learning ability. [284] proposes EMT, the first classification evaluation framework to investigate catastrophic forgetting in MLLMs. [32] presents a comprehensive benchmark CoIN, spanning 8 task categories and evaluating MLLMs from two perspectives: Instruction Following and General Knowledge, which assess the alignment

with human intention and knowledge preserved for reasoning, respectively. [296] constructs two datasets, UPMC-Food101-CMML and MM-IMDb-CMML to benchmark the novel CMML task, which means the data of certain modalities is missing during continual fine-tuning. UPMC-Food101-CMM contains 101 food categories and 61,142 training, 6,846 validation, and 22,716 test image-text pairs. MM-IMDb-CMML is a multi-label classification dataset across 27 distinct movie genres, consisting of 15,552 training, 2,608 validation and 7,799 test image-text pairs.