

Bring Your Own Knowledge: A Survey of Methods for LLM Knowledge Expansion

Anonymous ACL submission

Abstract

Adapting large language models (LLMs) to new and diverse knowledge is essential for their lasting effectiveness in real-world applications. This survey provides an overview of state-of-the-art methods for expanding the knowledge of LLMs, focusing on integrating various knowledge types, including factual information, domain expertise, language proficiency, and user preferences. We explore techniques, such as continual learning, model editing, and retrieval-based explicit adaptation, while discussing challenges like knowledge consistency and scalability. Designed as a guide for researchers and practitioners, this survey sheds light on opportunities for advancing LLMs as adaptable and robust knowledge systems.

1 Introduction

As large language models (LLMs) are increasingly deployed in real-world applications, their ability to adapt to evolving knowledge becomes crucial for maintaining relevance and accuracy. However, LLMs are typically trained once and thus only have knowledge up to a certain cutoff date, limiting their ability to stay updated with new information. This survey provides a comprehensive overview of methods that enable LLMs to incorporate various types of new knowledge, including factual, domain-specific, language, and user preference knowledge. We survey adaptation strategies, including continual learning, model editing, and retrieval-based approaches, and aim at providing guidelines for researchers and practitioners.

To remain effective, LLMs require updates across multiple dimensions. Factual knowledge consists of general truths and real-time information, while domain knowledge pertains to specialized fields, such as medicine or law. Language knowledge enhances multilingual capabilities, and preference knowledge aligns model behavior with user expectations and values. Ensuring that LLMs

	Continual Learning (\$4)	Model Editing (\$5)	Retrieval (\$6)
Knowledge Type			
Fact	✓	✓	✓
Domain	✓	✗	✓
Language	✓	✗	✗
Preference	✓	✓	✗
Applicability			
Large-scale data	✓	✗	✓
Precise control	✗	✓	✓
Computational cost	✗	✓	✓
Black-box applicable	✗	✗	✓

Table 1: We compare three key approaches for adapting LLMs — continual learning, model editing, and retrieval — based on their supported knowledge types and applicability across different criteria.

can integrate updates across these dimensions is essential for their sustained utility.

Existing LLM adaptation methods differ in approach and application. Continual learning enables incremental updates to models' parametric knowledge, mitigating catastrophic forgetting (McCloskey and Cohen, 1989) while ensuring long-term performance. Model editing allows for precise modifications of learned knowledge, providing controlled updates without requiring full retraining. Unlike these *implicit* knowledge expansion methods, which modify the model's internal parameters, retrieval-based approaches *explicitly* access external information dynamically during inference, reducing dependency on static parametric knowledge. The suitability of these methods for different knowledge types and their general applicability are summarized in Table 1. By leveraging these strategies, LLMs can maintain accuracy, contextual awareness, and adaptability to new information.

After placing our work into context (Section 2) and defining knowledge types covered in this paper (Section 3), we provide an overview of different knowledge expansion methods as detailed in Fig-

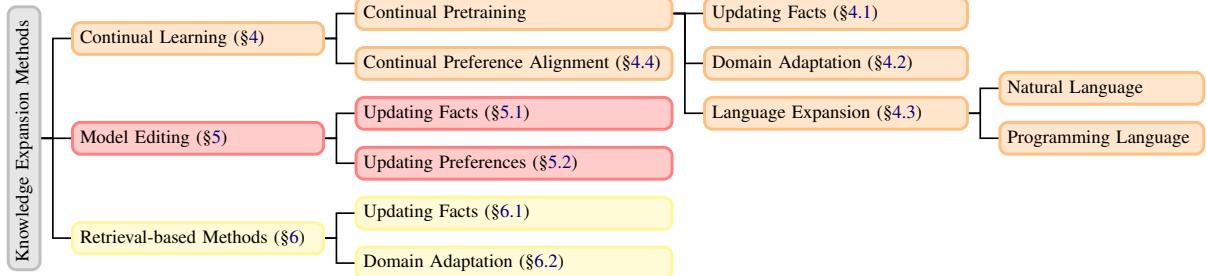


Figure 1: Taxonomy of current methods for expanding LLM knowledge. Due to space constraints, please refer to Appendix A.1 for a comprehensive review of methods and their corresponding citations.

ure 1. This work thus surveys diverse research efforts and may serve as a guide for researchers and practitioners aiming to develop and apply adaptable and robust LLMs. We highlight research opportunities and provide insights into optimizing adaptation techniques for various real-world applications.

2 Related Surveys

The main goal of our work is to provide researchers and practitioners a broad overview of various types of methods to adapt LLMs to diverse types of new knowledge. In this section, we explain how other more specialized surveys relate to our paper.

To the best of our knowledge, there is limited prior work that specifically focuses on continuous knowledge expansion for LLMs. Closest to our work, Zhang et al. (2023c) describe temporal factual knowledge updates, while we take a broader perspective by examining methods for adapting LLMs to unseen domain knowledge, expanding language coverage, and incorporating user preferences. Yao et al. (2023) and Zhang et al. (2024c) provide overviews of knowledge editing methodologies, categorizing approaches of knowledge editing. Similarly, Ke and Liu (2022), Wu et al. (2024b) and Wang et al. (2024b) offer a comprehensive overview of continual learning. In contrast, our survey shifts the focus towards a task-oriented perspective on knowledge expansion, detailing how various types of knowledge — including factual, domain-specific, language, and user preference knowledge — can be seamlessly integrated to ensure LLMs remain relevant and effective.

3 Knowledge Types

Integrating diverse types of knowledge into LLMs is essential to enhance their versatility and effectiveness. Depending on the use case, the type of knowledge that an LLM shall be adapted to, might

differ. In this paper, we distinguish four key types of knowledge, which cover a broad range of use cases of researchers and practitioners: **factual, domain, language, and preference knowledge**.

(1) We define **factual knowledge** as general truths or contextualized information about the world that can be expressed in factual statements. We adopt a broad, high-level definition, encompassing finer-grained categorizations, such as commonsense knowledge, cultural knowledge, temporal knowledge, and entity knowledge as subsets of factual knowledge, in contrast to prior works (Cao et al., 2024; Wu et al., 2024b) using more granular classifications. This inclusive perspective enables a comprehensive exploration of knowledge expansion techniques for LLMs, providing flexibility beyond predefined categories and taxonomies.

(2) We define **domain knowledge** as specialized information relevant to specific fields, such as medicine, law, or engineering, enabling LLMs to perform well in targeted applications. Since LLMs typically excel in general-domain tasks but struggle with specialized content, incorporating domain knowledge is crucial for bridging this gap and improving performance in specific fields.

(3) We define **language knowledge** as the ability of an LLM to understand, generate, and reason in specific natural or programming languages.¹ Its integration focuses on adapting models to new languages and enhancing performance in underrepresented ones for broader applicability.

(4) Finally, we define **preference knowledge** as the capability of LLMs to tailor their behavior to align with user-specific needs, preferences, or values. Preference knowledge integration involves

¹We distinguish language knowledge from linguistic knowledge as defined by Hernandez et al. (2024). Language knowledge refers to the multilingual capabilities of an LLM, whereas linguistic knowledge falls under factual knowledge, encompassing statements about syntax and grammar.

137 adapting LLM behavior to meet diverse and dy-
138 namic user expectations.
139

140 In the next sections, we survey knowledge ex-
141 pansion methods and explain for which of these
four knowledge types they are suitable.

142 4 Continual Learning

143 Continual learning (CL) is a machine learning
144 paradigm that mimics the human ability to con-
145 tinuously learn and accumulate knowledge without
146 forgetting previously acquired information (Chen
147 and Liu, 2018). In the context of knowledge ex-
148 pansion, CL allows LLMs to integrate new corpora
149 and incrementally update the knowledge stored in
150 their parameters. This ensures that LLMs remain
151 adaptable, relevant, and effective as facts, domains,
152 languages, and user preferences evolve over time.

153 In the era of LLMs, the training of language
154 models usually includes multiple stages: pretrain-
155 ing, instruction tuning, preference alignment, and
156 potentially fine-tuning on a downstream task (Shi
157 et al., 2024a). Depending on the stage, continual
158 learning can be categorized into continual pretrain-
159 ing, continual instruction tuning, continual prefer-
160 ence alignment, and continual end-task learning
161 (Ke et al., 2023; Shi et al., 2024a). For knowledge
162 expansion, the focus lies on continual pretraining
163 (CPT) and continual preference alignment (CPA).
164 In contrast, continual instruction tuning and contin-
165 ual end-task learning primarily aim to sequentially
166 fine-tune pretrained LLMs for acquiring new skills
167 and solving new tasks, which fall outside the scope
168 of this survey.

169 In the following sections, we review existing
170 studies that leverage continual pretraining for up-
171 dating facts, adapting domains, and expanding lan-
172 guages, and continual alignment for updating user
173 preferences.

174 4.1 Continual Pretraining for Updating Facts

175 This line of research focuses on updating a lan-
176 guage model’s outdated internal factual knowl-
177 edge by incrementally integrating up-to-date world
178 knowledge (Jang et al., 2022; Ke et al., 2023).

179 Early studies (Sun et al., 2020; Röttger and Pier-
180 rehumbert, 2021; Lazaridou et al., 2021; Dhingra
181 et al., 2022) empirically analyze continual pretrain-
182 ing on temporal data, demonstrating its potential
183 for integrating new factual knowledge. Jin et al.
184 (2022) and Jang et al. (2022) apply traditional con-
185 tinual learning methods to factual knowledge up-

186 dates in LLMs, evaluating their effectiveness in
187 continual knowledge acquisition. Similarly, Jang
188 et al. (2022) and Kim et al. (2024) classify world
189 knowledge into time-invariant, outdated, and new
190 categories — requiring knowledge retention, re-
191 moval, and acquisition, respectively — and bench-
192 mark existing continual pretraining methods for
193 knowledge updates.

194 Additionally, Hu et al. (2023) introduce a meta-
195 trained importance-weighting model to adjust per-
196 token loss dynamically, enabling LLMs to rapidly
197 adapt to new knowledge. Yu and Ji (2024) investi-
198 giate self-information updating in LLMs through
199 continual learning, addressing exposure bias by
200 incorporating fact selection into training losses.

201 4.2 Continual Pretraining for Domain 202 Adaptation

203 Continual domain adaptative pretraining (Ke et al.,
204 2022, 2023; Wu et al., 2024b) focuses on incre-
205 mentally adapting an LLM using a sequence of
206 unlabeled, domain-specific corpora. The objective
207 is to enable the LLM to accumulate knowledge
208 across multiple domains while mitigating cata-
209 strophic forgetting (McCloskey and Cohen, 1989)
210 of previously acquired domain knowledge or gen-
211 eral language understanding.

212 Gururangan et al. (2020) introduced the term of
213 domain-adaptive pretraining, demonstrating that
214 a second phase of pretraining on target domains
215 can effectively update an LLM with new domain
216 knowledge. It is important to note that further pre-
217 training can lead to catastrophic forgetting of gen-
218 eral concepts by overwriting essential parameters.
219 To mitigate this, recent works utilize *parameter-
220 isolation* methods which allocate different parame-
221 ter subsets to distinct tasks or domains and keep the
222 majority of parameters frozen (Razdaibiedina et al.,
223 2023; Wang et al., 2024d,e). DEMix-DAPT (Gu-
224 rurangan et al., 2022) replaces every feed-forward
225 network layer in the Transformer model with a do-
226 main expert mixture layer, containing one expert
227 per domain. When acquiring new knowledge, only
228 the newly added expert is trained while all others
229 remain fixed. Qin et al. (2022) propose ELLE for
230 efficient lifelong pretraining on various domains.
231 ELLE starts with a randomly initialized LLM and
232 expands the PLM’s width and depth to acquire new
233 knowledge more efficiently. Ke et al. (2022) intro-
234 duce a continual pretraining system which inserts
235 continual learning plugins to the frozen pretrained
236 language models that mitigate catastrophic forget-

ting while effectively learn new domain knowledge. Similarly, Lifelong-MoE (Chen et al., 2023) expands expert capacity progressively, freezing previously trained experts and applying output-level regularization to prevent forgetting.

In a later work, Ke et al. (2023) apply regularization to penalize changes to critical parameters learned from previous data, preventing catastrophic forgetting. Their approach computes the importance of LLM components, such as attention heads and neurons, in preserving general knowledge, applying soft-masking and contrastive loss during continual pretraining to maintain learned knowledge while promoting knowledge transfer.

4.3 Continual Pretraining for Language Expansion

Continual pretraining (CPT) has emerged as a pivotal strategy for adapting LLMs to new languages, or enhancing performance in underrepresented languages without full retraining (Wu et al., 2024b). Below, we discuss two major areas of expansion enabled by CPT: *natural language expansion* and *programming language expansion*.

Natural Language Expansion. Several recent studies have demonstrated the effectiveness of CPT in expanding language coverage. Glot500 (Imani et al., 2023) and EMMA-500 (Ji et al., 2024) enhance multilingual capabilities using CPT and vocabulary extension. Glot500, based on XLM-R (Ruder et al., 2019), and EMMA-500, built on LLaMA 2 (Touvron et al., 2023), expand language support up to 500 languages using extensive multilingual corpora. Similarly, Aya (Üstün et al., 2024) applies continual pretraining to the mT5 model (Xue et al., 2021) using a carefully constructed instruction dataset, achieving improved performance across 101 languages. Furthermore, LLaMAX (Lu et al., 2024) enhances multilingual translation by applying continual pretraining to the LLaMA model family. Supporting over 100 languages, it improves translation quality and promotes language inclusivity.

While covering many languages, many multilingual models exhibit suboptimal performance on medium- to low-resource languages (Ruder et al., 2019; Touvron et al., 2023; Imani et al., 2023). To bridge this performance gap, researchers have focused on expanding training corpora and strategically applying continual pretraining to enhance the multilingual capabilities of LLMs. Alabi et al.

(2022), Wang et al. (2023a), Fujii et al. (2024), and Zhang et al. (2024b) show that continual pre-training on one or more specific languages significantly improves performance across related languages. Blevins et al. (2024) extend this approach to the MoE paradigm for better parameter efficiency, while Zheng et al. (2024) investigate scaling laws for continual pretraining by training LLMs of varying sizes under different language distributions and conditions. Additionally, Tran (2020), Minixhofer et al. (2022), Dobler and de Melo (2023), Liu et al. (2024b), and Minixhofer et al. (2024) explore advanced tokenization and word embedding techniques to further improve LLMs’ multilingual performance in low-resource settings.

Programming Language Expansion. Going beyond natural languages, continual pretraining has demonstrated significant potential in enhancing the capabilities of LLMs for understanding and generating programming languages.

CERT, proposed by Zan et al. (2022), addresses the challenges of library-oriented code generation using unlabeled code corpora. It employs a two-stage framework to enable LLMs to effectively capture patterns in library-based code snippets. CodeTask-CL (Yadav et al., 2023) offers a benchmark for continual code learning, encompassing a diverse set of tasks such as code generation, summarization, translation, and refinement across multiple programming languages. Furthermore, continual pretrained models specifically for code understanding and programming from natural language prompts emerged with LLMs, such as CodeLLaMA (Grattafiori et al., 2023), Llama Pro (Wu et al., 2024a), CodeGemma (Team et al., 2024) and StarCoder 2 (Lozhkov et al., 2024), consistently outperform general-purpose LLMs of comparable or larger size on code benchmarks.

4.4 Continual Preference Alignment

Preference alignment ensures that large language models generate responses consistent with human values, improving usability, safety, and ethical behavior. While techniques like Reinforcement Learning from Human Feedback (RLHF) (Ziegler et al., 2019; Lambert et al., 2022) align LLMs with static preferences, societal values evolve, requiring continual preference alignment (CPA). It enables LLMs to adapt to emerging preferences while preserving previously learned values, ensuring relevance, inclusivity, and responsiveness to shifting

337 societal expectations. Despite its importance, CPA
338 remains a relatively underexplored area. Below,
339 we briefly discuss two representative works that
340 highlight the potential of this approach:

341 Zhang et al. (2023b) propose a non-
342 reinforcement learning approach for continual
343 preference alignment in LLMs. Their method uses
344 function regularization by computing an optimal
345 policy distribution for each task and applying it
346 to regularize future tasks, preventing catastrophic
347 forgetting while adapting to new domains. This
348 provides a single-phase, reinforcement learning-
349 free solution for maintaining alignment across
350 diverse tasks. Zhang et al. (2024a) introduce
351 Continual Proximal Policy Optimization (CPPO),
352 integrating continual learning into the RLHF
353 framework to accommodate evolving human
354 preferences. CPPO employs a sample-wise
355 weighting strategy to balance policy learning and
356 knowledge retention, consolidating high-reward
357 behaviors while mitigating overfitting and noise.

358 As the demand for responsive and inclusive AI
359 grows, CPA is key to keeping LLMs ethical and
360 aligned with evolving user needs, requiring further
361 research to reach its full potential.

362 4.5 Applicability and Limitations

363 Continual learning is a versatile framework for ex-
364 panding LLM knowledge across facts, domains,
365 languages, and preferences. It excels in large-scale
366 knowledge integration, retaining previously learned
367 knowledge, making it well-suited for tasks like do-
368 main adaptation and language expansion (Bu et al.,
369 2021; Jin et al., 2022; Cossu et al., 2024).

370 However, CL has notable limitations, including
371 a lack of precise control compared to model editing
372 (cf. Section 5) and retrieval-based methods (cf.
373 Section 6), inefficiency due to the computational
374 demands of retraining, and limited applicability
375 in black-box models. These challenges highlight
376 the need for alternative approaches like model edit-
377 ing and retrieval, which offer more targeted and
378 efficient updates.

379 5 Model Editing

380 Model editing offers a controllable and efficient
381 solution to update factual knowledge and user pref-
382 erences in LLMs. Introduced by Zhu et al. (2020),
383 De Cao et al. (2021) and Mitchell et al. (2022a),
384 it aims at modifying the model’s predictions for
385 specific inputs without affecting unrelated ones.

386 Yao et al. (2023) and Zhang et al. (2024c) define
387 four key evaluation metrics for model editing: (1)
388 **reliability**, ensuring the edited model produces the
389 target prediction for the target input; (2) **generaliza-**
390 **tion**, requiring the edited knowledge to apply to all
391 in-scope inputs — inputs that are directly related to
392 the target input, including rephrasings and seman-
393 tically similar variations; (3) **locality**, preserving
394 original outputs for unrelated out-of-scope inputs;
395 and (4) **portability**, extending the generalization
396 metric by assessing how well updated knowledge
397 transfers to complex rephrasings, reasoning chains,
398 and related facts.

399 While recent works (Mitchell et al., 2022b;
400 Madaan et al., 2022; Zhong et al., 2023; Zheng
401 et al., 2023) use *model editing* and *knowledge edit-
402 ing* interchangeably for updating factual knowl-
403 edge, we distinguish between them: model edit-
404 ing is a subset of knowledge editing that modifies
405 model parameters, whereas retrieval-based meth-
406 ods update knowledge dynamically without alter-
407 ing the model’s parameters (see Section 6).

5.1 Model Editing for Updating Facts

408 To address outdated or incorrect information
409 (Lazaridou et al., 2021), model editing research
410 focuses on selectively modifying this knowledge.
411 Below, we highlight key works in this area.

412 KnowledgeEditor (De Cao et al., 2021) uses a
413 hypernetwork to predict parameter shifts for mod-
414 ifying a fact, trained via constrained optimization
415 for locality. Similarly, MEND (Mitchell et al.,
416 2022a) trains a hypernetwork per LLM layer and
417 decomposes the fine-tuning gradient into a precise
418 one-step parameter update. Given the findings that
419 feed-forward layers in transformers function as
420 key-value memories (Geva et al., 2021), Dai et al.
421 (2022) introduce a knowledge attribution method
422 to identify these neurons and directly modify their
423 values via knowledge surgery.

425 Recent works employ a locate-and-edit strat-
426 egy for precise model editing. Using causal trac-
427 ing, Meng et al. (2022) identify middle-layer feed-
428 forward networks as key to factual predictions and
429 propose ROME, which updates facts by solving
430 a constrained least-squares problem in the MLP
431 weight matrix. MEMIT (Meng et al., 2023) extends
432 ROME to modify thousands of facts simultane-
433 ously across critical layers while preserving generaliza-
434 tion and locality. BIRD (Ma et al., 2023) intro-
435 duces bidirectional inverse relationship modeling
436 to mitigate the reverse curse (Berglund et al., 2003)

437 in model editing. While editing FFN layers has
438 proven effective, PMET (Li et al., 2024d) extends
439 editing to attention heads, achieving improved per-
440 formance. Wang et al. (2024g) further shift the
441 focus to conceptual knowledge, using ROME and
442 MEMIT to alter concept definitions, finding that
443 concept-level edits are reliable but have limited
444 influence on concrete examples.

445 5.2 Model Editing for Updating Preferences

446 Recent works expand model editing beyond factual
447 corrections to aligning LLMs with user preferences,
448 such as ensuring safety, reducing bias, and preserv-
449 ing privacy .

450 Wang et al. (2024c) use model editing to detox-
451 ify LLMs, ensuring safe responses to adversarial
452 inputs and preserving general LLM capabilities,
453 such as fluency, knowledge question answering,
454 and content summarization. Their results show that
455 model editing is promising for detoxification but
456 slightly affects general capabilities. Since LLMs
457 can exhibit social biases (Gallegos et al., 2024),
458 Chen et al. (2024a) propose fine-grained bias miti-
459 gation via model editing. Inspired by Meng et al.
460 (2022), they identify key layers responsible for bi-
461 ased knowledge and insert a feed-forward network
462 to adjust outputs with minimal parameter changes,
463 ensuring generalization, locality, and scalability.
464 For privacy protection, Wu et al. (2023) extend Dai
465 et al. (2022)'s work by identifying privacy neurons
466 that store sensitive information. Using gradient
467 attribution, they deactivate these neurons, reduc-
468 ing private data leakage while preserving model
469 performance. Moreover, Mao et al. (2024) apply
470 model editing techniques like MEND to modify
471 personality traits in LLMs, aligning responses to
472 opinion-based questions with target personalities.
473 While effective, this approach can degrade text gen-
474 eration quality.

475 5.3 Applicability and Limitations

476 Model editing complements continual learning
477 by allowing fine-grained knowledge updates with
478 lower computational costs. However, research has
479 primarily focused on structured, relational, and
480 instance-level knowledge, with limited exploration
481 of other knowledge types, multilingual generaliza-
482 tion, and cross-lingual transfer (Nie et al., 2024;
483 Wei et al., 2025).

484 Additionally, model editing faces several tech-
485 nical challenges, including limited locality and
486 gradual forgetting in large-scale edits (Bu et al.,

487 2019; Mitchell et al., 2022b; Gupta et al., 2024;
488 Li et al., 2024b), making it more suitable for mi-
489 nor updates. Additionally, it can impact general
490 LLM capabilities (Gu et al., 2024b; Wang et al.,
491 2024f) and downstream performance (Gupta et al.,
492 2024), potentially causing model collapse (Yang
493 et al., 2024b). Addressing these issues will enhance
494 model editing's role alongside continual learning
495 and retrieval, ensuring greater precision in dynamic
496 knowledge adaptation.

497 6 Retrieval-based Methods

498 Continual learning and model editing modify a
499 model's parameters to update its internal knowl-
500 edge, making them implicit knowledge expansion
501 methods (Zhang et al., 2023c). In contrast, retrieval-
502 based methods (Lewis et al., 2020) explicitly inte-
503 grate external knowledge, allowing models to over-
504 write outdated or undesired information without
505 parameter modifications. These methods leverage
506 external sources, such as databases, off-the-shelf
507 retriever systems, or the Internet, and thus provide
508 up-to-date or domain-specific knowledge (Zhang
509 et al., 2023c), making them effective for factual
510 updates and domain adaptation.

511 6.1 Retrieval-based Methods for Updating 512 Facts

513 Retrieval-based methods enhance LLMs by pair-
514 ing them with an updatable datastore, ensuring
515 access to current factual information. An early
516 approach, retrieval-augmented generation (RAG)
517 (Lewis et al., 2020), fine-tunes a pre-trained re-
518 triever end-to-end with the LLM to improve knowl-
519 edge retrieval. Similarly, kNN-LM (Khandelwal
520 et al., 2020) interpolates the LLM's output distribu-
521 tion with k-nearest neighbor search results from the
522 datastore, with later works optimizing efficiency
523 (He et al., 2021; Alon et al., 2022) and adapting it
524 for continual learning (Peng et al., 2023b).

525 For factual knowledge editing, Tandon et al.
526 (2022) store user feedback for post-hoc corrections,
527 while Mitchell et al. (2022b), Madaan et al. (2022),
528 and Dalvi Mishra et al. (2022) retrieve stored edits
529 to guide responses. Chen et al. (2024b) introduce
530 relevance filtering to efficiently handle multiple
531 edits. Retrieval-based in-context learning (Zheng
532 et al., 2023; Ram et al., 2023; Mallen et al., 2023;
533 Yu et al., 2023; Shi et al., 2024b; Bi et al., 2024)
534 enables dynamic factual updates.

535 For complex reasoning, retrieval supports multi-

hop question answering and iterative prompting: [Zhong et al. \(2023\)](#) propose iterative prompting for multi-hop knowledge editing, while [Gu et al. \(2024a\)](#) use a scope detector to retrieve relevant edits and improve question decomposition via entity extraction and knowledge prompts. Similarly, [Shi et al. \(2024c\)](#) enhance multi-hop question answering by retrieving fact chains from a knowledge graph with mutual information maximization and redundant fact pruning.

In multi-step decision-making, retrieval is combined with Chain-of-Thought (CoT) reasoning ([Trivedi et al., 2023](#); [Press et al., 2023](#)). Retrieval also aids post-generation fact-checking and refinement ([Gao et al., 2023](#); [Peng et al., 2023a](#); [Song et al., 2024](#)) by revising generated text or prompts based on retrieved facts.

For a more comprehensive review of retrieval-based factual knowledge updates, we refer to [Zhang et al. \(2023c\)](#).

6.2 Retrieval-based Methods for Domain Adaptation

Retrieval-based methods have been widely adopted for various domain-specific tasks, e.g., in science and finance. By integrating retrieved external knowledge, these models enhance their adaptability to specialized domains, improving decision-making, analysis, and information synthesis.

In the biomedical domain, retrieval-based approaches facilitate tasks, such as molecular property identification and drug discovery by integrating structured molecular data and information about biomedical entities like proteins, molecules, and diseases ([Wang et al., 2023b](#); [Liu et al., 2023](#); [Wang et al., 2024h](#); [Yang et al., 2024a](#)). For instance, [Wang et al. \(2023b\)](#) and [Li et al. \(2024a\)](#) introduce retrieval-based frameworks that extract relevant molecular data from databases to guide molecule generation. In protein research, retrieval-based approaches enhance protein representation and generation tasks ([Ma et al., 2024](#); [Sun et al., 2023](#)). Additionally, [Lozano et al. \(2023\)](#) develop a clinical question-answering system that retrieves relevant biomedical literature to provide more accurate responses in medical contexts.

The finance domain, characterized by its data-driven nature, also benefits from retrieval-based methods ([Li et al., 2024f,g](#)). [Zhang et al. \(2023a\)](#) enhance financial sentiment analysis by retrieving real-time financial data from external sources. Furthermore, financial question-answering also bene-

fits from retrieval-based methods, which involves extracting knowledge from professional financial documents. [Lin \(2024\)](#) introduces a PDF parsing method integrated with retrieval-augmented LLMs to retrieve relevant financial insights.

6.3 Applicability and Limitations

Despite their advantages, retrieval-based methods also come with several limitations. A major challenge is their reliance on external knowledge sources, which can introduce inconsistencies or outdated information if not properly curated ([Jin et al., 2024](#); [Xu et al., 2024](#)). Their effectiveness also depends on the quality and scope of the retrieval system ([Bai et al., 2024](#); [Liu et al., 2024a](#)); poor indexing or noisy retrieval may lead to irrelevant or misleading information. Another key issue is maintaining knowledge consistency across queries. Since retrieval-based methods do not update model parameters, contradictions can arise between retrieved facts and previously generated responses, affecting coherence ([Njeh et al., 2024](#); [Zhao et al., 2024](#); [Li et al., 2024c](#)).

Addressing these challenges is essential to improving retrieval-based approaches and ensuring their seamless integration with other LLM adaptation techniques.

7 Challenges, Opportunities, Guidelines

Solving Knowledge Conflicts. An inherent challenge of expanding a model’s knowledge is the emergence of knowledge conflicts, which can undermine the consistency and trustworthiness of LLMs ([Xu et al., 2024](#)). Studies have identified various types of conflicts following knowledge updates, including (i) temporal misalignment ([Luu et al., 2022](#)), where outdated and newly learned facts coexist inconsistently, (ii) model inconsistencies ([Huang et al., 2021](#)), where responses to similar queries vary unpredictably, and (iii) hallucinations ([Ji et al., 2023](#)), where the model generates fabricated or contradictory information. While some efforts have been made to address these issues ([Zhang and Choi, 2023](#); [Mallen et al., 2023](#); [Zhou et al., 2023](#); [Xie et al., 2024](#)), they remain an open challenge that requires further research and more robust solutions.

Minimizing Side Effects. Continual learning and model editing, both of which involve modifying model parameters, inevitably introduce side effects. A major challenge in continual learning

636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
is catastrophic forgetting (McCloskey and Cohen, 1989), where newly acquired knowledge overwrites previously learned information. In LLMs, the multi-stage nature of training exacerbates this issue, leading to cross-stage forgetting (Wu et al., 2024b), where knowledge acquired in earlier stages is lost as new training phases are introduced. For model editing, recent studies have shown that large-scale edits, particularly mass edits, can significantly degrade the model’s general capabilities, such as its language modeling performance (Wang et al., 2024f) or accuracy on general NLP tasks (Li et al., 2024e,b; Wang et al., 2024a). Effectively addressing these challenges is crucial for maximizing the potential of these methods for large-scale knowledge expansion while maintaining model stability and overall performance.

653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
Comprehensive Benchmarks. Although this paper explores the properties, strengths, and weaknesses of various methods for knowledge expansion, the discussion remains largely theoretical due to the lack of a comprehensive benchmark datasets for a uniform evaluation and a proper comparison. Existing works, such as Jang et al. (2022), Liska et al. (2022), and Kim et al. (2024), provide factual knowledge-based datasets and evaluate continual pretraining and/or retrieval-based methods. However, their experiments are limited in scale and fail to offer a comprehensive assessment. Developing benchmarks that encompass a variety of knowledge types and enable the evaluation of all methods would provide a more holistic and systematic understanding of their relative effectiveness.

669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
General Guideline. Selecting the appropriate method for knowledge expansion in LLM depends on the application context and type of knowledge that needs to be updated.

684
685
686
(i) For **factual knowledge**, model editing is ideal for precise, targeted updates, such as correcting specific facts, due to its efficiency and high level of control. Retrieval-based methods are effective for integrating dynamic or frequently changing facts, as they allow updates without modifying the model’s parameters, making them suitable for black-box applications. For large-scale factual updates, continual learning is preferred as it enables the incremental integration of new knowledge while preserving previously learned information.

(ii) For **domain knowledge**, both continual learning and retrieval-based methods are applicable. Continual learning excels in large-scale adap-

687
688
689
690
691
692
693
tation, using domain-specific corpora to ensure the model retains general knowledge while adapting to specialized contexts. Retrieval-based methods complement this by dynamically providing domain-specific information without requiring model modifications, making them valuable in scenarios where static updates are impractical.

694
695
696
697
698
699
(iii) For **language knowledge**, continual learning is the only method capable of supporting large-scale language expansion. It facilitates the integration of multilingual corpora and provides the foundational updates necessary for underrepresented or low-resource languages.

700
701
702
703
(iv) For **preference updates**, such as aligning models with evolving user values or ethical norms, continual alignment is typically achieved by combining continual learning techniques with preference optimization methods, such as reinforcement learning from human feedback. These approaches enable models to dynamically adapt to changing preferences while retaining alignment with previously learned values.

704
705
706
707
708
Summary. **Continual learning** is indispensable for large-scale updates like domain adaptation and language expansion, where foundational and incremental updates are required. **Model editing** excels at precise factual updates, while **retrieval-based methods** offer dynamic access to factual and domain knowledge without altering the model. A well-informed selection or combination of these methods ensures efficient and effective knowledge expansion tailored to specific use cases.

8 Conclusions

719
720
721
Adapting large language models to evolving knowledge is essential for maintaining their relevance and effectiveness. This survey explores three key adaptation methods — continual learning for large-scale updates, model editing for precise modifications, and retrieval-based approaches for external knowledge access without altering model parameters. We examine how these methods support updates across factual, domain-specific, language, and user preference knowledge while addressing challenges like scalability, controllability, and efficiency. By consolidating research and presenting a structured taxonomy, this survey provides insights into current strategies and future directions, promoting the development of more adaptable and efficient large language models.

736 Limitations

737 This survey provides a comprehensive overview of
738 knowledge expansion techniques for LLMs. How-
739 ever, due to page constraints, we had to limit its
740 scope and prioritize certain aspects:

741 First, the paper only provides a high-level
742 overview of each method rather than an in-depth
743 analysis. This can limit the understanding of the nu-
744ANCES and specific applications of each technique,
745 as well as implementation details.

746 Second, our work is a literature review of adap-
747 tation methods rather than an empirical study eval-
748 uating their actual performance. While we analyze
749 existing strategies, we do not benchmark or experi-
750 mentally compare their effectiveness, leaving room
751 for future studies to assess their practical impact
752 under real-world conditions.

753 Third, we focus solely on text-based models and
754 do not cover vision-language models, which inte-
755 grate multi-modal learning for textual and visual
756 understanding. While the methods covered in this
757 survey could be used to adapt the language en-
758 coders of such models in theory, extending these
759 adaptation methods to vision-language models re-
760 mains an open research direction.

761 Finally, this survey reflecting the current state of
762 research might become outdated as new research is
763 published, as the field of LLMs is rapidly evolving
764 and new methods for knowledge expansion are
765 continuously being developed.

766 References

767 Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius
768 Mosbach, and Dietrich Klakow. 2022. [Adapting pre-
769 trained language models to African languages via
770 multilingual adaptive fine-tuning](#). In *Proceedings of
771 the 29th International Conference on Computational
772 Linguistics*, pages 4336–4349, Gyeongju, Republic
773 of Korea. International Committee on Computational
774 Linguistics.

775 Uri Alon, Frank Xu, Junxian He, Sudipta Sengupta, Dan
776 Roth, and Graham Neubig. 2022. [Neuro-symbolic
777 language modeling with automaton-augmented re-
778 trieval](#). In *Proceedings of the 39th International
779 Conference on Machine Learning*, volume 162 of
780 *Proceedings of Machine Learning Research*, pages
781 468–485. PMLR.

782 Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jia-
783 heng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su,
784 Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024.
785 [MT-bench-101: A fine-grained benchmark for eval-
786 uating large language models in multi-turn dialogues](#).
787 In *Proceedings of the 62nd Annual Meeting of the*

788 *Association for Computational Linguistics (Volume 1:
789 Long Papers)*, pages 7421–7454, Bangkok, Thailand.
790 Association for Computational Linguistics.

791 Lukas Berglund, Meg Tong, Maximilian Kaufmann,
792 Mikita Balesni, Asa Cooper Stickland, Tomasz Kor-
793 bak, and Owain Evans. 2003. [The reversal curse:
794 Llms trained on “a is b” fail to learn “b is a”](#). In
795 *The Twelfth International Conference on Learning
796 Representations*.

797 Baolong Bi, Shenghua Liu, Lingrui Mei, Yiwei Wang,
798 Pengliang Ji, and Xueqi Cheng. 2024. [Decoding by
799 contrasting knowledge: Enhancing llms’ confidence
800 on edited facts](#). *arXiv preprint arXiv:2405.11613*.

801 Terra Blevins, Tomasz Limisiewicz, Suchin Gururan-
802 gan, Margaret Li, Hila Gonen, Noah A. Smith, and
803 Luke Zettlemoyer. 2024. [Breaking the curse of multi-
804 linguality with cross-lingual expert language models](#).
805 In *Proceedings of the 2024 Conference on Empiri-
806 cal Methods in Natural Language Processing*, pages
807 10822–10837, Miami, Florida, USA. Association for
808 Computational Linguistics.

809 Xingyuan Bu, Junran Peng, Junjie Yan, Tieniu Tan, and
810 Zhaoxiang Zhang. 2021. [GAIA: A Transfer Learning
811 System of Object Detection that Fits Your Needs](#). In
812 *2021 IEEE/CVF Conference on Computer Vision and
813 Pattern Recognition (CVPR)*, pages 274–283. IEEE
814 Computer Society.

815 Xingyuan Bu, Yuwei Wu, Zhi Gao, and Yunde Jia. 2019.
816 [Deep convolutional network with locality and spar-
817 sity constraints for texture classification](#). *Pattern
818 Recognition*, 91:34–46.

819 Boxi Cao, Hongyu Lin, Xianpei Han, and Le Sun. 2024.
820 [The life cycle of knowledge in big language models:
821 A survey](#). *Machine Intelligence Research*, 21(2):217–
822 238.

823 Ruizhe Chen, Yichen Li, Zikai Xiao, and Zuozhu Liu.
824 2024a. [Large language model bias mitigation from
825 the perspective of knowledge editing](#). *arXiv preprint
826 arXiv:2405.09341*.

827 Wuyang Chen, Yanqi Zhou, Nan Du, Yanping Huang,
828 James Laudon, Zhifeng Chen, and Claire Cui. 2023.
829 [Lifelong language pretraining with distribution-
830 specialized experts](#). In *Proceedings of the 40th Inter-
831 national Conference on Machine Learning*, volume
832 202 of *Proceedings of Machine Learning Research*,
833 pages 5383–5395. PMLR.

834 Yingfa Chen, Zhengyan Zhang, Xu Han, Chaojun Xiao,
835 Zhiyuan Liu, Chen Chen, Kuai Li, Tao Yang, and
836 Maosong Sun. 2024b. [Robust and scalable model
837 editing for large language models](#). In *Proceedings of
838 the 2024 Joint International Conference on Compu-
839 tational Linguistics, Language Resources and Evalua-
840 tion (LREC-COLING 2024)*, pages 14157–14172,
841 Torino, Italia. ELRA and ICCL.

842 Zhiyuan Chen and Bing Liu. 2018. [Continual learning
843 and catastrophic forgetting](#). In *Lifelong Machine
844 Learning*, pages 55–75. Springer.

845	Andrea Cossu, Antonio Carta, Lucia Passaro, Vincenzo Lomonaco, Tinne Tuytelaars, and Davide Bacciu.	Long Papers), pages 16477–16508, Toronto, Canada.	902
846	2024. Continual pre-training mitigates forgetting in language and vision. <i>Neural Networks</i> , 179:106492.	Association for Computational Linguistics.	903
847			
848			
849	Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei.	2022. Knowledge neurons in pretrained transformers. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8493–8502, Dublin, Ireland. Association for Computational Linguistics.	904
850			
851			
852			
853			
854			
855			
856	Bhavana Dalvi Mishra, Oyvind Tafjord, and Peter Clark.	2022. Towards teachable reasoning systems: Using a dynamic memory of user feedback for continual system improvement. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9465–9480, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	911
857			
858			
859			
860			
861			
862			
863			
864	Nicola De Cao, Wilker Aziz, and Ivan Titov.	2021. Editing factual knowledge in language models. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 6491–6506, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	912
865			
866			
867			
868			
869			
870	Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and William W. Cohen.	2022. Time-aware language models as temporal knowledge bases. <i>Transactions of the Association for Computational Linguistics</i> , 10:257–273.	913
871			
872			
873			
874			
875			
876	Konstantin Dobler and Gerard de Melo.	2023. FOCUS: Effective embedding initialization for monolingual specialization of multilingual models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13440–13454, Singapore. Association for Computational Linguistics.	914
877			
878			
879			
880			
881			
882			
883	Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota, Sakae Mizuki, Rio Yokota, and Naoaki Okazaki.	2024. Continual pre-training for cross-lingual LLM adaptation: Enhancing Japanese language capabilities. In <i>First Conference on Language Modeling</i> .	915
884			
885			
886			
887			
888			
889	Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed.	2024. Bias and fairness in large language models: A survey. <i>Computational Linguistics</i> , 50(3):1097–1179.	916
890			
891			
892			
893			
894			
895	Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu.	2023. RARR: Researching and revising what language models say, using language models. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16477–16508, Toronto, Canada.	917
896			
897			
898			
899			
900			
901			
902	Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy.	2021. Transformer feed-forward layers are key-value memories. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	918
903			
904			
905			
906			
907			
908			
909			
910			
911	Wenhan Xiong Grattafiori, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve.	2023. Code Llama: Open foundation models for code. <i>arXiv preprint arXiv:2308.12950</i> .	919
912			
913			
914			
915			
916	Hengrui Gu, Kaixiong Zhou, Xiaotian Han, Ning-hao Liu, Ruobing Wang, and Xin Wang.	2024a. PokeMQA: Programmable knowledge editing for multi-hop question answering. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8069–8083, Bangkok, Thailand. Association for Computational Linguistics.	920
917			
918			
919			
920			
921	Jia-Chen Gu, Hao-Xiang Xu, Jun-Yu Ma, Pan Lu, Zhen-Hua Ling, Kai-Wei Chang, and Nanyun Peng.	2024b. Model editing harms general abilities of large language models: Regularization to the rescue. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 16801–16819, Miami, Florida, USA. Association for Computational Linguistics.	922
922			
923			
924	Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli.	2024. Model editing at scale leads to gradual and catastrophic forgetting. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 15202–15232, Bangkok, Thailand. Association for Computational Linguistics.	925
925			
926			
927			
928			
929			
930			
931			
932	Suchin Gururangan, Mike Lewis, Ari Holtzman, Noah A. Smith, and Luke Zettlemoyer.	2022. DEMIx layers: Disentangling domains for modular language modeling. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5557–5576, Seattle, United States. Association for Computational Linguistics.	933
933			
934			
935			
936			
937			
938	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith.	2020. Don't stop pretraining: Adapt language models to domains and tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360, Online. Association for Computational Linguistics.	939
939			
940			
941			
942			
943			
944			
945			
946	Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick.	2021. Efficient nearest neighbor language models. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5703–5714, Online and Punta Cana,	947
947			
948			
949			
950			
951			
952			
953			

959	Dominican Republic. Association for Computational Linguistics.	1015
960		1016
961	Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. 2024. Linearity of relation decoding in transformer language models . In <i>The Twelfth International Conference on Learning Representations</i> .	1017
962		1018
963		1019
964		1020
965		1021
966		1022
967	Nathan Hu, Eric Mitchell, Christopher Manning, and Chelsea Finn. 2023. Meta-learning online adaptation of language models . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4418–4432, Singapore. Association for Computational Linguistics.	1023
968		
969		
970		
971		
972		
973	Yichong Huang, Xiachong Feng, Xiaocheng Feng, and Bing Qin. 2021. The factual inconsistency problem in abstractive text summarization: A survey . <i>arXiv preprint arXiv:2104.14839</i> .	1024
974		1025
975		1026
976		1027
977	Zeyu Huang, Yikang Shen, Xiaofeng Zhang, Jie Zhou, Wenge Rong, and Zhang Xiong. 2023. Transformer-patcher: One mistake worth one neuron . In <i>The Eleventh International Conference on Learning Representations</i> .	1028
978		1029
979		1030
980		
981		
982	Ayyoob Imani, Peiqin Lin, Amir Hossein Kargaran, Silvia Severini, Masoud Jalili Sabet, Nora Kassner, Chunlan Ma, Helmut Schmid, André Martins, François Yvon, and Hinrich Schütze. 2023. Glot500: Scaling multilingual corpora and language models to 500 languages . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1082–1117, Toronto, Canada. Association for Computational Linguistics.	1031
983		1032
984		1033
985		
986		
987		
988		
989		
990		
991		
992	Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2022. Towards continual knowledge learning of language models . In <i>International Conference on Learning Representations</i> .	1034
993		1035
994		1036
995		1037
996		
997	Shaoxiong Ji, Zihao Li, Indraneil Paul, Jaakko Paavola, Peiqin Lin, Pinzhen Chen, Dayyán O’Brien, Hengyu Luo, Hinrich Schütze, Jörg Tiedemann, et al. 2024. Emma-500: Enhancing massively multilingual adaptation of large language models . <i>arXiv preprint arXiv:2409.17892</i> .	1038
998		1039
999		1040
1000		1041
1001		1042
1002		
1003	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation . <i>ACM Comput. Surv.</i> , 55(12).	1043
1004		1044
1005		1045
1006		1046
1007		1047
1008	Jiajie Jin, Yutao Zhu, Yujia Zhou, and Zhicheng Dou. 2024. BIDER: Bridging knowledge inconsistency for efficient retrieval-augmented LLMs via key supporting evidence . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 750–761, Bangkok, Thailand. Association for Computational Linguistics.	1048
1009		1049
1010		1050
1011		1051
1012		
1013		
1014		
1015	Xisen Jin, Dejiao Zhang, Henghui Zhu, Wei Xiao, Shang-Wen Li, Xiaokai Wei, Andrew Arnold, and Xiang Ren. 2022. Lifelong pretraining: Continually adapting language models to emerging corpora . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4764–4780, Seattle, United States. Association for Computational Linguistics.	1052
1016		1053
1017		1054
1018		1055
1019		
1020		
1021		
1022		
1023		
1024	Zixuan Ke, Haowei Lin, Yijia Shao, Hu Xu, Lei Shu, and Bing Liu. 2022. Continual training of language models for few-shot learning . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10205–10216, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1056
1025		1057
1026		1058
1027		1059
1028		1060
1029		1061
1030		1062
1031	Zixuan Ke and Bing Liu. 2022. Continual learning of natural language processing tasks: A survey . <i>arXiv preprint arXiv:2211.12701</i> .	1063
1032		1064
1033		
1034	Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. 2023. Continual pre-training of language models . In <i>The Eleventh International Conference on Learning Representations</i> .	1065
1035		1066
1036		1067
1037		1068
1038	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. Generalization through memorization: Nearest neighbor language models . In <i>International Conference on Learning Representations</i> .	1069
1039		1070
1040		1071
1041		1072
1042		
1043	Yujin Kim, Jaehong Yoon, Seonghyeon Ye, Sangmin Bae, Namgyu Ho, Sung Ju Hwang, and Se-Young Yun. 2024. Carpe diem: On the evaluation of world knowledge in lifelong language models . In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5401–5415, Mexico City, Mexico. Association for Computational Linguistics.	1073
1044		1074
1045		1075
1046		1076
1047		1077
1048		1078
1049		1079
1050		1080
1051		1081
1052	Nathan Lambert, Louis Castricato, Leandro von Werra, and Alex Havrilla. 2022. Illustrating reinforcement learning from human feedback (rlhf) . <i>Hugging Face Blog</i> . Https://huggingface.co/blog/rlhf .	1082
1053		1083
1054		1084
1055		1085
1056	Angeliki Lazaridou, Adhiguna Kuncoro, Elena Grivovskaya, Devang Agrawal, Adam Liška, Tayfun Terzi, Mai Gimenez, Cyprien de Masson d’Autume, Tomas Kociský, Sebastian Ruder, Dani Yogatama, Kris Cao, Susannah Young, and Phil Blunsom. 2021. Mind the gap: Assessing temporal generalization in neural language models . In <i>Proceedings of the 35th International Conference on Neural Information Processing Systems, NIPS ’21</i> . Curran Associates Inc.	1086
1057		1087
1058		1088
1059		1089
1060		1090
1061		1091
1062		1092
1063		1093
1064		1094
1065	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks . In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20</i> . Curran Associates Inc.	1095
1066		1096
1067		1097
1068		1098
1069		1099
1070		1100
1071		1101
1072		1102

1073	Jiatong Li, Yunqing Liu, Wenqi Fan, Xiao-Yong Wei, Hui Liu, Jiliang Tang, and Qing Li. 2024a. Empowering Molecule Discovery for Molecule-Caption Translation With Large Language Models: A Chat- GPT Perspective . <i>IEEE Transactions on Knowledge & Data Engineering</i> , 36(11):6071–6083.	1128
1074		1129
1075	Dong Liu, Roger Waleffe, Meng Jiang, and Shivaram Venkataraman. 2024a. Graphsnapshot: Graph ma- chine learning acceleration with fast storage and re- trieval. <i>arXiv preprint arXiv:2406.17918</i> .	1130
1076		1131
1077		
1078		
1079	Qi Li, Xiang Liu, Zhenheng Tang, Peijie Dong, Zeyu Li, Xinglin Pan, and Xiaowen Chu. 2024b. Should we really edit language models? on the evaluation of edited language models. In <i>The Thirty-eighth An- nual Conference on Neural Information Processing Systems</i> .	1132
1080		1133
1081	Shengchao Liu, Weili Nie, Chengpeng Wang, Jiarui Lu, Zhuoran Qiao, Ling Liu, Jian Tang, Chaowei Xiao, and Animashree Anandkumar. 2023. Multi- modal molecule structure–text model for text-based retrieval and editing. <i>Nature Machine Intelligence</i> , 5(12):1447–1457.	1134
1082		1135
1083		1136
1084		1137
1085	Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schuetze. 2024b. OFA: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining. In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 1067–1097, Mexico City, Mexico. Asso- ciation for Computational Linguistics.	1138
1086		1139
1087		1140
1088		1141
1089		1142
1090		1143
1091		1144
1092		
1093		
1094	Shilong Li, Yancheng He, Hangyu Guo, Xingyuan Bu, Ge Bai, Jie Liu, Jiaheng Liu, Xingwei Qu, Yang- guang Li, Wanli Ouyang, Wenbo Su, and Bo Zheng. 2024c. GraphReader: Building graph-based agent to enhance long-context abilities of large language mod- els. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 12758–12786, Mi- ami, Florida, USA. Association for Computational Linguistics.	1145
1095		1146
1096		1147
1097		1148
1098		1149
1099	Alejandro Lozano, Scott L Fleming, Chia-Chun Chiang, and Nigam Shah. 2023. Clinfo. ai: An open-source retrieval-augmented large language model system for answering medical questions using scientific litera- ture. In <i>PACIFIC SYMPOSIUM ON BIocomput- ING 2024</i> , pages 8–23. World Scientific.	1150
1100		
1101		
1102		
1103		
1104	Zhaobo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. 2024e. Unveiling the pit- falls of knowledge editing for large language models. In <i>The Twelfth International Conference on Learning Representations</i> .	1151
1105		1152
1106		1153
1107		1154
1108		1155
1109	Yinquan Lu, Wenhao Zhu, Lei Li, Yu Qiao, and Fei Yuan. 2024. LLaMAX: Scaling linguistic horizons of LLM by enhancing translation capabilities beyond 100 languages. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10748–10772, Miami, Florida, USA. Association for Computational Linguistics.	1156
1110		1157
1111		1158
1112		1159
1113		1160
1114	Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. <i>arXiv preprint arXiv:2402.19173</i> .	1151
1115		1152
1116		1153
1117		1154
1118		1155
1119		
1120		
1121		
1122		
1123		
1124		
1125		
1126		
1127		
1128	Kelvin Luu, Daniel Khashabi, Suchin Gururangan, Kar- ishma Mandyam, and Noah A. Smith. 2022. Time waits for no one! analysis and challenges of tem- poral misalignment. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5944–5958, Seattle, United States. Association for Computational Lin- guistics.	1163
1129		1164
1130		1165
1131		1166
1132		1167
1133		1168
1134		1169
1135		1170
1136		1171
1137		
1138	Demiao Lin. 2024. Revolutionizing retrieval- augmented generation with enhanced pdf structure recognition. <i>arXiv preprint arXiv:2401.12599</i> .	1172
1139		1173
1140	Chang Ma, Haiteng Zhao, Lin Zheng, Jiayi Xin, Qin- tong Li, Lijun Wu, Zhihong Deng, Yang Young Lu, Qi Liu, Sheng Wang, and Lingpeng Kong. 2024. Re- trieved sequence augmentation for protein represen- tation learning. In <i>Proceedings of the 2024 Con- ference on Empirical Methods in Natural Language Processing</i> , pages 1738–1767, Miami, Florida, USA. Association for Computational Linguistics.	1174
1141		1175
1142		1176
1143		1177
1144		1178
1145		1179
1146		
1147		
1148		
1149		
1150		
1151		
1152		
1153		
1154		
1155		
1156		
1157		
1158		
1159		
1160		
1161		
1162		
1163		
1164		
1165		
1166		
1167		
1168		
1169		
1170		
1171		
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180		
1181		
1182		
1183		

1184	Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. 2022. Memory-assisted prompt editing to improve GPT-3 after deployment. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2833–2861, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	1241
1185		1242
1186		1243
1187		
1188	Ercong Nie, Bo Shao, Zifeng Ding, Mingyang Wang, Helmut Schmid, and Hinrich Schütze. 2024. Bmike-53: Investigating cross-lingual knowledge editing with in-context learning. <i>arXiv preprint arXiv:2406.17764</i> .	1244
1189		1245
1190		1246
1191	Chaima Njeh, Haïfa Nakouri, and Fehmi Jaafar. 2024. Enhancing RAG-retrieval to improve LLMs robustness and resilience to hallucinations. In <i>Hybrid Artificial Intelligent Systems</i> , pages 201–213. Springer Nature Switzerland.	1247
1192		1248
1193		
1194	Vaidehi Patil, Peter Hase, and Mohit Bansal. 2024. Can sensitive information be deleted from LLMs? objectives for defending against extraction attacks. In <i>The Twelfth International Conference on Learning Representations</i> .	1249
1195		1250
1196		1251
1197		1252
1198		1253
1199	Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, and Jianfeng Gao. 2023a. Check your facts and try again: Improving large language models with external knowledge and automated feedback. <i>arXiv preprint arXiv:2302.12813</i> .	1254
1200		1255
1201		1256
1202		1257
1203		1258
1204	Guangyue Peng, Tao Ge, Si-Qing Chen, Furu Wei, and Houfeng Wang. 2023b. Semiparametric language models are scalable continual learners. <i>arXiv preprint arXiv:2303.01421</i> .	1259
1205		1260
1206	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah Smith, and Mike Lewis. 2023. Measuring and narrowing the compositionality gap in language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5687–5711, Singapore. Association for Computational Linguistics.	1261
1207		1262
1208		1263
1209		1264
1210		
1211	Yujia Qin, Jiajie Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2022. ELLE: Efficient lifelong pre-training for emerging data. In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 2789–2810, Dublin, Ireland. Association for Computational Linguistics.	1265
1212		1266
1213		1267
1214		1268
1215		
1216	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. <i>Transactions of the Association for Computational Linguistics</i> , 11:1316–1331.	1269
1217		1270
1218		1271
1219		1272
1220		1273
1221	Anastasia Razdaibiedina, Yuning Mao, Rui Hou, Madian Khabsa, Mike Lewis, and Amjad Almahairi. 2023. Progressive prompts: Continual learning for language models. In <i>The Eleventh International Conference on Learning Representations</i> .	1274
1222		
1223		
1224		
1225	Paul Röttger and Janet Pierrehumbert. 2021. Temporal adaptation of BERT and performance on downstream document classification: Insights from social media. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 2400–2412, Punta Cana, Dominican Republic. Association for Computational Linguistics.	1275
1226		1291
1227		1292
1228		1293
1229		1294
1230		1295
1231		1296
1232		1297
1233		
1234	Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. 2022a. Fast model editing at scale. In <i>International Conference on Learning Representations</i> .	1286
1235		1287
1236		1288
1237		1289
1238	Eric Mitchell, Charles Lin, Antoine Bosselut, Christopher D Manning, and Chelsea Finn. 2022b. Memory-based model editing at scale. In <i>Proceedings of the</i>	1290
1239		
1240		

1298	Sebastian Ruder, Anders Søgaard, and Ivan Vulić. 2019.	CodeGemma Team, Heri Zhao, Jeffrey Hui, Joshua Howland, Nam Nguyen, Siqi Zuo, Andrea Hu, Christopher A Choquette-Choo, Jingyue Shen, Joe Kelley, et al. 2024. Codegemma: Open code models based on gemma . <i>arXiv preprint arXiv:2406.11409</i> .	1356
1299	Unsupervised cross-lingual representation learning .		1357
1300	In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts</i> , pages 31–38, Florence, Italy. Association for Computational Linguistics.		1358
1301			1359
1302			1360
1303			
1304	Haizhou Shi, Zihao Xu, Hengyi Wang, Weiyi Qin,	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models . <i>arXiv preprint arXiv:2307.09288</i> .	1361
1305	Wenyan Wang, Yibin Wang, Zifeng Wang, Sayna Ebrahimi, and Hao Wang. 2024a. Continual learning of large language models: A comprehensive survey .		1362
1306			1363
1307			1364
1308	Weijia Shi, Julian Michael, Suchin Gururangan, and Luke Zettlemoyer. 2022. Nearest neighbor zero-shot inference .	Ke Tran. 2020. From english to foreign languages: Transferring pre-trained language models . <i>arXiv preprint arXiv:2002.07306</i> .	1365
1309	In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3254–3265, Abu Dhabi, United Arab Emirates.		1366
1310	Association for Computational Linguistics.		
1311			
1312			
1313			
1314	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024b. REPLUG: Retrieval-augmented black-box language models .	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions .	1370
1315	In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8371–8384, Mexico City, Mexico. Association for Computational Linguistics.	In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10014–10037, Toronto, Canada. Association for Computational Linguistics.	1371
1316			1372
1317			1373
1318			1374
1319			1375
1320			1376
1321			1377
1322			
1323	Yucheng Shi, Qiaoyu Tan, Xuansheng Wu, Shaochen Zhong, Kaixiong Zhou, and Ninghao Liu. 2024c. Retrieval-enhanced knowledge editing in language models for multi-hop question answering .	Ahmet Üstün, Viraat Aryabumi, Zheng Yong, Wei-Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid, Freddie Vargas, Phil Blunsom, Shayne Longpre, Niklas Muennighoff, Marzieh Fadaee, Julia Kreutzer, and Sara Hooker. 2024. Aya model: An instruction fine-tuned open-access multilingual language model .	1378
1324	In <i>Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM ’24</i> , page 2056–2066. Association for Computing Machinery.	In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15894–15939, Bangkok, Thailand. Association for Computational Linguistics.	1379
1325			1380
1326			1381
1327			1382
1328			1383
1329			1384
1330			1385
1331	Chenglei Si, Zhe Gan, Zhengyuan Yang, Shuohang Wang, Jianfeng Wang, Jordan Lee Boyd-Graber, and Lijuan Wang. 2023. Prompting GPT-3 to be reliable .	Jianchen Wang, Zhouhong Gu, Zhuozhi Xiong, Hongwei Feng, and Yanghua Xiao. 2024a. The missing piece in model editing: A deep dive into the hidden damage brought by model editing . <i>arXiv preprint arXiv:2403.07825</i> .	1386
1332	In <i>The Eleventh International Conference on Learning Representations</i> .		1387
1333			1388
1334			
1335			
1336	Xiaoshuai Song, Zhengyang Wang, Keqing He, Guanting Dong, Yutao Mou, Jinxu Zhao, and Weiran Xu. 2024. Knowledge editing on black-box large language models . <i>arXiv preprint arXiv:2402.08631</i> .	Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. 2024b. A comprehensive survey of continual learning: theory, method and application . <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> .	1389
1337			1390
1338			1391
1339			1392
1340	Fang Sun, Zhihao Zhan, Hongyu Guo, Ming Zhang, and Jian Tang. 2023. Graphvf: Controllable protein-specific 3d molecule generation with variational flow .	Mengru Wang, Ningyu Zhang, Ziwen Xu, Zekun Xi, Shumin Deng, Yunzhi Yao, Qishen Zhang, Linyi Yang, Jindong Wang, and Huajun Chen. 2024c. Detoxifying large language models via knowledge editing .	1393
1341	<i>arXiv preprint arXiv:2304.12825</i> .	In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3093–3118, Bangkok, Thailand. Association for Computational Linguistics.	
1342			
1343			
1344	Yu Sun, Shuhuan Wang, Yukun Li, Shikun Feng, Hao Tian, Hua Wu, and Haifeng Wang. 2020. Ernie 2.0: A continual pre-training framework for language understanding .	Mingyang Wang, Heike Adel, Lukas Lange, Jannik Strötgen, and Hinrich Schütze. 2024d. Learn it or leave it: Module composition and pruning for continual learning .	1398
1345	<i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(05):8968–8975.	In <i>Proceedings of the 9th Workshop on Representation Learning for NLP (RepL4NLP-2024)</i> , pages 163–176, Bangkok, Thailand. Association for Computational Linguistics.	1407
1346			1408
1347			1409
1348			1410
1349	Niket Tandon, Aman Madaan, Peter Clark, and Yiming Yang. 2022. Learning to repair: Repairing model output errors after deployment using a dynamic memory of feedback .		1411
1350	In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 339–352, Seattle, United States. Association for Computational Linguistics.		1412
1351			
1352			
1353			
1354			
1355			

1413	Mingyang Wang, Heike Adel, Lukas Lange, Jannik Strötgen, and Hinrich Schuetze. 2024e. Rehearsal-free modular and compositional continual learning for language models. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)</i> , pages 469–480, Mexico City, Mexico. Association for Computational Linguistics.	Continual learning for large language models: A survey. <i>arXiv preprint arXiv:2402.01364.</i>	1471 1472
1414		Xinwei Wu, Junzhuo Li, Minghui Xu, Weilong Dong, Shuangzhi Wu, Chao Bian, and Deyi Xiong. 2023. DEPN: Detecting and editing privacy neurons in pre-trained language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2875–2886, Singapore. Association for Computational Linguistics.	1473 1474 1475 1476 1477 1478 1479
1415		Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2024. Adaptive chameleon or stubborn sloth: Revealing the behavior of large language models in knowledge conflicts. In <i>The Twelfth International Conference on Learning Representations</i> .	1480 1481 1482 1483 1484
1416		Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. Knowledge conflicts for LLMs: A survey. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 8541–8565, Miami, Florida, USA. Association for Computational Linguistics.	1485 1486 1487 1488 1489 1490 1491
1417		Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.	1492 1493 1494 1495 1496 1497 1498 1499
1418		Prateek Yadav, Qing Sun, Hantian Ding, Xiaopeng Li, Dejiao Zhang, Ming Tan, Parminder Bhatia, Xiaofei Ma, Ramesh Nallapati, Murali Krishna Ramanathan, Mohit Bansal, and Bing Xiang. 2023. Exploring continual learning for code generation models. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 782–792, Toronto, Canada. Association for Computational Linguistics.	1500 1501 1502 1503 1504 1505 1506 1507 1508
1419		Ling Yang, Zhilin Huang, Xiangxin Zhou, Minkai Xu, Wentao Zhang, Yu Wang, Xiawu Zheng, Wenming Yang, Ron O. Dror, Shenda Hong, and Bin CUI. 2024a. Prompt-based 3d molecular diffusion models for structure-based drug design.	1509 1510 1511 1512 1513
1420		Wanli Yang, Fei Sun, Jiajun Tan, Xinyu Ma, Du Su, Dawei Yin, and Huawei Shen. 2024b. The fall of ROME: Understanding the collapse of LLMs in model editing. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 4079–4087, Miami, Florida, USA. Association for Computational Linguistics.	1514 1515 1516 1517 1518 1519 1520
1421		Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubao Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. Editing large language models: Problems, methods, and opportunities. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 10222–10240, Singapore. Association for Computational Linguistics.	1521 1522 1523 1524 1525 1526 1527 1528
1422	Mingyang Wang, Heike Adel, Lukas Lange, Jannik Strötgen, and Hinrich Schütze. 2023a. NLNDE at SemEval-2023 task 12: Adaptive pretraining and source language selection for low-resource multi-lingual sentiment analysis. In <i>Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)</i> , pages 488–497, Toronto, Canada. Association for Computational Linguistics.		
1423			
1424			
1425			
1426			
1427			
1428			
1429			
1430	Mingyang Wang, Lukas Lange, Heike Adel, Jannik Strötgen, and Hinrich Schütze. 2024f. Better call SAUL: Fluent and consistent language model editing with generation regularization. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 7990–8000, Miami, Florida, USA. Association for Computational Linguistics.		
1431			
1432			
1433			
1434			
1435			
1436			
1437	Xiaohan Wang, Shengyu Mao, Shumin Deng, Yunzhi Yao, Yue Shen, Lei Liang, Jinjie Gu, Huajun Chen, and Ningyu Zhang. 2024g. Editing conceptual knowledge for large language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 706–724, Miami, Florida, USA. Association for Computational Linguistics.		
1438			
1439			
1440			
1441			
1442			
1443			
1444	Zichao Wang, Weili Nie, Zhuoran Qiao, Chaowei Xiao, Richard Baraniuk, and Anima Anandkumar. 2023b. Retrieval-based controllable molecule generation. In <i>The Eleventh International Conference on Learning Representations</i> .		
1445			
1446			
1447			
1448			
1449	Zifeng Wang, Zichen Wang, Balasubramaniam Srinivasan, Vassilis N. Ioannidis, Huzeifa Rangwala, and RISHITA ANUBHAI. 2024h. Biobridge: Bridging biomedical foundation models via knowledge graphs. In <i>The Twelfth International Conference on Learning Representations</i> .		
1450			
1451			
1452			
1453			
1454			
1455	Zihao Wei, Jingcheng Deng, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. 2025. MLaKE: Multilingual knowledge editing benchmark for large language models. In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 4457–4473, Abu Dhabi, UAE. Association for Computational Linguistics.		
1456			
1457			
1458			
1459			
1460			
1461			
1462	Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ying Shan, and Ping Luo. 2024a. LLaMA pro: Progressive LLaMA with block expansion. In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 6518–6537, Bangkok, Thailand. Association for Computational Linguistics.		
1463			
1464			
1465			
1466			
1467			
1468			
1469	Tongtong Wu, Linhao Luo, Yuan-Fang Li, Shirui Pan, Thuy-Trang Vu, and Gholamreza Haffari. 2024b.		
1470			

1529 Antonio Jimeno Yepes, Yao You, Jan Milczek, Sebastian Laverde, and Renyu Li. 2024. *Financial report chunking for effective retrieval augmented generation*. *arXiv preprint arXiv:2402.05131*.

1533 Pengfei Yu and Heng Ji. 2024. *Information association for language model updating by mitigating LM-logical discrepancy*. In *Proceedings of the 28th Conference on Computational Natural Language Learning*, pages 117–129, Miami, FL, USA. Association for Computational Linguistics.

1539 Zichun Yu, Chenyan Xiong, Shi Yu, and Zhiyuan Liu. 2023. *Augmentation-adapted retriever improves generalization of language models as generic plug-in*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2421–2436, Toronto, Canada. Association for Computational Linguistics.

1546 Daoguang Zan, Bei Chen, Dejian Yang, Zeqi Lin, Minsu Kim, Bei Guan, Yongji Wang, Weizhu Chen, and Jian-Guang Lou. 2022. *CERT: Continual pre-training on sketches for library-oriented code generation*. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 2369–2375. International Joint Conferences on Artificial Intelligence Organization.

1554 Boyu Zhang, Hongyang Yang, Tianyu Zhou, Muhammad Ali Babar, and Xiao-Yang Liu. 2023a. *Enhancing financial sentiment analysis via retrieval augmented large language models*. In *Proceedings of the Fourth ACM International Conference on AI in Finance, ICAIF ’23*, pages 349–356. Association for Computing Machinery.

1561 Han Zhang, Lin Gui, Yuanzhao Zhai, Hui Wang, Yu Lei, and Ruifeng Xu. 2023b. *Copf: Continual learning human preference through optimal policy fitting*. *arXiv preprint arXiv:2310.15694*.

1565 Han Zhang, Yu Lei, Lin Gui, Min Yang, Yulan He, Hui Wang, and Ruifeng Xu. 2024a. *Cppo: Continual learning for reinforcement learning with human feedback*. In *The Twelfth International Conference on Learning Representations*.

1570 Miaoran Zhang, Mingyang Wang, Jesujoba Alabi, and Dietrich Klakow. 2024b. *AAAdAM at SemEval-2024 task 1: Augmentation and adaptation for multilingual semantic textual relatedness*. In *Proceedings of the 18th International Workshop on Semantic Evaluation (SemEval-2024)*, pages 800–810, Mexico City, Mexico. Association for Computational Linguistics.

1577 Michael Zhang and Eunsol Choi. 2023. *Mitigating temporal misalignment by discarding outdated facts*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14213–14226, Singapore. Association for Computational Linguistics.

1583 Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, et al. 2024c. *A comprehensive study of knowledge editing for large language models*. *arXiv preprint arXiv:2401.01286*.

1589 Zihan Zhang, Meng Fang, Ling Chen, Mohammad-Reza Namazi-Rad, and Jun Wang. 2023c. *How do large language models capture the ever-changing world knowledge? a review of recent advances*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8289–8311, Singapore. Association for Computational Linguistics.

1595 Wayne Xin Zhao, Jing Liu, Ruiyang Ren, and Ji-Rong Wen. 2024. *Dense text retrieval based on pretrained language models: A survey*. *ACM Trans. Inf. Syst.*, 42(4).

1600 Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. *Can we edit factual knowledge by in-context learning?* In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4862–4876, Singapore. Association for Computational Linguistics.

1607 Wenzhen Zheng, Wenbo Pan, Xu Xu, Libo Qin, Li Yue, and Ming Zhou. 2024. *Breaking language barriers: Cross-lingual continual pre-training at scale*. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 7725–7738, Miami, Florida, USA. Association for Computational Linguistics.

1614 Zexuan Zhong, Zhengxuan Wu, Christopher Manning, Christopher Potts, and Danqi Chen. 2023. *MQuAKE: Assessing knowledge editing in language models via multi-hop questions*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15686–15702, Singapore. Association for Computational Linguistics.

1621 Wenxuan Zhou, Sheng Zhang, Hoifung Poon, and Muhaoo Chen. 2023. *Context-faithful prompting for large language models*. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 14544–14556, Singapore. Association for Computational Linguistics.

1627 Chen Zhu, Ankit Singh Rawat, Manzil Zaheer, Srinadh Bhojanapalli, Daliang Li, Felix Yu, and Sanjiv Kumar. 2020. *Modifying memories in transformer models*. *arXiv preprint arXiv:2012.00363*.

1631 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. 2019. *Fine-tuning language models from human preferences*. *arXiv preprint arXiv:1909.08593*.

A Appendix

A.1 Comprehensive Taxonomy of Methods

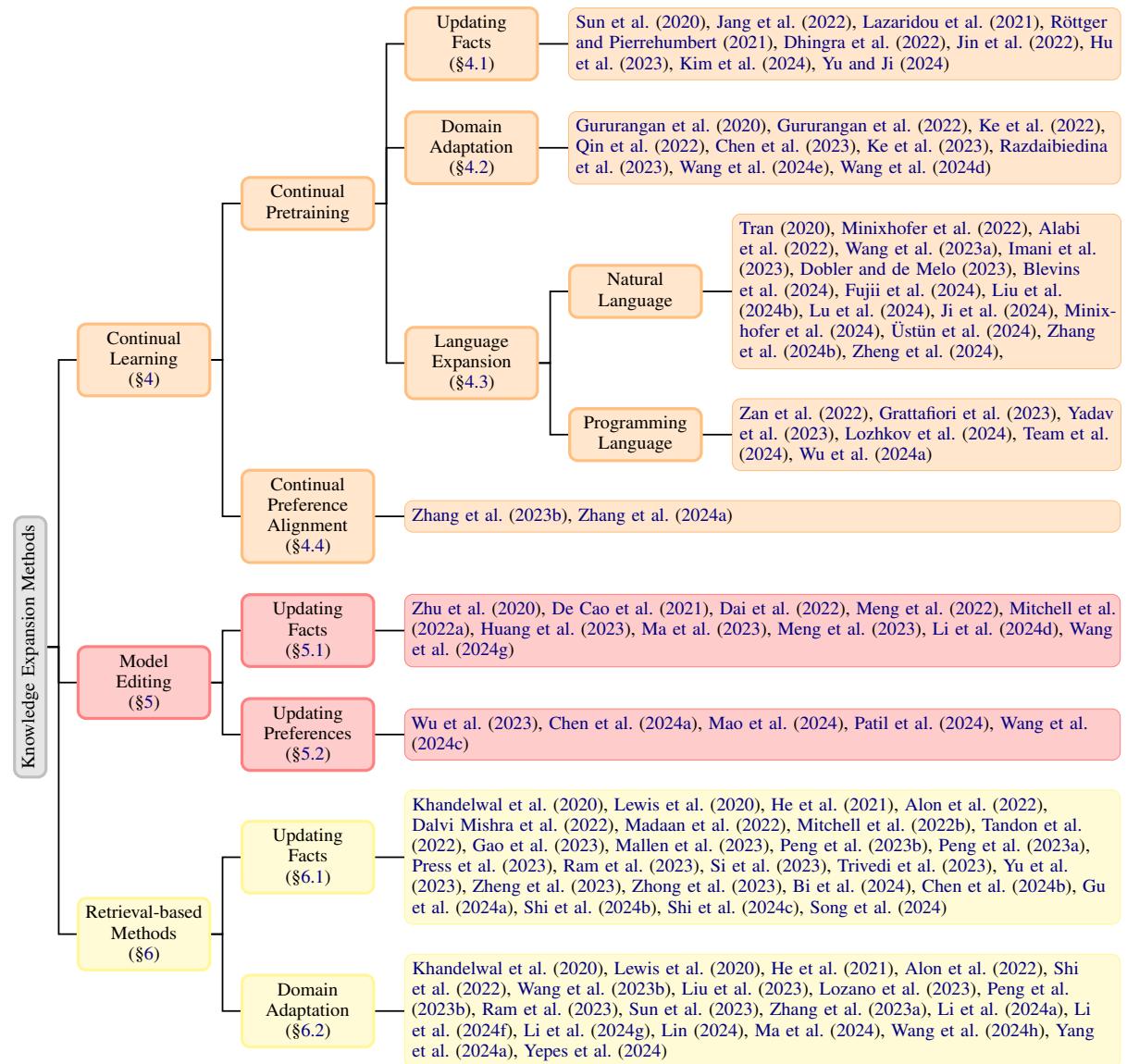


Figure 2: Taxonomy of methods for expanding LLM knowledge.