

Contextual Compression in Retrieval-Augmented Generation for Large Language Models: A Survey

Anonymous ACL submission

Abstract

Large Language Models (LLMs) showcase remarkable abilities, yet they struggle with limitations such as hallucinations, outdated knowledge, opacity, and inexplicable reasoning. To address these challenges, Retrieval-Augmented Generation (RAG) has proven to be a viable solution, leveraging external databases to improve the consistency and coherence of generated content, especially valuable for complex, knowledge-rich tasks, and facilitates continuous improvement by leveraging domain-specific insights. However, RAG is not without its limitations, including a limited context window, irrelevant information, and the high processing overhead for extensive contextual data. In this comprehensive work, we explore the evolution of Contextual Compression paradigms, providing an in-depth examination of the field. We also introduce a state-of-the-art evaluation framework and benchmark. Finally, we outline the current challenges and suggest potential research and development directions, paving the way for future advancements in this area.

1 Introduction

The pioneering accomplishments of large language models (LLMs) have galvanized research initiatives across both industrial and academic spheres. These LLMs showcase their capacity to converse with humans in a natural and articulate manner, excelling across various tasks such as document summarization, Q&A systems, conversational AI, and coding assistants. Despite their advancements, LLMs continue to struggle with tasks that require specialized knowledge or domain-specific expertise. (Kandpal et al., 2023). Notably, they may produce “hallucinations” (Zhang et al., 2023) when confronted with out-of-scope queries or requests that necessitate up-to-date knowledge. To address these challenges, Retrieval-Augmented Generation

(RAG) leverages external knowledge bases to retrieve relevant document snippets, utilizing semantic similarity metrics to identify the most pertinent information. By tapping into external knowledge sources, RAG successfully alleviates the issue of generating inaccurate content, thereby increasing the reliability of LLMs and paving the way for their widespread adoption in real-world applications.

However, RAG also has its challenges. One issue is that when retrieving relevant documents, the important information may be buried in a large amount of irrelevant text, leading to inefficient and poor responses. Another challenge is that current language models have a limited input length, which causes their performance to decline when processing lengthy documents, such as academic articles, research papers, or literary works. This constraint has fueled research into developing methods to increase the input length while maintaining the model’s accuracy and efficiency.

This paper aims to shed light on the latest advancements in contextual compression methods, with a focus on their application in retrieval-based systems. Our research involves a comprehensive review of methodologies, metrics, and benchmarks, which we systematically categorize into a novel taxonomy. Our taxonomy, as shown in Figure 1, presents a structured and comprehensive framework for categorizing and analyzing Contextual Compression techniques for LLMs. Our investigation involves a comprehensive analysis of established techniques, such as semantic compression, in-context auto-encoder compressors, and auto-compressors, among others. Furthermore, our research highlights the ongoing challenges in this field and provides a roadmap for future investigations. We emphasize the need for collective efforts to create a sustainable and environmentally responsible future for LLMs.



Figure 1: Taxonomy of Contextual Compression Methods for Large Language Models.

2 Methods

2.1 Semantic Compression

Semantic compression is a technique that helps identify common patterns of thought in a specific context by generalizing terms. It uses a "domain frequency dictionary" to establish the context and disambiguate multiple possible meanings of words. This approach, based on semantic networks, offers improvements over existing natural language processing techniques.

Semantic compression reduces the number of terms in a text document by replacing less frequent terms with more general terms (their hypernyms) using a semantic network and term frequency data. This compression minimizes information loss and enables efficient processing, especially in tasks involving vector space models (Baeza-Yates et al., 1999), (Erk and Padó, 2008). It also helps address linguistic (Sinha and Mihalcea, 2007) challenges like polysemy and synonymy (Krovetz and Croft, 1992) by replacing multiple rare terms with a single, more general concept. By using statistical analysis and frequency dictionaries, semantic compression can handle polysemic concepts more effectively and with lower error rates than other techniques. These efforts can be summarized into five approaches: *Context Distillation*, *Prompting*, *Efficient Attention Operations*, *Extrapolation and Interpolation*, and *Context Window Extension*.

2.1.1 Context Distillation

Recent studies have demonstrated that augmenting language models (LMs) with contextual information, such as task descriptions, illustrative examples, and explanatory notes (Chen et al., 2021), (Scheurer et al., 2022), can substantially enhance their performance capabilities. This approach can even facilitate zero-shot learning (Wei et al., 2021), (Victor et al., 2022) and enable models to tackle complex tasks by generating sequential reasoning steps (Nye et al., 2021), (Wei et al., 2022), (Zhou et al., 2022).

While LMs perform better with context tokens, this advantage disappears when the tokens are removed. Additionally, processing context tokens requires extra computation, which can be a drawback. The context tokens can also be very long, and it's unclear how to handle them when they exceed the context window size. These limitations are similar to human cognitive limitations (Wason and Evans, 1974), such as struggling with complex tasks and having limited working memory (Baddeley, 1992).

Humans overcome challenges through practice, which allows them to "distill" knowledge into habits and muscle memory. For example, learning to type a phone number becomes automatic with repetition, freeing up conscious reasoning for more complex tasks ¹. This process is essential

¹procedural learning vs. declarative learning - https://en.wikipedia.org/wiki/Procedural_knowledge

for building skills and knowledge, enabling us to tackle increasingly intricate challenges.

Researchers in NLP (Askeel et al., 2021), (Snell et al., 2022) are exploring techniques to fine-tune language models, such as context distillation and "Gisting". Context distillation involves generating "practice" questions, having the model reason step-by-step, and fine-tuning it to predict answers from simpler prompts. This helps the model internalize skills, like step-by-step addition (ref Figure 2). "Gisting" (Mu et al., 2024) compresses instructions into concise key-value attention prefixes, saving computational resources and generalizing well to new tasks. As depicted in Figure 3, the approach involves learning a gist model by incorporating gist tokens during instruction tuning, enabling the model to handle prompt compression and instruction following simultaneously.

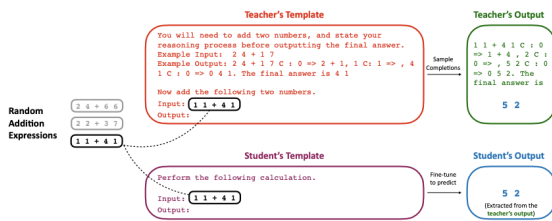


Figure 2: Internalization of step-by-step reasoning via context distillation (Snell et al., 2022)



Figure 3: Gisting - Each vertical rectangle here represents a stack of Transformer activations (Mu et al., 2024)

2.1.2 Prompting

Soft Prompts - As depicted in Figure 4, soft prompt tuning enables the adaptation of pre-trained Transformers without modifying their underlying parameters, as demonstrated in recent studies (Lester et al., 2021), (Zhong et al., 2021), and (Liu et al., 2022). It entails adding novel embeddings to the input sequence and fine-tuning only these new parameters while keeping the remainder of the model's architecture frozen. This approach is categorized as a parameter-efficient fine-tuning method (PEFT) (Lialin et al., 2023), and bears resemblance

to prefix tuning, which prepends task-specific vectors to the attention states instead of the input sequence (Li and Liang, 2021).

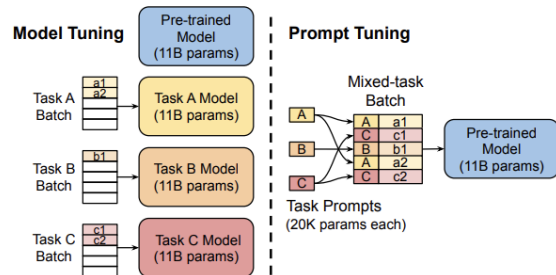


Figure 4: From 11 billion for a tuned model to just 20,480 for a tuned prompt, a reduction of over 5 orders of magnitude (Lester et al., 2021)

Prompt Compression - In their work, (Wingate et al., 2022) hypothesize using a soft prompt sp to compress information from a context ctx . They use a pre-trained LM p_{LM} to generate continuations $cty \sim p_{LM}(\cdot | ctx)$ based on the context, and then calibrate the model's outputs with the soft prompt sf , $p_{LM}(cty | sf)$ to the outputs based on the context ctx , $p_{LM}(cty | ctx)$. They find that soft prompts effectively preserve abstract knowledge and improve guided output. Nevertheless, this method necessitates distinct optimization for each novel context, lacking the ability to leverage knowledge across analogous contexts.

Task-Agnostic Prompt Compression - Current methods for compressing natural language prompts remove tokens or lexical units based on information entropy from a language model like LLaMa-7B. However, using information entropy as a compression metric has two limitations: 1) it only considers unidirectional context, which may miss important information, and 2) it doesn't perfectly align with the goal of prompt compression.

To address these issues, (Pan et al., 2024) propose a data distillation approach to compress prompts while retaining essential information. They introduce an extractive text compression dataset and frame prompt compression as a token classification problem (preserve or discard) (Refer to Figure 5). The key benefits are as follows:

1. *Comprehensive Information Capture:* By leveraging a Transformer encoder, the method captures essential details from the full bidirectional context.
2. *Reduced Latency:* Smaller models explicitly

204
205
206
207

208
209
210
211
212
213

214
215
216
217

218
219
220

221
222
223

224
225
226
227
228
229
230

231
232
233
234
235
236
237
238
239
240

learn the compression objective, leading to lower latency.

3. *Faithfulness*: The compressed prompt remains faithful to the original content.

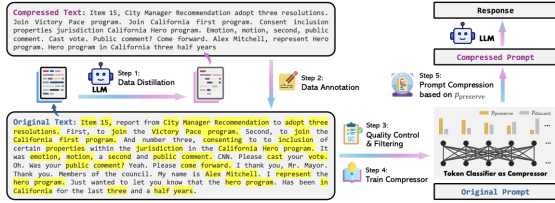


Figure 5: Overview of LLMLingua-2 (Pan et al., 2024)

2.1.3 Efficient Attention Operations

The self-attention mechanism in LLMs leads to an inference cost that scales quadratically with sequence length, prompting the development of various methods to alleviate this complexity. For example:

- *Transformer-XL* (Dai et al., 2019) - employs a recurrent architecture that operates on segments, paired with a novel positional encoding technique.
- *Longformer* (Beltagy et al., 2020) - introduces sparse attention, scaling linearly with sequence length.
- *FlashAttention* (Dao et al., 2022) - uses chunking and re-computation to avoid quadratic attention complexity.

However, these methods can be expensive to train and struggle with out-of-distribution content lengths (Ding et al., 2023). To address this, *LongLoRA* (Chen et al., 2023b) provides a computationally efficient fine-tuning method with minimal resource requirements. For further insights, refer to the study by (Huang et al., 2023).

2.1.4 Extrapolation and Interpolation

In the field of NLP, researchers are investigating methods to extend the capabilities of existing language models, initially trained on brief texts, to process longer sequences during inference (Anil et al., 2022). One approach is to alter positional embeddings, which are typically designed for shorter contexts. The Rotary Position Embeddings (RoPE) from LLaMA is a key foundation for several studies in this area. For example:

- *Position Interpolation (PI)* (Chen et al., 2021) applies a linear transformation to input positional indices.
- *YaRN* (Peng et al., 2023) leverages neural tangent kernel-inspired mechanisms to scale up the context window to 64,000 and 128,000 tokens.

2.1.5 Context Window Extension

Researchers (Fei et al., 2023) propose a semantic compression method that distills long texts into concise forms, retaining their meaning and broadening the context window (Figure 6). This method occurs before inputting tokens into pre-trained language models and is customizable and optimized for specific tasks. It outperforms existing methods in various tasks, including question answering, summarization, and few-shot learning, without requiring additional parameter updates or memory consumption, making it computationally efficient.

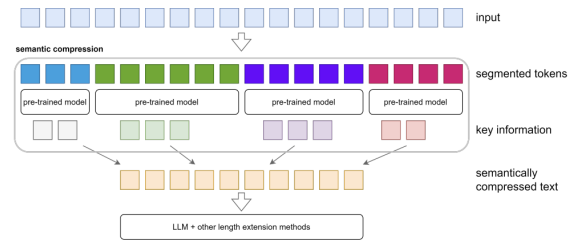


Figure 6: 1) clustering the input text into thematic groups, represented as a graph, to facilitate topic-based analysis, 2) tuning the thematic segments using pre-trained models to preserve crucial details, and 3) re-assembling the refined chunks in their original order - reducing the text length by approximately 6-8 times. Additionally, other techniques like extrapolation and interpolation can be used to further extend the length (Fei et al., 2023)

2.2 Pre-Trained Language Models (PLMs)

The development of PLMs has revolutionized the field of NLP. The first generation of PLMs, such as Skip-Gram (Mikolov et al., 2013b), word2vec (Mikolov et al., 2013a), and GloVe (Pennington et al., 2014), used shallow neural networks (Qiu et al., 2020) to obtain word embeddings. The second generation, including CoVe (McCann et al., 2017), ELMo (Peters et al., 2018), BERT (Devlin et al., 2018), and GPT (Radford et al., 2018), focused on learning dynamic word embeddings using transformers. The pre-training and fine-tuning approach has achieved remarkable success in various NLP tasks. Moreover, recent breakthroughs

in prompt learning (Liu et al., 2023a) have empowered PLMs to accomplish few-shot or zero-shot learning with minimal labeled data. Notable examples of successful PLMs include ChatGPT, GPT-4, Gemini, Claude, LLaMA-3, Mixtral, etc.

2.2.1 AutoCompressors

The authors of (Chevalier et al., 2023) propose teaching PLMs to compress text into summary vectors (Lester et al., 2021), which are significantly shorter than the original text (often 1-2 orders of magnitude shorter). These vectors have a two-pronged function: 1) they allow the LM to handle long documents by extending its context window with minimal computational overhead, and 2) they accelerate inference for pre-computed and cached text.

AutoCompressors, proposed by (Chevalier et al., 2023), are trained To distill key information into summary vectors, generated sequentially from extended documents (Figure 7). The approach builds upon the Recurrent Memory Transformers (RMT) architecture (Bulatov et al., 2022), introducing summary accumulation and training with randomly segmented inputs. This enhances long-range information retention and facilitates reasoning across multiple passages. AutoCompressors can be seeded with PLMs and fine-tuned on long sequences. They improve perplexity for long documents and demonstrate robust compression capabilities across different domains, making them valuable for various downstream applications.

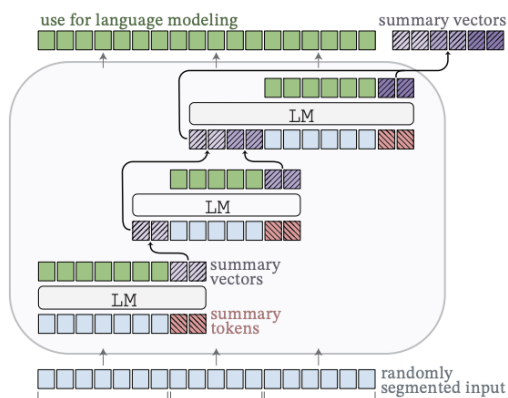


Figure 7: AutoCompressors recursively generate summary vectors from long documents, using them as soft prompts for subsequent segments (Chevalier et al., 2023)

2.2.2 LongNET

Overcoming sequence length limitations in language models has several advantages, including improved interactions with human language, better capture of complex causality and reasoning, and reduced catastrophic forgetting. However, scaling up sequence length poses a challenge in balancing computational complexity and model expressivity. RNN-style models and state space models (Gu et al., 2021), (Smith et al., 2022), (Fu et al., 2022), (Poli et al., 2023) have been proposed, but they have limitations from the perspective of parallelization and model adaptability (Fathi et al., 2023). An alternative approach is to reduce the complexity of Transformers (Vaswani et al., 2017), such as using sliding windows or convolution modules for attention, or sparse attention. LongNet (Ding et al., 2023), a novel approach, replaces the attention mechanism with "dilated attention", which achieves linear computational complexity and logarithmic dependency between tokens. This allows LongNet to efficiently scale sequence lengths to 1 billion tokens, overcoming the constraints of computation and memory.

2.2.3 In-Context Auto-Encoders

Modeling long-range dependencies is a hurdle for Transformer-based LMs (Vaswani et al., 2017) due to their self-attention mechanism. Previous research by (Beltagy et al., 2020), (Bulatov et al., 2022), and Ding (Ding et al., 2023) has attempted to cope with this issue through architectural innovations, but these approaches often struggle to maintain performance in long contexts, as underscored by (Liu et al., 2024). A novel approach, "context compression", is proposed by (Ge et al., 2023), which recognizes that an LLM can represent the same information in varying lengths. They introduce the In-context Autoencoder (ICAE), which compresses lengthy contexts into a fixed number of memory buffers using a learnable encoder and a fixed decoder (Figure 8). The ICAE is pre-trained using auto-encoding and language modeling objectives and fine-tuned using instruction data. The approach achieves 4x context compression while maintaining effective conditioning for the target LLM, enabling faster and more memory-efficient inference.

2.2.4 RECOMP

In their work, (Xu et al., 2024) introduce RECOMP, an intermediary step for Retrieval-augmented Lan-

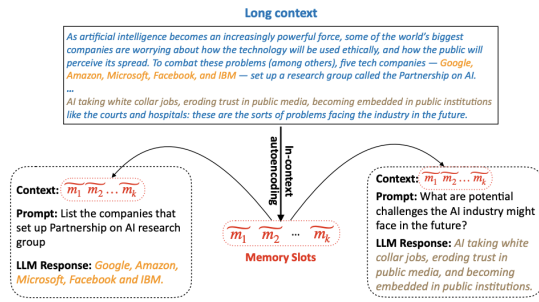


Figure 8: Condensing an extended context into a compact memory representation, which can be leveraged by the target LLM to respond to diverse prompts. (Ge et al., 2023)

guage Models (RALMs) (Izacard et al., 2022), (Borgeaud et al., 2022). RECOMP compresses retrieved documents into concise textual summaries before integrating them during inference, reducing computational costs and alleviating the burden on LMs to process lengthy documents. The aim is to produce summaries that balance brevity and fidelity to the original evidence documents, guiding the RALM to produce targeted outputs when the summary is used as a prefix to the input (illustrated in Figure 9). To achieve this, the authors train two types of compressors:

1. *Extractive Compressor*: This compressor filters out irrelevant sentences, retaining only the most pertinent ones from the retrieved document set.
2. *Abstractive Compressor*: This compressor produces a summary by fusing information from multiple retrieved documents.

Both compressors employ a multi-document query-based summarization approach (Xu and Lapata, 2020), summarizing evidence documents concerning the input query. The authors develop training strategies that maximize performance on the target task to guarantee accurate output. Contrastive learning is employed to train the extractive compressor enabling it to select key sentences effectively, while the abstractive compressor is distilled (West et al., 2021) from a large language model (like GPT-3 or GPT-4), achieving strong summarization performance. This approach holds promise for enhancing the efficiency and efficacy of RALMs.

2.3 Retrievers

The retriever (Chase, 2017-) is an interface that processes an unstructured query and returns a curated

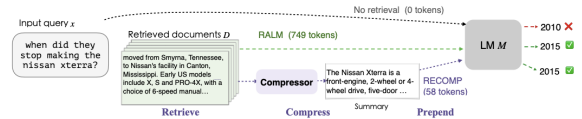


Figure 9: RECOMP’s document compression technique generates a summary that serves as input to a language model, facilitating correct answer generation while minimizing encoding costs. (Xu et al., 2024)

list of documents in response. Contextual compression aims to address the challenges of retrieval by compressing the retrieved context to only include relevant information. In this context, "compressing" encompasses both condensing the content of individual documents and eliminating irrelevant documents altogether. The Contextual Compression Retriever uses a *base retriever* and a *Document Compressor* to process queries. The base retriever retrieves the initial documents, which are then passed through the Document Compressor to shorten the list of documents by either reducing the contents of individual documents or excluding entire documents altogether.

2.3.1 LLMChainExtractor

In this approach, the base retriever is wrapped with a *ContextualCompressionRetriever*. Additionally, an *LLMChainExtractor* serves as the base compressor. The *LLMChainExtractor* iterates over the initially retrieved documents and extracts only the relevant content for the given query. It achieves this by making an additional LLM call for each retrieved document and summarizing the relevant information

2.3.2 EmbeddingsFilter

Making an additional LLM call for each retrieved document can be both costly and slow. However, the *EmbeddingsFilter* offers a more economical and faster alternative. By embedding both the documents and the query, it selectively returns only those documents that exhibit sufficiently similar embeddings to the query. This approach optimizes retrieval efficiency while maintaining relevance.

2.3.3 DocumentCompressorPipeline

The *DocumentCompressorPipeline* allows a seamless combination of multiple compressors in a sequence. Alongside these compressors, we can incorporate *BaseDocumentTransformers* into our pipeline. Unlike contextual compressors, these transformers don’t alter the content significantly

but perform specific transformations on a set of documents. For instance, *TextSplitters* can divide documents into smaller segments, while the *EmbeddingsRedundantFilter* identifies and filters out redundant documents based on embedding similarity. This modular approach enhances flexibility and adaptability in document processing. e.g.

- *Splitter*: create small chunks
- *Redundant filter*: remove similar docs — embedded
- *Relevant filter*: relevant to query

3 Metrics and Benchmarks

3.1 Metrics

Evaluating language model inference efficiency involves considering various metrics that capture different performance aspects, including accuracy, zero-shot capabilities, compression ratio, and inference time. Within the framework of RAG-based solutions, the "Triad of Metrics" ² - Groundedness, Context Relevance, and Answer Relevance - are also employed for evaluation. Achieving satisfactory performance across these metrics helps ensure that the language model application is reliable and free from hallucinations.

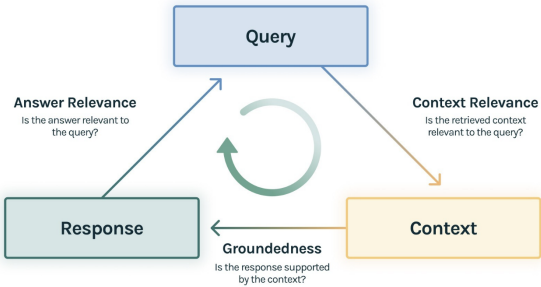


Figure 10: RAG-Triad

3.1.1 Compression Ratio

The compression ratio measures the reduction in size from the original uncompressed context to the compressed context. A higher compression ratio means that the compression is more efficient, as it achieves a greater reduction in size while preserving the context's coherence.

²RAG Triad (Figure 10): https://www.trulens.org/trulens_eval/getting_started/core_concepts/rag_triad/

3.1.2 Inference Time

Inference time, also known as latency, measures how long it takes for a Large Language Model (LLM) to process input data and generate responses. This metric is crucial for real-world applications that require quick handling of user queries or processing of large data volumes in real-time.

3.1.3 Context Relevance

In RAG applications, the first step is retrieval, and it's crucial to ensure that the retrieved context chunks are relevant to the input query. Irrelevant information in the context can lead to hallucinations in the LLM's answer. To evaluate context relevance, the structure of the serialized record can be analyzed.

3.1.4 Groundedness

After retrieving the context, an LLM transforms it into an answer. However, LLMs can sometimes stray from the facts and generate responses that are not entirely accurate. To ensure the groundedness of the application, the response can be broken down into individual claims and verified by searching for supporting evidence within the retrieved context.

3.1.5 Answer Relevance

Furthermore, our response must still effectively address the original question. We can assess this by evaluating the relevance of the final response to the user's input.

3.1.6 Others

RAG evaluation also encompasses four key abilities that reflect the model's adaptability and efficiency: noise robustness, negative rejection, information integration, and counterfactual robustness (Chen et al., 2024), (Liu et al., 2023b). The model's quality scores are heavily influenced by its ability to leverage these capabilities in diverse challenges and complex scenarios:

1. *Noise Robustness*: This metric gauges a model's capacity to distinguish between relevant and irrelevant documents, even when the latter are tangentially related to the question.
2. *Negative Rejection*: The metric measures a model's capacity to recognize when the retrieved documents are insufficient to answer a question, and to withhold a response accordingly.
3. *Information Integration*: Information integration tests a model's proficiency in combining

509	relevant information from multiple documents	Despite its significance, the theoretical and empir-	555
510	to provide well-informed answers to challeng-	ical foundations of this trade-off remain poorly	556
511	ing questions.	understood. Future investigations should focus on	557
512	4. <i>Counterfactual Robustness</i> : Counterfactual	conducting exhaustive examinations to drive the	558
513	robustness measures a model’s skill in identi-	creation of sophisticated compression techniques	559
514	fying and ignoring flawed or misleading infor-	that can meet the demands of increasingly complex	560
515	mation in documents, regardless of its aware-	data sets, enabling researchers to create tailored	561
516	ness of potential errors.	methods that effectively navigate the design space	562
517	In brief, context relevance and noise robustness are	and optimize performance.	563
518	crucial for evaluating the retrieval process, while		
519	answer groundedness, answer relevance, negative	4.3 Dynamic Contextual Compression	564
520	rejection, information integration, and counterfac-	Contemporary compression approaches still utilize	565
521	tual robustness are vital for assessing the quality of	manual compressors, such as retrievers, which of-	566
522	generated text.	ten require an empirical methodology driven by	567
523	3.2 Benchmarks and Datasets	input data or task specifications. This can be a prac-	568
524	The primary objective of these benchmarks and	tical hindrance to adoption, especially in scenarios	569
525	datasets is to assess the trade-offs between com-	like context distillation, where finding suitable stu-	570
526	pressed and uncompressed contexts in terms of	dent templates within computational constraints	571
527	effectiveness, efficiency, and accuracy, covering a	can be time-consuming and require multiple trials.	572
528	broad range of NLP tasks and applications.		
529	3.2.1 Common Benchmarks and Datasets	4.4 Explainability	573
530	RAG’s primary function revolves around answer-	Compressing pre-trained language models can	574
531	ing questions, encompassing various formats such	make them hard to understand (lacking explain-	575
532	as single-hop and multi-hop queries, multiple-	ability). To fix this, using explainable compression	576
533	choice options, and domain-specific inquiries, as	methods can help make models more interpretable,	577
534	well as lengthy scenarios that leverage RAG’s ca-	easier to evaluate, and more reliable in real-life	578
535	pabilities. Moreover, RAG is constantly evolving	scenarios.	579
536	to tackle additional tasks, including extracting rel-		
537	evant information, generating conversational dia-	5 Conclusion	580
538	logue, and searching for code snippets, documenta-	This in-depth analysis explores the domain of con-	581
539	tions and even interpreting them. For more details,	textual compression techniques, with a focus on	582
540	refer to the study by (Gao et al., 2023).	their application to LLMs. Our study encompasses	583
541		a broad range of compression methods, evaluation	584
542	4 Challenges and Future Directions	metrics, and benchmark datasets, providing a com-	585
543	4.1 More advanced Methods	prehensive understanding of the field. By exam-	586
544	Research on contextual compression for LLMs is	ining the complexities of contextual compression,	587
545	still in its early stages. While previous studies have	we identify the key challenges and opportunities	588
546	shown compressed contexts, they still lag behind	that arise in this area. As research in this field	589
547	uncompressed contexts in terms of performance.	continues to advance, the development of special-	590
548	By exploring more advanced compression methods	ized methodologies tailored to the needs of LLMs	591
549	tailored for LLMs, we can potentially bridge this	is crucial for unlocking their full potential across	592
550	performance gap and enhance the performance of	various domains. This survey aims to serve as a	593
551	uncompressed contexts.	valuable resource, providing a detailed overview	594
552	4.2 Performance-Size Trade-offs	of the current landscape and encouraging further	595
553	Previous research highlights the importance of bal-	investigation into this vital topic.	596
554	ancing LLM performance with context size, consid-		
	ering hardware limitations and practical constraints.	Limitations	597
		While this survey provides a comprehensive	598
		overview of contextual compression techniques for	599
		large language models, there are several limitations	600
		to acknowledge. Firstly, the field of contextual	601

602	compression is rapidly evolving, and this survey	Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun.	652
603	may not capture the very latest advancements in	2024. Benchmarking large language models in	653
604	the area. Additionally, the focus on large language	retrieval-augmented generation. In <i>Proceedings of</i>	654
605	models may not be representative of other types of	<i>the AAAI Conference on Artificial Intelligence</i> , vol-	655
606	language models or AI systems, which may have	ume 38, pages 17754–17762.	656
607	different compression requirements. Furthermore,	Shouyuan Chen, Sherman Wong, Liangjian Chen, and	657
608	the survey’s reliance on existing evaluation metrics	Yuandong Tian. 2023a. Extending context window	658
609	and benchmark datasets may not fully capture the	of large language models via positional interpolation.	659
610	complexities and nuances of contextual compres-	<i>arXiv preprint arXiv:2306.15595</i> .	660
611	sion. Moreover, the need for advanced methodolo-	Yanda Chen, Ruiqi Zhong, Sheng Zha, George Karypis,	661
612	gies specifically designed for LLMs highlights the	and He He. 2021. Meta-learning via language model	662
613	potential limitations of current approaches, which	in-context tuning. <i>arXiv preprint arXiv:2110.07814</i> .	663
614	may not be scalable or effective for future LLM	Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai,	664
615	architectures. Finally, the survey’s scope is limited	Zhijian Liu, Song Han, and Jiaya Jia. 2023b. Lon-	665
616	to contextual compression, and future research may	glora: Efficient fine-tuning of long-context large lan-	666
617	uncover new challenges and opportunities at the	guage models. <i>arXiv preprint arXiv:2309.12307</i> .	667
618	intersection of compression and other aspects of	Alexis Chevalier, Alexander Wettig, Anirudh Ajith,	668
619	LLMs.	and Danqi Chen. 2023. Adapting language	669
		models to compress contexts. <i>arXiv preprint</i>	670
		<i>arXiv:2305.14788</i> .	671
620	References	Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Car-	672
621	Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor	bonell, Quoc V Le, and Ruslan Salakhutdinov.	673
622	Lewkowycz, Vedant Misra, Vinay Ramasesh, Am-	2019. Transformer-xl: Attentive language mod-	674
623	brose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam	els beyond a fixed-length context. <i>arXiv preprint</i>	675
624	Neyshabur. 2022. Exploring length generalization in	<i>arXiv:1901.02860</i> .	676
625	large language models. <i>Advances in Neural Informa-</i>	Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and	677
626	<i>tion Processing Systems</i> , 35:38546–38556.	Christopher Ré. 2022. Flashattention: Fast and	678
627	Amanda Askell, Yuntao Bai, Anna Chen, Dawn Drain,	memory-efficient exact attention with io-awareness.	679
628	Deep Ganguli, Tom Henighan, Andy Jones, Nicholas	<i>Advances in Neural Information Processing Systems</i> ,	680
629	Joseph, Ben Mann, Nova DasSarma, et al. 2021. A	35:16344–16359.	681
630	general language assistant as a laboratory for align-	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	682
631	ment. <i>arXiv preprint arXiv:2112.00861</i> .	Kristina Toutanova. 2018. Bert: Pre-training of deep	683
632	Alan Baddeley. 1992. Working memory. <i>Science</i> ,	bidirectional transformers for language understand-	684
633	255(5044):556–559.	ing. <i>arXiv preprint arXiv:1810.04805</i> .	685
634	Ricardo Baeza-Yates, Berthier Ribeiro-Neto, et al. 1999.	Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang,	686
635	<i>Modern information retrieval</i> , volume 463. ACM	Shaohan Huang, Wenhui Wang, Nanning Zheng,	687
636	press New York.	and Furu Wei. 2023. Longnet: Scaling trans-	688
637	Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020.	formers to 1,000,000,000 tokens. <i>arXiv preprint</i>	689
638	Longformer: The long-document transformer. <i>arXiv</i>	<i>arXiv:2307.02486</i> .	690
639	<i>preprint arXiv:2004.05150</i> .	Katrin Erk and Sebastian Padó. 2008. A structured	691
640	Sebastian Borgeaud, Arthur Mensch, Jordan Hoff-	vector space model for word meaning in context. In	692
641	mann, Trevor Cai, Eliza Rutherford, Katie Mill-	<i>Proceedings of the 2008 conference on empirical</i>	693
642	ican, George Bm Van Den Driessche, Jean-Baptiste	<i>methods in natural language processing</i> , pages 897–	694
643	Lespiau, Bogdan Damoc, Aidan Clark, et al. 2022.	906.	695
644	Improving language models by retrieving from tril-	Mahan Fathi, Jonathan Pilault, Pierre-Luc Bacon,	696
645	lions of tokens. In <i>International conference on ma-</i>	Christopher Pal, Orhan Firat, and Ross Goroshin.	697
646	<i>chine learning</i> , pages 2206–2240. PMLR.	2023. Block-state transformer. <i>arXiv preprint</i>	698
647	Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev.	<i>arXiv:2306.09539</i> .	699
648	2022. Recurrent memory transformer. <i>Advances</i>	Weizhi Fei, Xueyan Niu, Pingyi Zhou, Lu Hou, Bo Bai,	700
649	<i>in Neural Information Processing Systems</i> , 35:11079–	Lei Deng, and Wei Han. 2023. Extending context	701
650	11091.	window of large language models via semantic com-	702
651	Harrison Chase. 2017-. <i>LangChain</i> .	pression. <i>arXiv preprint arXiv:2312.09571</i> .	703

704	Daniel Y Fu, Tri Dao, Khaled K Saab, Armin W	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning:	760
705	Thomas, Atri Rudra, and Christopher Ré. 2022.	Optimizing continuous prompts for generation. <i>arXiv</i>	761
706	Hungry hungry hippos: Towards language mod-	<i>preprint arXiv:2101.00190</i> .	762
707	eling with state space models. <i>arXiv preprint</i>		
708	<i>arXiv:2212.14052</i> .		
709	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	Vladislav Lialin, Vijeta Deshpande, and Anna	763
710	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen	Rumshisky. 2023. Scaling down to scale up: A guide	764
711	Wang. 2023. Retrieval-augmented generation for	to parameter-efficient fine-tuning. <i>arXiv preprint</i>	765
712	large language models: A survey. <i>arXiv preprint</i>	<i>arXiv:2303.15647</i> .	766
713	<i>arXiv:2312.10997</i> .		
714	Tao Ge, Jing Hu, Xun Wang, Si-Qing Chen, and Furu	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paran-	767
715	Wei. 2023. In-context autoencoder for context com-	jape, Michele Bevilacqua, Fabio Petroni, and Percy	768
716	pression in a large language model. <i>arXiv preprint</i>	Liang. 2024. Lost in the middle: How language mod-	769
717	<i>arXiv:2307.06945</i> .	els use long contexts. <i>Transactions of the Association</i>	770
718	Albert Gu, Karan Goel, and Christopher Ré. 2021. Effi-	<i>for Computational Linguistics</i> , 12:157–173.	771
719	ciently modeling long sequences with structured state		
720	spaces. <i>arXiv preprint arXiv:2111.00396</i> .	Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang,	772
721	Yunpeng Huang, Jingwei Xu, Zixu Jiang, Junyu Lai,	Hiroaki Hayashi, and Graham Neubig. 2023a. Pre-	773
722	Zenan Li, Yuan Yao, Taolue Chen, Lijuan Yang,	train, prompt, and predict: A systematic survey of	774
723	Zhou Xin, and Xiaoxing Ma. 2023. Advancing trans-	prompting methods in natural language processing.	775
724	former architecture in long-context large language	<i>ACM Computing Surveys</i> , 55(9):1–35.	776
725	models: A comprehensive survey. <i>arXiv preprint</i>		
726	<i>arXiv:2311.12351</i> .	Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengx-	777
727	Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas	iao Du, Zhilin Yang, and Jie Tang. 2022. P-tuning:	778
728	Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-	Prompt tuning can be comparable to fine-tuning	779
729	Yu, Armand Joulin, Sebastian Riedel, and Edouard	across scales and tasks . In <i>Proceedings of the 60th</i>	780
730	Grave. 2022. Atlas: Few-shot learning with retrieval	<i>Annual Meeting of the Association for Computational</i>	781
731	augmented language models. <i>arXiv preprint</i>	<i>Linguistics (Volume 2: Short Papers)</i> , pages 61–68,	782
732	<i>arXiv:2208.03299</i> .	Dublin, Ireland. Association for Computational Lin-	783
733	Huiqiang Jiang, Qianhui Wu, , Xufang Luo, Dongsheng	guistics.	784
734	Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2023a.	Yi Liu, Lianzhe Huang, Shicheng Li, Sishuo Chen, Hao	785
735	LongLLMLingua: Accelerating and enhancing llms	Zhou, Fandong Meng, Jie Zhou, and Xu Sun. 2023b.	786
736	in long context scenarios via prompt compression.	Recall: A benchmark for llms robustness against	787
737	<i>ArXiv preprint</i> , abs/2310.06839.	external counterfactual knowledge. <i>arXiv preprint</i>	788
738	Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing	<i>arXiv:2311.08147</i> .	789
739	Yang, and Lili Qiu. 2023b. LLMLingua: Compress-	Bryan McCann, James Bradbury, Caiming Xiong, and	790
740	ing prompts for accelerated inference of large lan-	Richard Socher. 2017. Learned in translation: Con-	791
741	guage models . In <i>Proceedings of the 2023 Confer-</i>	textualized word vectors. <i>Advances in neural infor-</i>	792
742	<i>ence on Empirical Methods in Natural Language Pro-</i>	<i>mation processing systems</i> , 30.	793
743	<i>cessing</i> , pages 13358–13376. Association for Com-		
744	putational Linguistics.	Tomas Mikolov, Kai Chen, Greg Corrado, and Jef-	794
745	Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric	frey Dean. 2013a. Efficient estimation of word	795
746	Wallace, and Colin Raffel. 2023. Large language	representations in vector space. <i>arXiv preprint</i>	796
747	models struggle to learn long-tail knowledge. In <i>In-</i>	<i>arXiv:1301.3781</i> .	797
748	<i>ternational Conference on Machine Learning</i> , pages	Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Cor-	798
749	15696–15707. PMLR.	rado, and Jeff Dean. 2013b. Distributed representa-	799
750	Robert Krovetz and W Bruce Croft. 1992. Lexical am-	tations of words and phrases and their compositionality.	800
751	biguity and information retrieval. <i>ACM Transactions</i>	<i>Advances in neural information processing systems</i> ,	801
752	<i>on Information Systems (TOIS)</i> , 10(2):115–141.	26.	802
753	Brian Lester, Rami Al-Rfou, and Noah Constant. 2021.	Jesse Mu, Xiang Li, and Noah Goodman. 2024. Learn-	803
754	The power of scale for parameter-efficient prompt	ing to compress prompts with gist tokens. <i>Advances</i>	804
755	tuning . In <i>Proceedings of the 2021 Conference on</i>	<i>in Neural Information Processing Systems</i> , 36.	805
756	<i>Empirical Methods in Natural Language Processing</i> ,	Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari,	806
757	pages 3045–3059, Online and Punta Cana, Domini-	Henryk Michalewski, Jacob Austin, David Bieber,	807
758	can Republic. Association for Computational Lin-	David Dohan, Aitor Lewkowycz, Maarten Bosma,	808
759	guistics.	David Luan, et al. 2021. Show your work: Scratch-	809
		pads for intermediate computation with language	810
		models. <i>arXiv preprint arXiv:2112.00114</i> .	811
		Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia,	812
		Xufang Luo, Jue Zhang, Qingwei Lin, Victor Ruhle,	813
		Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu,	814

815	and Dongmei Zhang. 2024. LLMLingua-2: Data distillation for efficient and faithful task-agnostic prompt compression . <i>ArXiv preprint</i> , abs/2403.12968.	869
816		870
817		871
818	Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. Yarn: Efficient context window extension of large language models. <i>arXiv preprint arXiv:2309.00071</i> .	872
819		873
820		
821		
822	Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In <i>Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)</i> , pages 1532–1543.	874
823		875
824		876
825		877
826		878
827	Matthew E Peters, Mark Neumann, Luke Zettlemoyer, and Wen-tau Yih. 2018. Dissecting contextual word embeddings: Architecture and representation. <i>arXiv preprint arXiv:1808.08949</i> .	879
828		880
829		881
830		882
831	Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua Bengio, Stefano Ermon, and Christopher Ré. 2023. Hyena hierarchy: Towards larger convolutional language models. In <i>International Conference on Machine Learning</i> , pages 28043–28078. PMLR.	883
832		884
833		885
834		886
835		887
836		888
837	Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. 2020. Pre-trained models for natural language processing: A survey. <i>Science China Technological Sciences</i> , 63(10):1872–1897.	889
838		890
839		891
840		892
841		893
842	Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training.	894
843		895
844		896
845	Jérémy Scheurer, Jon Ander Campos, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. 2022. Learning from natural language feedback. In <i>ACL Workshop on Learning with Natural Language Supervision</i> .	897
846		898
847		899
848		900
849		901
850	Ravi Sinha and Rada Mihalcea. 2007. Unsupervised graph-based word sense disambiguation using measures of word semantic similarity. In <i>International conference on semantic computing (ICSC 2007)</i> , pages 363–369. IEEE.	902
851		903
852		904
853		905
854		906
855	Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. 2022. Simplified state space layers for sequence modeling. <i>arXiv preprint arXiv:2208.04933</i> .	907
856		908
857		909
858	Charlie Snell, Dan Klein, and Ruiqi Zhong. 2022. Learning by distilling context. <i>arXiv preprint arXiv:2209.15189</i> .	910
859		911
860		912
861	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	913
862		914
863		915
864		916
865		917
866	Sanh Victor, Webson Albert, Raffel Colin, Bach Stephen, Sutawika Lintang, Alyafeai Zaid, Chaffin Antoine, Stiegler Arnaud, Raja Arun, Dey Manan, et al. 2022. Multitask prompted training enables zero-shot task generalization. In <i>International Conference on Learning Representations</i> .	918
867		919
868		920
		921
		922
	Peter C Wason and J St BT Evans. 1974. Dual processes in reasoning? <i>Cognition</i> , 3(2):141–154.	
	Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. <i>arXiv preprint arXiv:2109.01652</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	
	Peter West, Chandra Bhagavatula, Jack Hessel, Jena D Hwang, Liwei Jiang, Ronan Le Bras, Ximing Lu, Sean Welleck, and Yejin Choi. 2021. Symbolic knowledge distillation: from general language models to commonsense models. <i>arXiv preprint arXiv:2110.07178</i> .	
	David Wingate, Mohammad Shoeybi, and Taylor Sorensen. 2022. Prompt compression and contrastive conditioning for controllability and toxicity reduction in language models. <i>arXiv preprint arXiv:2210.03162</i> .	
	Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. RECOMP: Improving retrieval-augmented LMs with context compression and selective augmentation . In <i>The Twelfth International Conference on Learning Representations</i> .	
	Yumo Xu and Mirella Lapata. 2020. Coarse-to-fine query focused multi-document summarization. In <i>Proceedings of the 2020 Conference on empirical methods in natural language processing (EMNLP)</i> , pages 3632–3645.	
	Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. 2023. Siren’s song in the ai ocean: a survey on hallucination in large language models. <i>arXiv preprint arXiv:2309.01219</i> .	
	Zexuan Zhong, Dan Friedman, and Danqi Chen. 2021. Factual probing is [MASK]: Learning vs. learning to recall . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5017–5033, Online. Association for Computational Linguistics.	
	Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. 2022. Least-to-most prompting enables complex reasoning in large language models. <i>arXiv preprint arXiv:2205.10625</i> .	

923 Wangchunshu Zhou, Yuchen Eleanor Jiang, Peng Cui,
924 Tiannan Wang, Zhenxin Xiao, Yifan Hou, Ryan Cot-
925 terell, and Mrinmaya Sachan. 2023. Recurrentgpt:
926 Interactive generation of (arbitrarily) long text. *arXiv*
927 *preprint arXiv:2305.13304*.