

---

# Beyond Proxy Metrics: A New Evaluation Framework for LLM Compression by Directly Measuring Generative Faithfulness

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 Existing evaluation frameworks for Large Language Model (LLM) compression  
2 rely on proxy metrics like Perplexity and curated benchmarks, leading to a signifi-  
3 cant gap between reported scores and real-world performance. To bridge this gap,  
4 we introduce a new evaluation paradigm that abandons proxies to directly measure  
5 generative faithfulness. Our framework has two core innovations: (1) Conditional  
6 Generation Accuracy (CGA), a metric that evaluates a compressed model’s ability  
7 to replicate an original model’s next-token predictions under a teacher-forcing  
8 regime, and (2) an evaluation set of real user queries, ensuring alignment with  
9 practical applications. Using this framework, we conduct a large-scale evaluation  
10 of 9 mainstream compression techniques across model sizes (7B–32B) and con-  
11 text lengths (8K–24K). Our findings challenge prior claims: KV cache dropping  
12 methods severely underperform, while a sparse attention baseline unexpectedly  
13 surpasses popular quantization techniques. Moreover, we uncover distinct scaling  
14 laws: quantization accuracy degrades sharply with longer contexts, whereas the  
15 fidelity of sparse attention improves. To promote transparent and reproducible  
16 research, we will release a public leaderboard. Click [anonymous github](#) to access  
17 our implementation of all the compression methods and the evaluation pipeline for  
18 all the metrics discussed.

19 

## 1 Introduction

20 The rapid adoption of Large Language Models (LLMs) has created an urgent demand for efficient  
21 inference under constrained GPU resources. Model compression has emerged as a promising solution,  
22 with many training-free or calibration-light approaches claiming substantial inference speedups with  
23 only marginal accuracy loss. However, a persistent disconnect remains between reported benchmark  
24 performance and real-world user experience. This discrepancy motivates a critical re-evaluation of  
25 existing evaluation frameworks for LLM compression.

26 Current evaluation practices predominantly rely on Perplexity (PPL) and Question Answering (QA)  
27 benchmarks, both of which suffer from fundamental limitations.

28 **Limitations of Perplexity.** Perplexity is plagued by two key issues. First, its unbounded value range  
29 and model-specific baselines make fair cross-model comparison difficult. More importantly, PPL is  
30 easily exploitable, certain compression strategies, such as KV cache dropping, can be over-optimized  
31 to perform well on PPL while producing outputs that diverge significantly from the original model,  
32 thereby undermining practical usability. Our experiments confirm that such methods may excel in  
33 PPL metrics yet fail to preserve output fidelity.

34 **Limitations of QA Benchmarks.** In response to the shortcomings of PPL, recent works have shifted  
35 toward QA-style downstream evaluations. While these efforts attempt to expand the evaluation space  
36 by incorporating additional dimensions (e.g., fairness, privacy, commonsense reasoning), they remain  
37 fundamentally constrained by rigid formats such as multiple-choice or binary classification. These  
38 formats poorly reflect the nature of real-world model interactions, which often involve generating  
39 long-form, coherent, and contextually grounded text, especially in multi-step reasoning tasks. More-  
40 over, reliance on answer accuracy introduces new biases and may incentivize overfitting to narrow  
41 evaluation schemas.

42 In summary, the core problem lies in the reliance on *proxy metrics* (e.g., PPL, multiple-choice  
43 accuracy) and *proxy data distributions* (e.g., curated benchmarks). This proxy-based paradigm is  
44 the root cause of the mismatch between benchmark scores and real-world utility and poses a serious  
45 barrier to meaningful progress in LLM compression.

46 To address this issue, we introduce a novel evaluation framework that eliminates proxies in favor of  
47 direct, fidelity-based evaluation grounded in real-world usage scenarios. Our framework is built on  
48 two core innovations:

49 **Conditional Generation Accuracy (CGA): A Direct Evaluation Metric.** We propose CGA as  
50 a robust and interpretable metric for measuring the fidelity of compressed models. Unlike prior  
51 work that relied on edit distance or BERTScore, which are affected by output length sensitivity and  
52 pretraining bias respectively, CGA directly measures a compressed model’s ability to replicate the  
53 original model’s output distribution. Under a teacher-forcing paradigm, we feed the ground-truth  
54 response (generated by the original model) to the compressed model token-by-token, and evaluate  
55 prediction accuracy for the next token (Acc@1, Acc@5). This approach provides a precise and  
56 faithful assessment of how well the compressed model preserves the original generation behavior.

57 **Real User Queries as the Evaluation Distribution.** Instead of synthetic or curated benchmarks,  
58 we utilize authentic user queries sourced from public platforms like ShareGPT. These queries better  
59 reflect actual deployment scenarios—they are open-ended, diverse, and often lack a single “correct”  
60 answer. By using the original model’s output as the reference, our framework remains effective even  
61 in ambiguous, multi-modal contexts. To enable a more granular analysis, we leverage GPT-4o [xx] to  
62 classify tens of thousands of these queries into distinct domains such as mathematics, programming,  
63 medicine, and summarization. In parallel, we stratify these queries by context length. This dual  
64 categorization allows us to assess how well a compressed LLM preserves its capabilities not only in  
65 specific areas of expertise but also across different input scales.

66 Using this framework, we conducted a large-scale, head-to-head evaluation of 10 mainstream LLM  
67 compression methods spanning four categories. Our pipeline supports model sizes from 7B to 32B  
68 and context lengths from a few dozen tokens to 24K, yielding scores that are inherently bounded  
69 within [0, 1], ensuring robust comparability. Key findings include:

## Key Takeaways

### Cross-Method Insights

- Overall performance ranking: low-precision attention > INT4 quantization > 50% pruning >> KV cache dropping.
- Differences between categories are far greater than within-category variation.
- Among INT4 quantization methods, GPTQ consistently outperforms AWQ in both accuracy and generalization.
- KV cache dropping methods (e.g., SnapKV, H2O) underperform significantly, often producing outputs that are poorly aligned with the original model, raising concerns about their practical viability.
- Notably, Top-10% attention, a sparse attention approximation, surpasses INT4 quantization methods and FlashAttention FP8 in performance.

### Scalability with Model Size

- Low-precision attention (FlashAttention FP8, SageAttention) and INT4 quantization (GPTQ, AWQ) benefit substantially from larger model scales.
- Unstructured pruning (SparseGPT, Wanda) and Top-10% attention exhibit consistent performance across model sizes.
- KV cache dropping methods show a significant preference for smaller models.

### Scalability with Long Contexts

- All INT4 quantization methods degrade sharply in accuracy as context length increases.
- FlashAttention FP8 deteriorates significantly with long contexts, while Top-10% Attention maintains or improves accuracy.
- KV Cache compression methods degrade markedly as context grows, contradicting their claims of long-context suitability.
- Pruning-based methods (SparseGPT, Wanda) and SageAttention show the most stable performance across varying context lengths.

70

71 To promote transparency and accelerate progress in the field, we will open-source our evaluation  
72 pipeline and release a public leaderboard upon paper acceptance. We hope this work provides a  
73 rigorous, interpretable, and scalable foundation for fair benchmarking in LLM compression research.

## 74 2 Related Works

75 **LLM Compression.** The deployment of LLMs on resource-constrained hardware has catalyzed  
76 extensive research into model compression. Key lossy techniques include low-precision attention,  
77 quantization, pruning, and KV cache compression, with state-of-the-art methods in each category  
78 often claiming significant efficiency gains at minimal performance cost.

79 *Low-Precision Attention.* The advent of specialized hardware like NVIDIA’s Hopper architecture has  
80 made FP8 a viable option for accelerating attention mechanisms. FlashAttention utilizes IO-aware  
81 algorithms to speed up exact attention computation, and its FP8 implementation offers substantial  
82 throughput gains. However, this can introduce noise that degrades performance on sensitive down-  
83 stream tasks. In response, hybrid-precision methods like SageAttention have been proposed, which  
84 selectively use INT8 or INT4 for query and key (QK) projections while retaining FP8 for value and  
85 output (VO) computations to preserve accuracy.

86 Another line of work focuses on sparsity. Methods like MoBA and NSA use heuristics to approximate  
87 attention by focusing only on a subset of KV blocks. While practical, these are generally considered  
88 to be outperformed by Top-K Attention, which, despite not offering a direct speedup, serves as a  
89 valuable upper bound on the performance of sparse attention methods by using true scores to select  
90 the most relevant KV blocks.

91     *Quantization.* Quantization has become a cornerstone of LLM compression. GPTQ is a widely-  
92     used post-training quantization (PTQ) method that leverages approximate second-order Hessian  
93     information to perform layer-wise weight quantization, achieving high accuracy at INT4 and INT3  
94     bit-widths. While effective, the quantization process itself is computationally intensive. In contrast,  
95     AWQ proposes a faster method that first scales weights based on the magnitude of their corresponding  
96     activations before quantization. This dependence on activation statistics, however, can make AWQ  
97     more sensitive to the distribution of the calibration data, potentially impacting its generalization  
98     capabilities.

99     *Pruning.* Pruning methods aim to improve efficiency by removing redundant model parameters.  
100    While simple magnitude-based pruning is a common baseline, it risks removing weights crucial for  
101    model performance. To mitigate this, more sophisticated methods have been developed. SparseGPT  
102    introduced a calibration process to achieve high levels of sparsity (e.g., 60%) with less accuracy loss.  
103    Wanda refines this by removing weights with the smallest magnitudes after multiplying them by  
104    the norms of their corresponding input activations, providing a more robust measure of a weight’s  
105    contribution.

106    *KV Cache Dropping.* To accelerate the decoding phase of generative inference, several methods focus  
107    on compressing the KV cache. H2O identifies and retains the most influential “heavy-hitter” KV  
108    pairs based on their cumulative attention scores, evicting less important pairs to stay within a fixed  
109    token budget. Building on this, SnapKV further compresses the prompt’s KV cache, aiming to reduce  
110    memory and computational overhead even more.

111    **Compressed LLMs Evaluation.** The evaluation of compressed LLMs has evolved rapidly. Early  
112    work predominantly relied on PPL, but this metric has been widely criticized for its susceptibility to  
113    being “gamed” and its poor correlation with downstream task performance.

114    This led to a shift towards using comprehensive QA benchmarks. Researchers have proposed  
115    increasingly complex evaluation suites that measure performance across multiple dimensions, such  
116    as fairness, privacy, and ethics, as seen in works like [1] and [4]. Others, such as LLM-Kick [2], have  
117    focused on aggregating numerous existing QA tasks into a single, extensive leaderboard. However,  
118    these benchmarks are typically limited to structured formats like multiple-choice questions. This  
119    format diverges significantly from the open-ended, generative nature of real-world user interactions,  
120    introducing a new form of evaluation bias and creating the risk that compression methods will be  
121    over-optimized for these specific test formats.

122    Sharing our motivation to move beyond proxy metrics like PPL and QA benchmarks, a concurrent  
123    work [3] directly compares the outputs of compressed and original models. This approach has  
124    explored metrics such as edit distance and BERTScore, but both have inherent drawbacks. Edit  
125    distance is unreliable for longer generations, as minor initial errors can cascade into large, misleading  
126    divergence scores. BERTScore, while designed to capture semantic similarity, introduces its own set  
127    of potential biases by relying on an external pretrained model.

128 **3 A Direct, Fidelity-Based Evaluation Framework**

129 Let  $\mathcal{F}$  and  $\tilde{\mathcal{F}}$  denote an original LLM and its compressed counterpart, respectively. For a given user  
130 prompt  $X$ , the models generate output sequences  $Y = \mathcal{F}(X)$  and  $\tilde{Y} = \tilde{\mathcal{F}}(X)$ . An ideal evaluation  
131 would measure the "quality" of  $\tilde{\mathcal{F}}$  by comparing  $\tilde{Y}$  to  $Y$  over a distribution of prompts  $X$  that reflects  
132 real-world usage.

133 Existing evaluation frameworks fall short on two fronts: the data distribution and the comparison  
134 metric.

135 **Proxy Data Distributions.** Current benchmarks rely on proxy data distributions that fail to capture  
136 the complexity of real user interactions. QA benchmarks, for instance, use structured prompts ( $X_{QA}$ )  
137 that require selecting from a fixed set of answers, e.g.,  $\{A, B, C\}$ . This constrained format cannot  
138 assess a model's performance on the ambiguous, open-ended, and multi-turn conversational prompts  
139 common in practice. PPL benchmarks often use large text corpora (e.g., WikiText) which, while vast,  
140 are stylistically monolithic and unrepresentative of conversational queries. This mismatch between  
141 benchmark data and real-world data introduces a fundamental evaluation bias.

142 **Proxy Comparison Metrics.** Instead of directly comparing the output sequences  $Y$  and  $\tilde{Y}$ , prevailing  
143 methods rely on proxy metrics. The final score  $S$  in these frameworks is not a direct measure  
144 of similarity, but a comparison of proxy values derived from the model outputs. For PPL, the score  
145 compares the perplexity values, not the text itself:

$$S_{PPL} \propto \sum_{X_{PPL}} \text{sim}(\text{PPL}(\mathcal{F}(X_{PPL})), \text{PPL}(\tilde{\mathcal{F}}(X_{PPL}))) \quad (1)$$

146 Similarly, QA benchmarks compare final scores on a task, abstracting away the generative process  
147 entirely:

$$S_{QA} \propto \sum_{X_{QA}} \mathbf{1}(\text{score}(\mathcal{F}(X_{QA})) = \text{score}(\tilde{\mathcal{F}}(X_{QA}))) \quad (2)$$

148 The use of these proxies—perplexity and multiple-choice accuracy—obscures the evaluation. A  
149 compressed model might achieve a similar PPL or QA score to the original while generating text that  
150 is factually incorrect, stylistically different, or incoherent. This indirection is the primary source of  
151 the gap between benchmark results and real world user experience.

152 **3.1 Aligning with Real-World Data Distributions**

153 To ensure our evaluation reflects practical use cases, we curated a dataset from real-world user  
154 interactions. We sourced tens of thousands of conversations from ShareGPT, a platform containing  
155 diverse, multilingual user prompts. To facilitate a granular analysis of model capabilities, we further  
156 processed this data in two ways:

157 **Domain-Specific Subsets.** We employed GPT-4o to classify and filter the prompts into ten distinct  
158 categories. These include five knowledge-intensive domains (mathematics, programming, business,  
159 law, and medicine), three non-English languages (Japanese, French, and Chinese), and two specific  
160 tasks (summarization and translation). Each category contains approximately 100 high-quality  
161 prompts, allowing for a targeted assessment of performance in specialized areas.

162 **Long-Context Strata.** To analyze performance under long-context scenarios, we segment the data  
163 into three distinct length-based tiers: 8K, 16K, and 24K tokens. Each tier consists of approximately  
164 20 prompts, creating a testbed that mirrors the demands of real-world, long-context applications.

165 This dual-strategy approach ensures our evaluation data is not only aligned with the distribution  
166 of genuine user queries but also structured to yield precise insights into model performance across  
167 various domains and context scales.

168 **3.2 Conditional Generation Accuracy: A Direct Fidelity Metric**

169 The ultimate measure of a compressed model's quality is its ability to replicate the generative behavior  
170 of the original model. Existing proxy metrics, such as task-specific scores, fail to capture the nuanced

171 aspects of generation, including tone, style, and the reasoning process embedded in the output  
172 sequence.

173 To overcome this, we introduce Conditional Generation Accuracy (CGA), a metric that directly  
174 compares the output distributions of the compressed and original models. The standard auto-regressive  
175 generation process for the original model  $\mathcal{F}$  and a compressed model  $\tilde{\mathcal{F}}$  at step  $i$  is given by:

$$y_i = \mathcal{F}(X, y_1, y_2, \dots, y_{i-1}) \quad \text{and} \quad \tilde{y}_i = \tilde{\mathcal{F}}(X, \tilde{y}_1, \tilde{y}_2, \dots, \tilde{y}_{i-1}), \quad (3)$$

176 where a naive comparison of  $y_i$  and  $\tilde{y}_i$  is unreliable, as a single-token divergence early in the sequence  
177 can lead to a cascade of ‘‘errors’’, making the remainder of the generated text entirely different. This  
178 error accumulation problem renders simple metrics like edit distance or BERTScore ineffective for  
179 long sequences.

180 To isolate the predictive capability of the compressed model at each step, we employ a teacher-forcing  
181 strategy. Instead of conditioning the compressed model’s next-token prediction on its own previously  
182 generated tokens ( $\tilde{y}_{<i}$ ), we condition it on the ground-truth sequence from the original model ( $y_{<i}$ ):

$$\hat{y}_i = \tilde{\mathcal{F}}(X, y_1, y_2, \dots, y_{i-1}). \quad (4)$$

183 Here,  $\hat{y}_i$  represents the compressed model’s prediction at step  $i$  given the ideal context from the  
184 original model. This approach decouples the evaluation at each step from errors made in prior steps,  
185 eliminating the issue of cascading failures and allowing for a pure measure of the compressed model’s  
186 predictive fidelity.

187 The final CGA score is the mean accuracy over all tokens in the sequence, averaged across all prompts  
188 in the dataset:

$$S_{\text{CGA}} = \frac{1}{|D|} \sum_{X \in D} \frac{1}{|Y|} \sum_{i=1}^{|Y|} \mathbf{1}(\hat{y}_i = y_i) = \frac{1}{|D|} \sum_{X \in D} \frac{1}{|Y|} \sum_{i=1}^{|Y|} \mathbf{1}(\tilde{\mathcal{F}}(X, y_{<i}) = \mathcal{F}(X, y_{<i})). \quad (5)$$

189 As the formula shows, CGA directly measures the alignment between the compressed and original  
190 models’ next-token predictions under the golden output distribution. It uses no proxies, providing the  
191 most direct and fundamental assessment of generative faithfulness.

## 192 4 Experiments

193 Our experimental evaluation is structured as follows. First, we validate our proposed Conditional  
194 Generation Accuracy (CGA) metric by benchmarking it against standard metrics and human-aligned  
195 judgments (Section 4.1). Second, we conduct a comprehensive, head-to-head comparison of all  
196 selected compression methods to assess their overall performance (Section 4.2). Finally, we analyze  
197 the scaling properties of each method with respect to both model size (7B to 32B) and context length  
198 (8K to 24K) in Section 4.3 and Section 4.4, respectively.

Table 1: Validation of evaluation metrics on a 7B model. While proxy metrics like PPL and MMLU often fail to reflect true performance degradation, our proposed CGA score aligns closely with human-aligned GPT Scores. Misleading scores from proxy metrics are underlined.

Metric	INT4 Quant		Low-Precision Attn			KV Cache Comp.		Pruning		Baseline
	AWQ	GPTQ	Flash FP8	Sage	Top-10%	SnapKV	H2O	Sparse	Wanda	
PPL	17.76	16.79	116.27	16.00	15.99	15.57	18.56	21.37	20.89	15.99
MMLU	0.718	0.718	0.570	0.718	0.718	0.401	0.680	0.651	0.658	0.719
GPT Score	0.833	0.849	0.042	0.882	0.895	0.042	0.000	0.289	0.303	1.0
PPL (norm)	0.172	0.454	<u>0.000</u>	0.997	1.000	<u>1.000</u>	0.077	0.005	0.008	1.0
MMLU (norm)	<u>0.999</u>	<u>0.999</u>	0.793	<u>0.999</u>	<u>0.999</u>	0.558	0.946	0.907	0.915	1.0
CGA (ours)	0.916	0.922	0.586	0.987	0.953	0.538	0.465	0.784	0.737	1.0

### 199 4.1 Validating the CGA Metric

200 We begin by validating our central claim: that CGA aligns more closely with human-perceived quality  
201 than do traditional proxy metrics like Perplexity (PPL) or QA accuracy. For this analysis, we used a  
202 7B model and employed GPT-4o judgments as a high-quality proxy for human perception.

Figure 1: Correlation matrix of evaluation metrics. Darker cells indicate stronger positive correlation. Our proposed Conditional Generation Accuracy (CGA) demonstrates the highest correlation (0.923) with the human preference benchmark (GPT Score), significantly exceeding that of proxy metrics like PPL and MMLU.

203 **Evaluation Setup**

- 204 • **PPL:** We calculated the average score on the WikiText-2 dataset. To facilitate comparison, raw  
205 PPL scores were normalized to a [0, 1] range by applying a sigmoid function to the PPL difference  
206 between the compressed and original models.
- 207 • **QA Score:** We used the MMLU benchmark in a zero-shot setting, evaluated via lm-eval-harness  
208 (v0.4.9). The final score is the ratio of the compressed model’s accuracy to the baseline’s.
- 209 • **GPT Score:** We prompted GPT-4o to perform a pairwise comparison between the outputs of the  
210 original and compressed models for a given prompt, with their presentation order randomized. The  
211 final score is the win rate of the compressed model. The prompt template is available in Appendix  
212 [xx].

213 **Results.** As shown in Table [xx], both PPL and MMLU scores can be misleading. PPL, in particular,  
214 shows poor correlation with GPT Score, especially for KV cache compression methods, which it  
215 rates favorably despite significant output degradation. While the MMLU score correctly identifies  
216 poorly performing methods, it quickly saturates near 1.0 for most methods, failing to distinguish  
217 between moderately and highly faithful models. For instance, H2O achieves an MMLU score of 0.94,  
218 masking a substantial drop in generative quality as judged by GPT-4o.

219 Further analysis in Table [xx] shows the Pearson correlation coefficient between each metric and the  
220 GPT Score. The results confirm that CGA has a significantly stronger correlation with human-aligned  
221 judgments than PPL or MMLU. While the GPT Score itself is a useful proxy, its inherent variance,  
222 prompt sensitivity, and high cost make it unsuitable as a scalable, primary evaluation metric. CGA, in  
223 contrast, offers a robust, reproducible, and well-correlated alternative.

224 — Define custom colors for ranked highlighting —

Table 2: Refined comparison of 9 compression methods on Qwen2.5-7B-Instruct. Performance is measured by Conditional Generation Accuracy (CGA). In each row, cells are highlighted to show the performance ranking: 1st , 2nd , and 3rd .

Category	Sub-dataset	Compression Method								
		AWQ	GPTQ	Flash FP8	Sage	Top-10%	SnapKV	H2O	Sparse	Wanda
Reasoning	code	0.9228	0.9504	0.4576	0.9844	0.9472	0.7466	0.6323	0.8150	0.8165
	math	0.9567	0.9426	0.5457	0.9829	0.9671	0.6350	0.6193	0.6455	0.8752
Knowledge	fact	0.8773	0.8789	0.6078	0.9874	0.9010	0.6403	0.5217	0.7855	0.7486
	law	0.9023	0.9055	0.1899	0.9850	0.9205	0.6593	0.5782	0.7635	0.7755
Language	business	0.8994	0.9102	0.6753	0.9861	0.9401	0.6456	0.5929	0.7511	0.7700
	medicine	0.9169	0.9299	0.4307	0.9795	0.9619	0.7007	0.6070	0.7929	0.8063
Skills	fr	0.9165	0.9319	0.3199	0.9874	0.9425	0.7171	0.5794	0.7706	0.7811
	ch	0.8986	0.9094	0.6236	0.9866	0.9291	0.6876	0.4982	0.7289	0.7496
Skills	jp	0.8914	0.8935	0.7486	0.9874	0.9254	0.6608	0.5301	0.8547	0.8827
	en2zh	0.9497	0.9341	0.7273	0.9874	0.9801	0.4152	0.4812	0.8086	0.8313
	zh2en	0.9383	0.9608	0.7158	0.9809	0.9824	0.4800	0.4647	0.8415	0.8939
<b>Average</b>		0.9091	0.9212	0.5796	0.9860	0.9428	0.5997	0.5709	0.7797	0.8200

225 **4.2 Head-to-Head Comparison of Compression Methods**

226 We next conducted a large-scale benchmark of 9 compression methods on a 7B model using our CGA  
227 framework. This direct comparison provides a clear overview of their relative performance. The  
228 complete results are summarized in Table ??, with corresponding results for 14B and 32B models  
229 available in Appendix [xx].

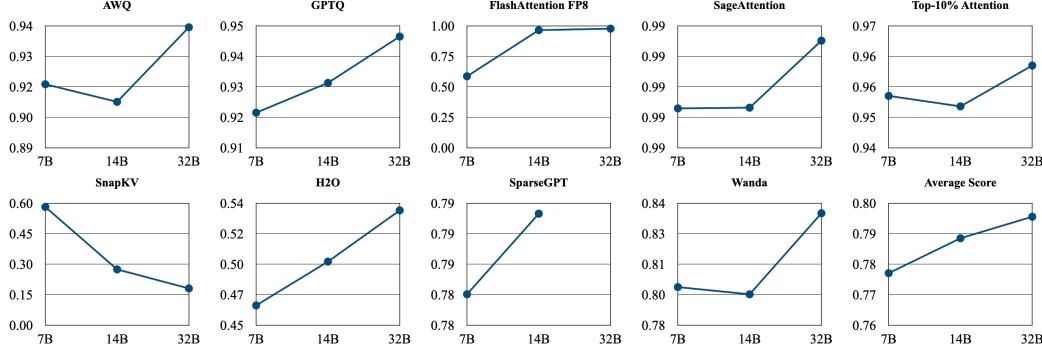


Figure 2: Caption

230 The results reveal several key insights:

- 231 • There are significant performance gaps between method categories. KV cache compression  
232 methods, in particular, dramatically underperform relative to their published claims. At the default  
233 20% budget, their CGA scores are the lowest, suggesting their generalization capabilities warrant  
234 reassessment.
- 235 • The comparison between FlashAttention FP8 and SageAttention confirms that naively replacing  
236 BF16 with FP8 for all attention computations leads to a severe performance drop. SageAtten-  
237 tion’s hybrid-precision approach successfully mitigates this degradation, aligning with its reported  
238 findings.
- 239 • Top-10% Attention, a sparse attention baseline, delivered surprisingly strong performance. Despite  
240 its high effective sparsity (90%) on shorter sequences, it ranked second only to SageAttention,  
241 highlighting the viability of sparse attention as a high-fidelity compression strategy.

#### 242 4.3 Scaling with Model Size

243 Theoretically, larger models exhibit greater parameter redundancy, suggesting they should be more  
244 amenable to compression. To test this hypothesis, we evaluated how the fidelity of each compression  
245 method scales with model size, using the 7B, 14B, and 32B versions of the Qwen2.5-Instruct model.  
246 The results are plotted in Figure [xx].

247 Most methods, including quantization and low-precision attention, become more faithful as model  
248 size increases, confirming the initial hypothesis. While AWQ, Top-10% Attention, and SparseGPT  
249 show a minor, anomalous dip at the 13B scale, the overall trend remains positive. SnapKV is the  
250 notable exception; its fidelity degrades as the model grows larger. This analysis not only provides a  
251 new dimension for comparing compression techniques but also empirically supports the notion that  
252 parameter efficiency decreases with model scale.

#### 253 4.4 Scaling with Context Length

254 Low-precision attention and KV cache compression are often motivated by their potential to improve  
255 efficiency in long-context scenarios. Our framework is naturally suited to evaluating performance on  
256 prompts of varying lengths. We stratified our test data from ShareGPT into three tiers based on token  
257 count (8K, 16K, and 24K) to assess how the fidelity of each method scales as the context window  
258 expands.

259 The results, shown in Figure 3, reveal three distinct scaling behaviors:

- 260 • All quantization methods (GPTQ, AWQ), KV cache dropping methods (H2O, SnapKV), SparseGPT,  
261 and SageAttention show a clear decline in fidelity as context length increases.
- 262 • FlashAttention FP8 and Wanda maintain relatively consistent performance across all context  
263 lengths.
- 264 • Top-10% Attention is unique in that its fidelity measurably improves with longer contexts.

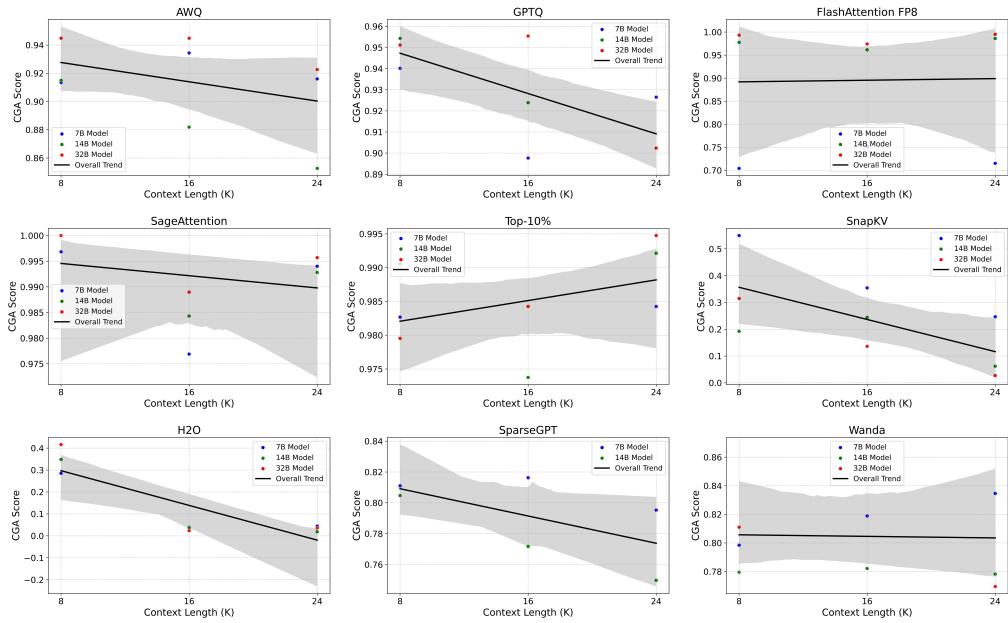


Figure 3: Scaling of compression method fidelity with increasing context length. Methods exhibit distinct behaviors: degradation (quantization, KV cache dropping), stability (Wanda, FlashAttention FP8), or improvement (Top-10% Attention).

265 These findings demonstrate that even methods within the same category possess different scaling  
 266 properties, a critical factor for selecting an appropriate compression strategy for long-context  
 267 applications.

## 268 5 Conclusion and Limitations

269 We have proposed and validated a new evaluation framework for LLM compression that uses real-  
 270 world user queries and a direct fidelity metric, Conditional Generation Accuracy (CGA), to measure  
 271 performance. Our experiments show that CGA aligns significantly better with human-aligned  
 272 judgments than proxy metrics like PPL and QA benchmark accuracy. Using this framework, we  
 273 conducted a comprehensive analysis of mainstream compression methods, revealing performance  
 274 trade-offs across model sizes and context lengths that provide a more realistic assessment.

275 Our approach has two primary limitations. First, its reliance on a token-by-token evaluation makes it  
 276 more computationally intensive than traditional benchmarks. Second, by averaging accuracy across  
 277 the entire sequence, CGA treats all tokens as equally important, which may not fully capture the  
 278 nuanced impact of compression on output quality. Investigating methods to weigh tokens by their  
 279 semantic or structural importance is a promising area for future work.

280 **References**

- 281 [1] J. Hong, J. Duan, C. Zhang, Z. Li, C. Xie, K. Lieberman, et al. Decoding compressed trust:  
282 scrutinizing the trustworthiness of efficient llms under compression. In *ICML*, 2024.
- 283 [2] A. Jaiswal, Z. Gan, X. Du, B. Zhang, Z. Wang, and Y. Yang. Compressing llms: The truth is  
284 rarely pure and never simple. In *ICLR*, 2024.
- 285 [3] Q. Wang, M. Wang, N. Feldhus, S. Ostermann, Y. Cao, H. Schütze, et al. Through a compressed  
286 lens: Investigating the impact of quantization on llm explainability and interpretability. *arXiv*  
287 preprint *arXiv:2505.13963*, 2025.
- 288 [4] Z. Xu, A. Gupta, T. Li, O. Bentham, and V. Srikumar. Beyond perplexity: Multi-dimensional  
289 safety evaluation of llm compression. In *EMNLP*, 2024.