

# Sentinel: Decoding Context Utilization via Attention Probing for Efficient LLM Context Compression

Anonymous ACL submission

## Abstract

Retrieval-augmented generation (RAG) often suffers from long and noisy retrieved contexts. Prior context compression methods rely on predefined importance metrics or supervised compression models, rather than on the model’s own inference-time behavior. We propose Sentinel, a lightweight sentence-level compression framework that treats context compression as an understanding decoding problem. Sentinel probes native attention behaviors of a frozen LLM with a lightweight readout to decode which parts of the context are actually utilized when answering a query, rather than using attention as a direct relevance score. We empirically observe that decoded relevance signals exhibit sufficient consistency across model scales to support effective compression with compact proxy models. On LongBench, Sentinel with a 0.5B proxy model achieves up to 5× compression while matching the QA performance of 7B-scale baselines, and despite being trained only on English QA data, generalizes effectively to Chinese and out-of-domain settings.

## 1 Introduction

Large language models (LLMs) have achieved impressive performance across open-domain question answering, reasoning, and dialogue tasks (Brown et al., 2020; OpenAI, 2024). To scale their capabilities to knowledge-intensive applications, Retrieval-Augmented Generation (RAG) has emerged as a powerful paradigm that augments model inputs with retrieved evidence from external corpora (Lewis et al., 2020; Guu et al., 2020; Shi et al., 2024). However, long retrieved contexts are often noisy, redundant, or exceed model input limits, making context compression essential for both efficiency and effectiveness (Liu et al., 2024; Yoran et al., 2024).

Existing context compression methods can be broadly divided into two categories. Metric-based approaches estimate the utility of context using

predefined or model-derived importance metrics, such as perplexity, self-information, mutual information, or query–context similarity (Jiang et al., 2023, 2024a; Li et al., 2023). While lightweight and training-free, these methods estimate relevance via heuristic or proxy importance scores, which are only indirectly related to the model’s inference-time behavior. In contrast, data-driven approaches learn compression decisions using external supervision or generator feedback to optimize downstream task performance (Pan et al., 2024; Xu et al., 2024; Hwang et al., 2024). Although effective, these approaches treat context compression as an optimization problem external to the model’s inference process, introducing additional training cost and often tying compression behavior to specific training objectives or generator feedback.

Recent mechanistic studies of Transformer-based LLMs have shown that decoder-only models exhibit structured context-utilization behaviors, with specialized attention heads supporting query–context alignment and evidence retrieval (Wu et al., 2024; Jin et al., 2024; Huang et al., 2025). These findings suggest that LLMs actively form query-aware contextual understanding during inference, rather than passively consuming retrieved context. Despite these insights, existing compression methods rarely exploit model-internal understanding signals directly. A key challenge is that such signals are not readily accessible in a stable and lightweight manner in decoder-only LLMs. Moreover, naively using raw attention as a compression signal often proves unreliable, as attention patterns mix informative behaviors with various forms of noise.

Motivated by these observations, we reformulate context compression as an understanding decoding problem, where compression decisions are derived from how an LLM internally utilizes context when answering a query, rather than from external importance heuristics. We propose Sentinel, a

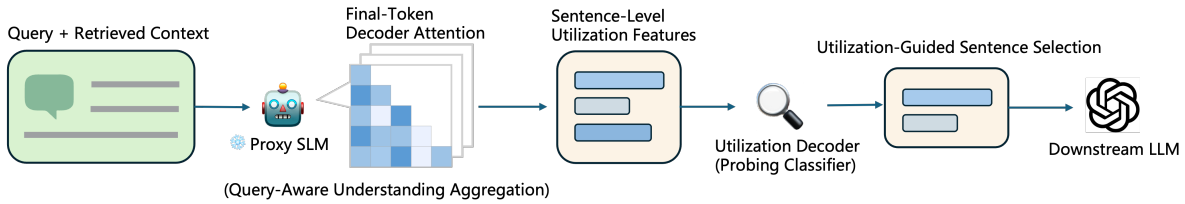


Figure 1: **Sentinel Framework Overview.** Sentinel decodes query-aware context utilization from native attention behaviors of a frozen LLM. By probing sentence-level attention features aggregated at a single decoding step, Sentinel identifies relevant context without training compression models or performing full autoregressive generation.

lightweight context compression framework that decodes query-aware context utilization from native attention behaviors of LLMs. Sentinel requires neither training a dedicated compression model nor running full autoregressive generation. Instead, it probes attention behaviors of a frozen proxy LLM under carefully designed prompts, enabling efficient sentence-level compression in a single forward pass. By grounding compression decisions in the model’s own inference-time behavior, Sentinel provides an interpretable and scalable alternative that leverages the LLM’s native understanding capabilities, rather than learning task-specific compression policies.

Empirically, Sentinel achieves up to 5× input compression on LongBench while attaining question-answering performance comparable to 7B-scale compression baselines, using only a 0.5B proxy model. The probing classifier relies solely on existing English QA data, yet the resulting compression strategy generalizes effectively to out-of-domain English LongBench tasks and exhibits robust cross-lingual generalization on Chinese benchmarks. Across multiple model families and scales, including Qwen-2.5, Qwen3, and LLaMA-3 variants, Sentinel shows consistent compression behavior under a unified and lightweight probing paradigm.

**Our contributions are as follows:**

- We introduce a new formulation of context compression that formulates it as an *understanding decoding problem*, grounding compression decisions in how LLMs internally utilize context when answering a query, rather than in external importance metrics or generation-based objectives.
- We propose **Sentinel**, a lightweight context compression framework that decodes query-aware context utilization from the native attention behaviors of frozen LLMs via a probing-based paradigm, eliminating the need for ded-

icated compression models or full autoregressive generation.

- We empirically observe that the query–context relevance signals decoded by Sentinel remain consistent across proxy model scales and families, enabling compact proxy models to approximate the context compression behavior of much larger LLMs.
- We conduct extensive experiments on long-context benchmarks, demonstrating that Sentinel achieves substantial compression while improving downstream QA performance, with mechanistic analyses supporting consistency with prior findings on retrieval-oriented behaviors in LLMs.

**2 Methodology**

**2.1 Context Compression as Understanding**

We propose Sentinel, a lightweight framework that approaches context compression from an understanding-centric perspective. Rather than training a compression model end-to-end or estimating sentence importance using external heuristics, Sentinel treats context compression as a problem of decoding how a language model internally understands and utilizes retrieved context when answering a query.

Formally, given a query  $q$  and a retrieved context  $C = \{s_1, s_2, \dots, s_n\}$  composed of sentences, the goal of context compression is to select a subset  $C' \subseteq C$  that preserves the information actually used by the model to answer  $q$ . Under this formulation, compression decisions are grounded in the model’s internal inference behavior, rather than in generation outputs or predefined utility metrics.

**2.2 Aggregating Query-Aware Understanding in LLMs**

Recent studies suggest that decoder-only LLMs actively form query-aware contextual understanding during inference, which is reflected in their

164	internal attention behaviors (Wu et al., 2024; Jin	which maps each feature vector $\mathbf{v}_i$ to a scalar rele-	212
165	et al., 2024; Huang et al., 2025). However, directly	ance score via a sigmoid function:	213
166	exploiting such understanding signals for context		
167	compression is non-trivial, as context utilization is	$\hat{y}_i = \sigma(\mathbf{w}^\top \mathbf{v}_i + b),$	214
168	typically distributed across multiple decoding steps	where $\mathbf{w}$ and $b$ denote the probe parameters.	215
169	and raw attention patterns are highly noisy.	The resulting probability is interpreted as the de-	216
170	An encouraging observation is that, under ap-	gree to which sentence $s_i$ is utilized by the model	217
171	propriate prompting, global contextual understand-	when answering the query. We choose a linear	218
172	ing can be implicitly aggregated at specific stages	probe for two reasons. First, it minimizes the risk	219
173	of the inference process, with the first genera-	of learning new behaviors beyond those already	220
174	tion step encoding a compressed representation	encoded in the model. Second, it enables direct	221
175	of the entire input (Jiang et al., 2024c). From an	interpretability of individual attention heads, allow-	222
176	information-theoretic perspective, this aggregation	ing us to analyze which attention head contribute	223
177	can be viewed as over-squashing, where prior con-	positively or negatively to context utilization.	224
178	text is progressively compressed into the final token		
179	representation (Barbero et al., 2024).	<b>2.4 Weak Supervision for Probing Context</b>	225
180	Following this insight, we feed the query and re-	<b>Utilization</b>	226
181	trieved context into a compact decoder-only proxy	To decode the model’s internal context utilization	227
182	model with instruction-following capability, apply	behavior, we train the probing classifier using weak	228
183	a QA-style prompt that encourages semantic com-	supervision derived from question answering data.	229
184	pression at the final position, and extract decoder	Importantly, this supervision is not intended to an-	230
185	attention from the final decoding token as a com-	notate sentence importance in general, but to iden-	231
186	compact carrier of query-aware understanding signals.	tify which sentences the model treats as evidence	232
187	<b>2.3 Decoding Context Utilization via Probing</b>	when it genuinely relies on retrieved context to	233
188	We decode context utilization by probing sentence-	answer a query.	234
189	level attention features extracted from the proxy	<b>2.4.1 Probing Data Construction</b>	235
190	model, without modifying or fine-tuning the model.	We construct probing examples from existing QA	236
191	<b>2.3.1 Sentence-Level Attention Features</b>	datasets that provide answer span annotations	237
192	For each query–context input, we extract the de-	within retrieved contexts, covering both single-	238
193	coder attention tensor from the final decoding to-	hop and multi-hop question answering settings.	239
194	ken, capturing attention scores across layers, heads,	For each QA instance, sentences containing the	240
195	and input tokens. To obtain sentence-level represen-	gold answer span are labeled as positive, while all	241
196	tations, we aggregate the attention weights directed	other sentences in the retrieved context are labeled	242
197	toward the tokens of each sentence, and normalize	as negative. This weak supervision allows us to	243
198	them by the total attention mass over the context	scale probing without manual relevance annotation	244
199	span. This normalization removes the influence of	and exposes the probe to diverse reasoning pat-	245
200	prompt and query tokens and yields comparable	terns, ranging from localized factual evidence to	246
201	relevance signals across sentences.	distributed multi-hop evidence.	247
202	Averaging the normalized attention weights over	<b>2.4.2 Selecting Context-Reliant Samples</b>	248
203	tokens within each sentence produces a feature vec-	To purify supervision, we retain only QA examples	249
204	tor for sentence $s_i$ , denoted as $\mathbf{v}_i \in R^{LH}$ , where	that require retrieved context for correct answering.	250
205	each dimension corresponds to a specific attention	Specifically, we keep examples where the model	251
206	head at a particular layer.	fails without access to the retrieved context but	252
207	<b>2.3.2 Probing Context Utilization</b>	succeeds when the context is provided. This filter-	253
208	To decode sentence relevance from attention fea-	ing echoes prior work that probes model behavior	254
209	tures, we train a lightweight probing classifier on	via intervention-based output changes (Meng et al.,	255
210	top of the sentence-level representations. Specifi-	2022). It ensures that positive sentences genuinely	256
211	cally, we adopt logistic regression as a linear probe,	contribute essential information for answering, re-	257
		ducing contamination from internal memorization	258
		or hallucinated knowledge. By filtering for retrieval	259

260	dependency, we focus training on cases where relevance must be decoded from the provided context, aligning with the goal of decoding context utilization rather than internal recall.	308
261		309
262		310
263		311
264	<b>2.4.3 Robustness via Sentence Shuffling</b>	312
265	To mitigate positional biases (Liu et al., 2024), especially common in multi-document retrieval settings, we apply sentence shuffling during training by randomly permuting sentence order within each passage. This simple perturbation encourages the classifier to rely on semantic relevance rather than fixed positions, improving generalization to real-world RAG inputs with noisy or varied structure.	313
266		314
267		315
268		316
269		317
270		318
271		319
272		320
273	<b>2.5 Inference-Time Context Compression</b>	321
274	At inference time, given a query–context pair $(q, C)$ , Sentinel runs a single forward pass of a compact proxy model with a fixed QA-style prompt, extracts final-token decoder attention, and computes sentence-level attention features. A trained probing classifier assigns relevance scores to sentences, based on which a top-ranked subset $C' \subseteq C$ is selected under a length budget and passed to the downstream LLM for answer generation.	322
275		323
276		324
277		325
278		326
279		327
280		328
281		329
282		330
283	<b>3 Experiments</b>	331
284	<b>Datasets</b> We evaluate Sentinel on both the English and Chinese subsets of LongBench (Bai et al., 2024). Following our focus on query-conditioned context compression, we report results on question answering tasks and query-conditioned summarization (e.g., QMSum), which involve an explicit query. When using Qwen-2.5-7B as the downstream LLM, we exclude Few-shot, Code, and Synthetic tasks, as they lack an explicit query or rely on strict prompt structure. For fair comparison with prior work, we additionally report results on the full LongBench benchmark when comparing against baseline methods. Detailed dataset descriptions are provided in the appendix A.	332
285		333
286		334
287		335
288		336
289		337
290		338
291		339
292		340
293		341
294		342
295		343
296		344
297		345
298	<b>Probing Data</b> We train the probing classifier on a small set of English QA examples spanning both single-hop and multi-hop reasoning. In the default setting, we sample 3K QA instances, each yielding one positive sentence containing the gold answer span and one negative sentence from the same context, resulting in 6K sentence-level training examples. We retain only context-reliant examples, where correct answering requires access to the retrieved context. Sentence-level attention features	346
299		347
300		348
301		349
302		350
303		351
304		352
305		353
306		354
307		355
	are extracted using a fixed QA-style prompt by collecting decoder attention from the final decoding token. Additional implementation details are provided in Appendix B.	
	<b>Probing Classifier Training</b> We train a logistic regression probe on attention-derived features using standard cross-validation and regularization. Additional training details are in Appendix B.	
	<b>Compression Strategy</b> Sentinel performs length-controlled context compression by ranking sentences with the probing classifier and selecting a top-ranked subset under a specified budget. Selected sentences are concatenated in their original order and passed to the downstream LLM.	
	We consider two budget settings, both measured using the target model’s tokenizer: (i) a fixed token budget $B$ (e.g., 2000 tokens), where sentences are selected until the budget is reached; and (ii) a compression ratio $\tau \in [0.1, 0.5]$ , where the retained sentences do not exceed a fraction of the original context length.	
	<b>Proxy Model Setup</b> Unless otherwise specified, Sentinel uses Qwen-2.5-0.5B-Instruct as the default proxy model for attention-based feature extraction and sentence relevance probing, with a chunk size of 1024 tokens.	
	<b>Evaluation Models</b> Following LLMingua setup, we use ChatGPT (gpt-3.5-turbo) as the primary model for evaluation. To assess the generality of our method, we also experiment with Qwen-2.5-7B-Instruct in our main results. All evaluations follow the LongBench prompt and decoding setup (Bai et al., 2024), as detailed in Appendix I.	
	<b>Baselines</b> We compare Sentinel against representative context compression baselines spanning both metric-based and data-driven approaches. Metric-based baselines include LLMingua-1 (Jiang et al., 2023), LongLLMingua (Jiang et al., 2024b), and Selective Context (Li et al., 2023). Data-driven baselines include LLMingua-2 (Pan et al., 2024) and CPC (Liskavets et al., 2024). We include an attention-based heuristic baseline, denoted as <b>Raw Attention</b> , which represents a class of methods that directly use aggregated decoder attention weights as relevance scores, such as QUITO (Wang et al., 2024) and AttentionRAG (Fang et al., 2025). We also include non-learning baselines including Ran-	

Methods	LongBench-En (GPT-3.5-Turbo, 2000-token constraint)							Compression Stats	
	SingleDoc	MultiDoc	Summ.	FewShot	Synth.	Code	AVG	Tokens	1/τ
Selective-Context (LLaMA-2-7B-Chat)	16.2	34.8	24.4	15.7	8.4	49.2	24.8	1,925	5x
LLMLingua (LLaMA-2-7B-Chat)	22.4	32.1	24.5	61.2	10.4	56.8	34.6	1,950	5x
LLMLingua-2 (XLM-RoBERTa-Large-0.6B)	29.8	33.1	25.3	66.4	21.3	<u>58.9</u>	39.1	1,954	5x
LongLLMLingua (LLaMA-2-7B-Chat)	39.0	42.2	<b>27.4</b>	69.3	<b>53.8</b>	56.6	48.0	1,809	6x
CPC (Mistral-7B-Instruct-v0.2)	<b>42.6</b>	<b>48.6</b>	23.7	69.4	<u>52.8</u>	<b>60.0</b>	<b>49.5</b>	1,844	5x
Sentinel (Qwen-2.5-0.5B-Instruct)	40.1	47.4	25.8	<b>69.9</b>	46.3	58.0	47.89	1,885	5x
Sentinel (Qwen-2.5-1.5B-Instruct)	<u>40.6</u>	<u>48.1</u>	<u>26.0</u>	69.1	49.0	57.6	<u>48.4</u>	1,883	5x
Original Prompt	39.7	38.7	26.5	67.0	37.8	54.2	44.0	10,295	-

Table 1: Performance on English LongBench using GPT-3.5-Turbo as the inference model. Best results are in **bold**, second-best are underlined.

Methods	LongBench-En (2000-token constraint)				LongBench-Zh (2000-token constraint)			Overall AVG
	SingleDoc	MultiDoc	Summ.	En-AVG	SingleDoc	MultiDoc	Zh-AVG	
Empty	10.72	22.26	16.46	16.48	17.71	13.54	15.62	16.05
Random	28.22	30.68	20.33	26.41	43.18	17.22	30.20	28.30
Raw Attention (Qwen-2.5-0.5B-Instruct)	34.92	38.96	21.32	31.74	51.72	17.29	34.50	33.12
Sentinel (Qwen-2.5-0.5B-Instruct)	37.73	<b>46.16</b>	<b>23.03</b>	<b>35.64</b>	<b>62.24</b>	<b>18.57</b>	<b>40.41</b>	<b>38.02</b>
Original Prompt	<b>38.84</b>	44.74	22.76	35.45	60.06	18.21	39.14	37.30

Table 2: LongBench results under a 2000-token context constraint, evaluated using Qwen-2.5-Instruct-7B as the downstream LLM. Results are reported on both English and Chinese subsets. The **Summ.** column corresponds to query-conditioned summarization tasks (QMSum).

dom Selection and Empty Context. Full descriptions are provided in Appendix C.

**Metrics** We follow the LongBench evaluation protocol and adopt task-specific metrics for each task category: QA-F1 for Single-Document QA, Multi-Document QA, ROUGE-L for Summarization. All metrics are computed using the official evaluation scripts.

### 3.1 Results on LongBench

We evaluate Sentinel under two settings: (1) English LongBench tasks using GPT-3.5-Turbo as the inference model, and (2) both English and Chinese LongBench tasks using Qwen-2.5-7B-Instruct. All results are reported under a 2,000-token input constraint.

**Strong Performance with Compact Proxies under GPT-3.5-Turbo** Table 1 presents results on the English subset of LongBench using GPT-3.5-Turbo as the inference model. Using a compact 0.5B proxy model, Sentinel achieves strong performance across all evaluated task categories.

Despite its small proxy size, Sentinel consistently outperforms task-agnostic compression methods such as LLMLingua and LLMLingua-2, while using only a **0.5B proxy to perform competitively with 7B-scale compression baselines**, including CPC and LongLLMLingua.

These results demonstrate that effective query-aware sentence selection can be achieved using

compact proxy models, without training dedicated compression networks or relying on large-scale generators. Additional results on Chinese tasks under GPT-3.5-Turbo are reported in Appendix D.

**Effective Compression and Cross-Lingual Generalization under Qwen-2.5-7B-Instruct** Table 2 reports results on both English and Chinese LongBench subsets using Qwen-2.5-7B-Instruct as the downstream LLM. Across all evaluated tasks, Sentinel substantially outperforms Random, Empty Context, and Raw Attention baselines, confirming the effectiveness of decoding query-aware context utilization signals.

On the English subset, Sentinel surpasses the Original Prompt baseline on average and improves performance on most task categories under a strict 2,000-token budget. Notably, despite being trained exclusively on English QA data, Sentinel maintains strong performance on the Chinese subset and achieves clear improvements over Raw Attention, demonstrating **robust cross-lingual generalization** of the decoded relevance signals.

Overall, these results show that Sentinel enables effective query-conditioned context compression by decoding model-internal understanding signals, and that the resulting relevance estimates generalize reliably across languages and inference settings.

### 3.2 Effect of Proxy Model Family and Size

We evaluate Sentinel using proxy models from three model families, including Qwen-2.5 (Qwen

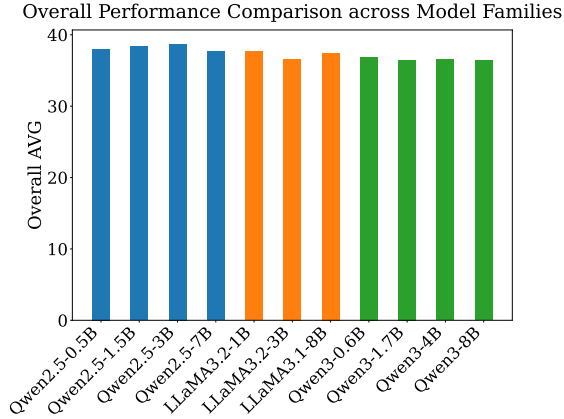


Figure 2: Impact of proxy model family and scale on Sentinel performance under a 2k-token context (LongBench Overall AVG)

et al., 2025), Qwen-3 (Yang et al., 2025), and LLaMA-3 (Grattafiori et al., 2024), spanning parameter scales from 0.5B to 8B. As shown in Figure 2, Sentinel achieves comparable overall performance across proxy models of different sizes and families, indicating that the query-context relevance signals decoded by Sentinel are consistent across model scales. Notably, increasing the proxy model size does not lead to systematic improvements in compression performance. This suggests that the relevance signals exploited by Sentinel are already sufficiently expressed in compact proxy models, rather than emerging only with increased model capacity. As a result, Sentinel can be deployed efficiently using small and computationally inexpensive proxy models without sacrificing compression quality. Detailed per-model and per-task results are provided in Appendix E.

### 3.3 Ablation

We conduct ablation studies to analyze where Sentinel’s performance gains come from and to verify that they arise from decoding model-internal understanding signals rather than from probe capacity, large-scale supervision, or specific inference heuristics. By default, Sentinel is instantiated on Qwen-2.5-0.5B-Instruct. Unless otherwise specified, all experiments use Qwen-2.5-7B-Instruct as the downstream LLM and are evaluated on LongBench.

#### 3.3.1 Effect of Probing Data Size

A key design principle of Sentinel is to decode relevance signals already encoded in frozen LLMs, rather than to learn a new compression behavior from supervision. To evaluate the dependence of Sentinel on probing data size, we vary the number

Probing Size	Overall AVG
500	38.03
1000	38.24
2000	37.92
3000	38.02

Table 3: Overall performance under different probing data sizes.

Feature	HotpotQA	SQuAD	NewsQA	Overall AUC	Overall AVG
All Layers	0.9228	0.9987	0.9838	0.9700	38.02
Selected	0.9171	0.9943	0.9832	0.9662	36.56
Last Layer	0.8606	0.9538	0.9588	0.9121	37.24

Table 4: AUC comparison of different attention feature extraction strategies.

of QA examples used for probing from 500 to 3000.

As shown in Table 3, downstream performance remains nearly invariant across this range, with no clear trend of improvement as more probing data is added. This indicates that Sentinel does not rely on large amounts of supervision: the attention-based relevance signals are already present in the model, and the probing classifier mainly acts as a lightweight readout of these signals. Notably, even a small probing set is sufficient to support effective context compression. Detailed task-level results are provided in Appendix F.

#### 3.3.2 Attention Feature Ablations

We evaluate three attention-based feature construction strategies using Qwen-2.5-0.5B-Instruct as the proxy model: (i) aggregating attention from all decoder layers, (ii) using only the final decoder layer, and (iii) selecting a compact subset of attention heads via mRMR (Ding and Peng, 2005), constrained to at most one layer’s worth of heads (see Appendix F for details).

As shown in Table 4, aggregating attention across all layers consistently achieves the strongest AUC and downstream performance. Using only the final decoder layer leads to a noticeable performance drop, while the selected-head variant preserves most of the performance with substantially fewer features. These results suggest that query-aware relevance signals are distributed across multiple layers rather than being dominated by the final decoder layer, supporting our choice of aggregating attention information across the full model.

#### 3.3.3 Compression Ratio Variants

We further evaluate robustness under different compression ratios  $\tau \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ , where smaller  $\tau$  corresponds to more aggressive pruning. As shown in Figure 3, the Raw Attention baseline degrades sharply with increasing compression,

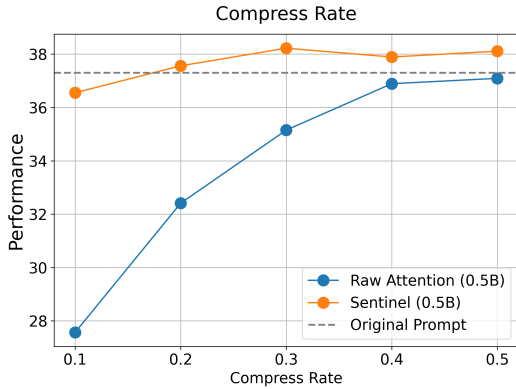


Figure 3: Compression ratio ablation on Qwen-2.5-7B-Instruct with a 0.5B proxy.

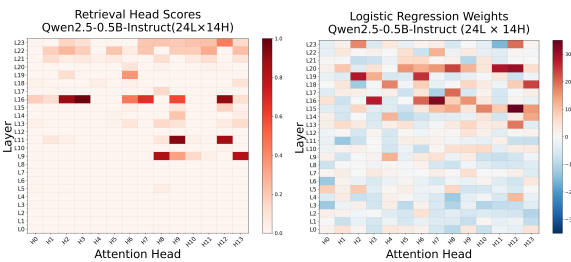


Figure 4: Comparison of attention head importance patterns identified by retrieval-head analysis (left) and Sentinel probing (right).

collapsing when  $\tau < 0.4$ . In contrast, Sentinel achieves strong performance across all compression levels, consistently outperforming the Raw Attention baseline and surpassing the original prompt performance over a wide range of  $\tau$ , including the extreme setting of  $\tau = 0.2$ . This result indicates that decoding query-aware context utilization signals yields more robust context compression than relying on raw attention magnitudes under aggressive pruning, enabling Sentinel to retain the most critical contextual information even under severe compression. Full task-level results are reported in Appendix F.

#### 4 Analysis: Interpreting Context Utilization Decoded from Attention Heads

**Alignment with Retrieval-Oriented Attention Heads** Prior mechanistic analyses have shown that retrieval-oriented and context-utilization behaviors in decoder-only LLMs are concentrated in a sparse subset of attention heads (Wu et al., 2024; Jin et al., 2024), whose activation depends on the input tokens and contexts.

Using a fundamentally different identification mechanism, Sentinel recovers a meaningful subset of these retrieval-oriented structures. Specifically,

when comparing the top-14 positively weighted heads identified by Sentinel with the top-14 retrieval heads obtained by reproducing the retrieval-head identification procedure of prior work (Wu et al., 2024) on the same Qwen-2.5-0.5B-Instruct model, we observe five overlapping heads. As illustrated in Figure 4, these overlapping heads are not randomly distributed, but are predominantly located in middle-to-late layers, with a noticeable concentration around layer 16.

**Efficient and Ensemble-Based Decoding of Context Utilization** This observed overlap and layer-wise concentration indicate that the proposed probing approach captures core context-utilization behaviors identified by prior mechanistic studies. Whereas prior work identifies retrieval heads individually through autoregressive decoding and token-level copy analysis, Sentinel assigns weights to all attention heads and decodes context utilization through their weighted aggregation. This ensemble-style decoding provides a more efficient and compression-oriented access to model-internal understanding signals.

**Beyond Positive Retrieval Heads** More importantly, Sentinel goes beyond identifying only positively contributing heads. By assigning signed weights through probing, Sentinel captures both heads that support evidence utilization and heads that systematically interfere with it. We observe that heads assigned negative weights are frequently associated with structurally dominant but semantically uninformative attention patterns, such as attention sinks (Bondarenko et al., 2021; Son et al., 2024). Further analysis of negatively weighted heads is provided in Appendix G.

**Robustness under Dynamic and Multi-Functional Head Behavior** This distinction is particularly important given that retrieval heads are dynamically activated depending on input tokens and contexts (Wu et al., 2024), reflecting the multi-functional nature of attention heads. A single head may support evidence retrieval in some contexts while exhibiting non-retrieval behaviors in others (Zheng et al., 2024). By aggregating both positive and negative contributions across all heads, Sentinel mitigates instability caused by context-dependent role switching and counterbalances spurious activations, resulting in a more robust decoding of context utilization than approaches that rely exclusively on positively

Methods	LongBench-En (2000-token constraint)				LongBench-Zh (2000-token constraint)			Overall AVG
	SingleDoc	MultiDoc	Summ.	En-AVG	SingleDoc	MultiDoc	Zh-AVG	
Top 14 Retrieval Heads (Qwen-2.5-0.5B-Instruct)	36.53	43.16	22.27	33.99	58.59	18.53	38.56	36.28
Sentinel (Qwen-2.5-0.5B-Instruct)	37.73	46.16	23.03	35.64	62.24	18.57	40.41	38.02

Table 5: Comparison between retrieval-head-based compression and Sentinel on the LongBench benchmark, where retrieval-head compression scores sentences using attention from the top-14 retrieval heads.

identified retrieval heads. Table 5 further supports this analysis: Sentinel consistently outperforms retrieval-head-based compression on LongBench across English, Chinese, and overall averages.

## 5 Related Work

**Metric-Based Context Compression** Metric-based approaches estimate context utility using predefined importance scores, such as self-information, mutual information, or query-context similarity, and select tokens or sentences accordingly without training a dedicated compression model. Representative token-level methods include LLMLingua (Jiang et al., 2023) and LongLLM-Lingua (Jiang et al., 2024a), which prune tokens based on perplexity or query-conditioned probability estimates. QUITO-X (Cao et al., 2024) further introduces mutual information-based scoring to guide context selection. At a coarser granularity, Selective Context (Li et al., 2023) removes low-information content based on token-level self-information scores.

Some recent methods (Wang et al., 2024; Fang et al., 2025) leverage decoder attention as an importance heuristic, using vanilla or heuristically aggregated attention weights as a proxy for relevance. However, attention is used in a post-hoc scoring manner, rather than being decoded to reflect how the model internally utilizes context during inference.

While effective and training-free, these methods rely on predefined or heuristically applied importance scores that are not explicitly tied to decoding-time context utilization.

**Data-Driven Context Compression** Data-driven approaches learn compression decisions from external supervision, typically by training a ranking or classification model to predict which tokens or sentences should be retained. Token-level methods such as LLMLingua-2 (Pan et al., 2024) leverage distilled labels from large language models to train lightweight compressors. At the sentence level, methods such as RECOMP (Xu et al., 2024) train compressors to produce

extractive or abstractive summaries that improve downstream performance, while EXIT (Hwang et al., 2024) learns a sentence-level classifier to select query-relevant sentences. Other works, such as CPC (Liskavets et al., 2024), Refiner (Li et al., 2024), and FineFilter (Zhang et al., 2025), further incorporate query-aware ranking, structure-aware reranking, or multi-hop reasoning objectives. Although these methods often achieve strong performance, they introduce additional training cost and data dependency, which can limit their adaptability across tasks and models.

**Our Approach: Utilization-Driven Context Compression** In contrast to prior work, we propose Sentinel, a lightweight and model-agnostic framework that adopts a utilization-driven perspective on context compression. Rather than relying on predefined metrics or externally supervised compression models, Sentinel decodes how a frozen LLM internally utilizes context when answering a query. Crucially, Sentinel does not treat attention weights as direct relevance scores. While raw attention answers the question of where the model attends, Sentinel probes attention behaviors with a lightweight readout to decode which parts of the context are actually utilized by the model during inference.

## 6 Conclusion

We present **Sentinel**, a lightweight context compression framework that decodes how LLMs utilize context when answering a query. By probing native attention behaviors of a frozen proxy model, Sentinel enables effective sentence-level compression without training a dedicated compressor or relying on full autoregressive generation. Empirically, Sentinel achieves up to  $5\times$  compression on LongBench while matching or improving QA performance compared to strong 7B-scale baselines, using only a 0.5B proxy. These results suggest that model-internal utilization signals provide a principled and efficient foundation for query-aware context compression.

## 649 Limitations

650 **Query-Conditioned Scope.** Sentinel is designed  
651 for query-conditioned context compression, where  
652 relevance is defined with respect to an explicit  
653 question or instruction. Tasks that lack a clear  
654 query signal or rely on strict prompt structure, such  
655 as free-form summarization or code completion,  
656 fall outside the current scope of our framework.  
657 An interesting direction for future work is to ex-  
658 plore whether such tasks can be reformulated with  
659 auxiliary or synthetic queries, enabling utilization-  
660 driven compression beyond explicit QA settings.

661 **Evaluation Backbones.** Our downstream evalua-  
662 tion focuses on two decoder LLMs, GPT-3.5-Turbo  
663 and Qwen-2.5-7B-Instruct, covering both propri-  
664 etary and open-source settings. While Sentinel  
665 demonstrates consistent behavior across multiple  
666 proxy model families and scales, further evalua-  
667 tion with additional inference backbones could help  
668 characterize its behavior under different instruction  
669 styles and decoding configurations.

## 670 References

671 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu,  
672 Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao  
673 Liu, Aohan Zeng, Lei Hou, and 1 others. 2024. Long-  
674 bench: A bilingual, multitask benchmark for long  
675 context understanding. In *Proceedings of the 62nd  
676 Annual Meeting of the Association for Computational  
677 Linguistics (Volume 1: Long Papers)*, pages 3119–  
678 3137.

679 Federico Barbero, Andrea Banino, Steven Kapturowski,  
680 Dharshan Kumaran, João Madeira Araújo, Oleksandr  
681 Vitvitskyi, Razvan Pascanu, and Petar Veličković.  
682 2024. Transformers need glasses! information over-  
683 squashing in language tasks. *Advances in Neural  
684 Information Processing Systems*, 37:98111–98142.

685 Yelysei Bondarenko, Markus Nagel, and Tijmen  
686 Blankevoort. 2021. Understanding and overcoming  
687 the challenges of efficient transformer quantization.  
688 In *Proceedings of the 2021 Conference on Empirical  
689 Methods in Natural Language Processing*, pages  
690 7947–7969. Association for Computational Linguis-  
691 tics.

692 Tom Brown, Benjamin Mann, Nick Ryder, Melanie  
693 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind  
694 Neelakantan, Pranav Shyam, Girish Sastry, Amanda  
695 Askell, and 1 others. 2020. Language models are  
696 few-shot learners. *Advances in neural information  
697 processing systems*, 33:1877–1901.

698 Zhiwei Cao, Qian Cao, Yu Lu, Ningxin Peng, Luyang  
699 Huang, Shanbo Cheng, and Jinsong Su. 2024. Retain-  
700 ing key information under high compression ratios:

Query-guided compressor for LLMs. In *Proceedings  
of the 62nd Annual Meeting of the Association for  
Computational Linguistics (Volume 1: Long Papers)*,  
pages 12685–12695. 701  
702  
703  
704

Chris Ding and Hanchuan Peng. 2005. Minimum re-  
dundancy feature selection from microarray gene ex-  
pression data. *Journal of bioinformatics and compu-  
tational biology*, 3(02):185–205. 705  
706  
707  
708

Yixiong Fang, Tianran Sun, Yuling Shi, and Xiaodong  
Gu. 2025. Attentionrag: Attention-guided context  
pruning in retrieval-augmented generation. *arXiv  
preprint arXiv:2503.10720*. 709  
710  
711  
712

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-  
ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh  
Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-  
tra, Archie Sravankumar, Artem Korenev, Arthur  
Hinsvark, and 542 others. 2024. *The llama 3 herd of  
models*. *Preprint*, arXiv:2407.21783. 713  
714  
715  
716  
717  
718  
719  
720

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-  
pat, and Mingwei Chang. 2020. Retrieval augmented  
language model pre-training. In *International confer-  
ence on machine learning*, pages 3929–3938. PMLR. 721  
722  
723  
724

Yanwen Huang, Yong Zhang, Ning Cheng, Zhitao Li,  
Shaojun Wang, and Jing Xiao. 2025. Dynamic  
attention-guided context decoding for mitigating con-  
text faithfulness hallucinations in large language mod-  
els. *arXiv preprint arXiv:2501.01059*. 725  
726  
727  
728  
729

Taeho Hwang, Sukmin Cho, Soyeong Jeong, Hoyun  
Song, SeungYoon Han, and Jong C Park. 2024. Exit:  
Context-aware extractive compression for enhanc-  
ing retrieval-augmented generation. *arXiv preprint  
arXiv:2412.12559*. 730  
731  
732  
733  
734

Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing  
Yang, and Lili Qiu. 2023. LLMLingua: Compressing  
prompts for accelerated inference of large language  
models. In *Proceedings of the 2023 Conference on  
Empirical Methods in Natural Language Processing*,  
pages 13358–13376. 735  
736  
737  
738  
739  
740

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dong-  
sheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu.  
2024a. LongLLMLingua: Accelerating and enhanc-  
ing LLMs in long context scenarios via prompt com-  
pression. In *Proceedings of the 62nd Annual Meeting  
of the Association for Computational Linguistics (Vol-  
ume 1: Long Papers)*, pages 1658–1677. 741  
742  
743  
744  
745  
746  
747

Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dong-  
sheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu.  
2024b. LongLLMLingua: Accelerating and enhanc-  
ing LLMs in long context scenarios via prompt com-  
pression. In *Proceedings of the 62nd Annual Meeting  
of the Association for Computational Linguistics (Vol-  
ume 1: Long Papers)*, pages 1658–1677. 748  
749  
750  
751  
752  
753  
754

Ting Jiang, Shaohan Huang, Zhongzhi Luan, Deqing  
Wang, and Fuzhen Zhuang. 2024c. Scaling sentence 755  
756

757	embeddings with large language models. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 3182–3196.		
758		Yang, Jiayi Yang, Jingren Zhou, and 25 others. 2025. <a href="#">Qwen2.5 technical report</a> . <i>Preprint</i> , arXiv:2412.15115.	812
759			813
			814
760	Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen, Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and Jun Zhao. 2024. Cutting off the head ends the conflict: A mechanism for interpreting and mitigating knowledge conflicts in language models. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pages 1193–1215.	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-joon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. REPLUG: Retrieval-augmented black-box language models. 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.	815
761			816
762			817
763			818
764			819
765			820
766			821
767	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. <i>Advances in neural information processing systems</i> , 33:9459–9474.	Seungwoo Son, Wonpyo Park, Woohyun Han, Kyuyeon Kim, and Jaeho Lee. 2024. Prefixing attention sinks can mitigate activation outliers for large language model quantization. <i>arXiv preprint arXiv:2406.12016</i> .	823
768			824
769			825
770			826
771			827
772			
773			
774	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. 2023. Compressing context to enhance inference efficiency of large language models. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 6342–6353.	Wenshan Wang, Yihang Wang, Yixing Fan, Huaming Liao, and Jiafeng Guo. 2024. Quito: Accelerating long-context reasoning through query-guided context compression. <i>arXiv preprint arXiv:2408.00274</i> .	828
775			829
776			830
777			831
778			
779	Zhonghao Li, Xuming Hu, Aiwei Liu, Kening Zheng, Sirui Huang, and Hui Xiong. 2024. <i>Refiner</i> : Restructure retrieved content efficiently to advance question-answering capabilities. In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 8548–8572.	Wenhao Wu, Yizhong Wang, Guangxuan Xiao, Hao Peng, and Yao Fu. 2024. Retrieval head mechanistically explains long-context factuality. <i>arXiv preprint arXiv:2404.15574</i> .	832
780			833
781			834
782			835
783			
784			
785	Barys Liskavets, Maxim Ushakov, Shuvendu Roy, Mark Klibanov, Ali Etemad, and Shane Luke. 2024. Prompt compression with context-aware sentence encoding for fast and improved llm inference. <i>arXiv preprint arXiv:2409.01227</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> .	836
786			837
787			838
788			839
789			840
790	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. <a href="#">Qwen3 technical report</a> . <i>Preprint</i> , arXiv:2505.09388.	841
791			842
792			843
793			844
794			845
			846
			847
795	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. Locating and editing factual associations in gpt. <i>Advances in Neural Information Processing Systems</i> , 35:17359–17372.	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2024. Making retrieval-augmented language models robust to irrelevant context. In <i>The Twelfth International Conference on Learning Representations</i> .	848
796			849
797			850
798			851
799	OpenAI. 2024. <a href="#">Gpt-4 technical report</a> . <i>Preprint</i> , arXiv:2303.08774.	Qianchi Zhang, Hainan Zhang, Liang Pang, Hongwei Zheng, Yongxin Tong, and Zhiming Zheng. 2025. Finefilter: A fine-grained noise filtering mechanism for retrieval-augmented large language models. <i>arXiv preprint arXiv:2502.11811</i> .	852
800			853
801	Zhuoshi Pan, Qianhui Wu, Huiqiang Jiang, Menglin Xia, Xufang Luo, Jue Zhang, Qingwei Lin, Victor Rühle, Yuqing Yang, Chin-Yew Lin, H. Vicky Zhao, Lili Qiu, and Dongmei Zhang. 2024. LLMingua-2: Data distillation for efficient and faithful task-agnostic prompt compression. In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 963–981.	Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Mingchuan Yang, Bo Tang, Feiyu Xiong, and Zhiyu Li. 2024. Attention heads of large language models: A survey. <i>arXiv preprint arXiv:2409.03752</i> .	854
802			855
803			856
804			857
805			
806			
807			
808	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	<b>A Dataset Details</b>	862
809		We provide a detailed description of the datasets used in our experiments, based on the English subset of LongBench (Bai et al., 2024). LongBench is	863
810			864
811			865

866	a long-context benchmark covering diverse tasks	– LCC: Predict the next line of code given	909
867	designed to evaluate the capabilities of language	a code block, without an explicit natural	910
868	models in understanding and reasoning over ex-	language query.	911
869	tended textual inputs. It consists of six task cat-	– REPOBENCH-P: Predict the next line of	912
870	egories, each comprising multiple representative	a function given multi-file code context	913
871	datasets:	and the function signature.	914
872	• <b>Single-Document QA:</b>	<b>Excluded Task Categories</b> We exclude Long-	915
873	– NARRATIVEQA: Answer questions	Bench task categories that violate the assumptions	916
874	based on a single narrative document,	of query-conditioned context compression, includ-	917
875	such as a story or movie script.	ing Code, Synthetic, Few-shot, and Summarization	918
876	– QASPER: Answer questions grounded in	tasks. These tasks either lack explicit queries, de-	919
877	a scientific paper.	pend on strict prompt structure or surface-level	920
878	– MULTIFIELDQA-EN: Answer factual	consistency, or rely on evaluation metrics insensi-	921
879	questions from a long structured ency-	tive to compression quality. We therefore focus on	922
880	clopedic entry.	question answering tasks, where query-conditioned	923
881	• <b>Multi-Document QA:</b>	relevance estimation is well defined.	924
882	– HOTPOTQA, 2WIKIMULTIHOPQA,	<b>B Additional Probing Data and Training</b>	925
883	and MUSIQUE: Multi-hop QA tasks	<b>Details</b>	926
884	requiring reasoning across multiple	<b>Probing Data Composition</b> The probing classi-	927
885	passages to answer complex factoid	fier is trained on 3,000 QA examples sampled from	928
886	questions.	three widely used QA datasets: NewsQA (50%),	929
887	• <b>Summarization:</b>	SQuAD (20%), and HotpotQA (30%), covering	930
888	– GOVREPORT: Summarize long govern-	both single-hop and multi-hop reasoning scenarios.	931
889	ment reports.	Each QA example yields two sentence-level in-	932
890	– QMSUM: Query-based summarization	stances: one positive sentence containing the gold	933
891	of meeting transcripts.	answer span and one negative sentence sampled	934
892	– MULTINEWS: Summarize multi-source	from the same retrieved context, resulting in a total	935
893	news articles.	of 6,000 sentence-level training examples.	936
894	• <b>Few-shot Reasoning:</b>	<b>Context Length Statistics</b> We report the context	937
895	– TREC: Classify question types.	length distribution for completeness. In NewsQA,	938
896	– TRIVIAQA: Answer trivia-style factual	30.1% of examples contain 0–500 tokens and	939
897	questions.	69.9% contain 500–1,000 tokens when tokenized	940
898	– SAMSUM: Summarize short dialogues.	with the Qwen-2.5 tokenizer. In SQuAD, 99.3% of	941
899	– LSHT: Classify Chinese news headlines	examples fall within the 0–500 token range. For	942
900	into topic categories.	HotpotQA, all examples are restricted to 0–500 to-	943
901	• <b>Synthetic Retrieval:</b>	kens by limiting unrelated content in the retrieved	944
902	– PASSAGECOUNT: Count the number of	context.	945
903	unique paragraphs among potentially du-	<b>Sentence Segmentation</b> Retrieved contexts	946
904	plicated inputs.	are segmented into sentences using spaCy’s	947
905	– PASSAGERETRIEVAL-EN: Identify the	sentencizer. Sentence boundaries are used	948
906	source paragraph corresponding to a	consistently for both positive and negative sentence	949
907	given abstract.	extraction as well as for sentence-level attention	950
908	• <b>Code Completion:</b>	aggregation.	951
		<b>Prompt Template</b> Sentence-level attention fea-	952
		tures are extracted using a fixed QA-style prompt	953
		applied to each query–context pair. The prompt	954
		format is shown below:	955

Given the following information: {context}  
 Answer the following question based on the  
 given information with one or few words:  
 {question}  
 Answer:

**Attention Feature Extraction** For each prompted input, we collect decoder attention weights from the final decoding token across all layers and attention heads. Attention weights directed to tokens belonging to each sentence are aggregated and normalized to form fixed-length sentence-level feature vectors, which are then used as input to the probing classifier.

**Context-Reliant Sample Selection** To improve the quality of weak supervision, we retain only *context-reliant* QA examples, where access to the retrieved context is necessary for correct answering. Specifically, we compare model predictions with and without the retrieved context, using exact match (EM) or F1 as appropriate for each dataset.

For NewsQA and SQuAD, we retain examples where the model fails to answer correctly without context (memory-based EM = 0) but succeeds when the context is provided (context-based EM = 1). For HotpotQA, we retain samples with memory-based F1  $\leq 0.2$  and context-based F1  $\geq 0.5$ , reflecting its multi-hop and partial-match evaluation setting.

**Probing Classifier Training** We train a logistic regression (LR) model on attention-derived features, using 5-fold cross-validation with balanced accuracy as the scoring metric. We perform grid search over regularization strengths  $C \in \{0.01, 0.1, 1.0, 10.0, 100.0\}$ , and use the lib-linear solver with  $\ell_2$  regularization, class-balanced weighting, and a maximum of 2,000 iterations. The best model is selected based on AUC on the validation set.

## C Baseline Descriptions

We compare Sentinel against the following baseline methods, grouped by their design paradigms:

- **LLMLingua-1/2** (Jiang et al., 2023; Pan et al., 2024): Token-level compression methods based on saliency estimation via perplexity and LLM distillation. These methods are task-agnostic and do not condition on the query.
- **Selective-Context** (Li et al., 2023): A sentence-level, task-agnostic method that

scores context segments based on general informativeness, independent of the question.

- **LongLLMLingua** (Jiang et al., 2024b): A query-aware, multi-stage compression system using query-conditioned perplexity scoring, document reordering, and adaptive compression ratios.
- **CPC** (Liskavets et al., 2024): A contrastively trained sentence-ranking model that selects sentences based on semantic similarity to the query in embedding space. It is query-aware and trained on synthetic QA data.
- **Raw Attention** (Wang et al., 2024; Fang et al., 2025): A non-learning baseline that selects sentences by averaging attention weights from the final decoder token. This mimics attention-based heuristics used in prior work such as QUITO and AttentionRAG.
- **Random Selection**: Sentences are sampled uniformly at random until the token budget is met. Serves as a lower-bound reference.
- **Empty Context**: The model receives only the question without any retrieved context, serving as a zero-context baseline.

All baselines are evaluated under the same token budget and LLM generation setting for fair comparison.

## D Additional Chinese Results with GPT-3.5-Turbo

To assess the cross-lingual robustness of our method, we evaluate Sentinel on LongBench-Zh using GPT-3.5-Turbo as the inference model. We compare against LLMLingua and LLMLingua-2 baselines, which are evaluated under a 3,000-token input constraint. Sentinel uses only 2,000 tokens but consistently outperforms the baselines across all task categories, as shown in Table 7.

## E Additional Results on Proxy Model Size and Family

This section provides detailed experimental results of Sentinel using proxy models from different families and parameter sizes, complementing the aggregated analysis presented in the main paper. Table 6 reports the full breakdown across all LongBench tasks.

Methods	LongBench-En (2000-token constraint)				LongBench-Zh (2000-token constraint)			Overall AVG
	SingleDoc	MultiDoc	Summ.	En-AVG	SingleDoc	MultiDoc	Zh-AVG	
Sentinel (Qwen2.5-0.5B-Instruct)	37.73	46.16	23.03	35.64	62.24	18.57	40.41	38.02
Sentinel (Qwen2.5-1.5B-Instruct)	39.48	46.07	23.10	36.22	62.02	18.91	40.47	38.34
Sentinel (Qwen2.5-3B-Instruct)	39.53	47.97	23.06	36.85	62.04	19.23	40.63	38.74
Sentinel (Qwen2.5-7B-Instruct)	38.79	45.56	22.52	35.62	60.88	18.43	39.66	37.64
Sentinel (Llama-3.2-1B-Instruct)	39.43	44.96	21.90	35.43	60.64	19.18	39.91	37.67
Sentinel (Llama-3.2-3B-Instruct)	36.03	44.46	22.00	34.17	59.24	18.89	39.06	36.62
Sentinel (Llama-3.1-8B-Instruct)	36.58	45.15	22.90	34.87	60.84	19.07	39.95	37.41
Sentinel (Qwen3-0.6B)	38.12	42.55	22.77	34.48	60.04	18.51	39.27	36.88
Sentinel (Qwen3-1.7B)	36.52	42.06	22.29	33.62	60.79	17.96	39.38	36.50
Sentinel (Qwen3-4B)	37.15	43.17	22.67	34.33	59.68	17.74	38.71	36.52
Sentinel (Qwen3-8B)	36.31	42.19	22.15	33.55	60.74	17.77	39.26	36.40
Original Prompt	38.84	44.74	22.76	35.45	60.06	18.21	39.14	37.30

Table 6: Detailed Sentinel performance across proxy model families and sizes.

Methods	LongBench-Zh (GPT-3.5-Turbo, 3000-token constraint)						Compression Stats	
	SingleDoc	MultiDoc	Summ.	FewShot	Synth.	AVG	Tokens	$1/\tau$
LLMLingua	35.2	20.4	11.8	24.3	51.4	28.6	3,060	5x
LLMLingua-2	46.7	23.0	15.3	32.8	72.6	38.1	3,023	5x
Evaluated under 2000-token constraint								
Sentinel (Qwen-2.5-0.5B-Instruct)	<b>64.8</b>	<u>25.1</u>	14.3	<u>38.0</u>	<u>89.0</u>	<u>46.2</u>	1,932	5x
Sentinel (Qwen-2.5-1.5B-Instruct)	<u>63.3</u>	24.9	14.8	<b>40.3</b>	<b>95.0</b>	<b>47.6</b>	1,929	5x
Original Prompt	61.2	<b>28.7</b>	<b>16.0</b>	29.2	77.5	42.5	14,940	-

Table 7: Performance comparison on LongBench-Zh using GPT-3.5-Turbo. LLMLingua baselines are evaluated under a 3,000-token budget. Sentinel uses only 2,000 tokens but consistently outperforms the baselines, demonstrating effective compression across languages.

## F Ablation Details

**Effect of Probing Data Size.** We evaluate how training size affects probing quality. As shown in Table 8, performance remains stable across 500–3000 training examples, with only marginal gains. This suggests that even a small probing set can support effective compression.

**Feature Selection Details** To construct a compact attention-based feature set, we use the Minimum Redundancy Maximum Relevance (mRMR) algorithm. We first compute mutual information between each feature (i.e., attention head statistics) and the binary relevance label, selecting the most informative one. We then iteratively add features that maximize relevance while minimizing redundancy, measured via Pearson correlation with already selected features. The number of features is capped at the number of heads in a single decoder layer to ensure compactness and interpretability.

**Compression Ratio.** Table 9 reports results with varying compression ratios ( $\tau \in \{0.1, 0.2, 0.3, 0.4, 0.5\}$ ), under a fixed chunk size of 1024. Sentinel remains robust even at high compression, while Raw attention

deteriorates significantly.

## G Analysis of Negatively Weighted Attention Heads

To better understand the role of attention heads assigned negative weights by Sentinel (Qwen-2.5-0.5B-Instruct), we analyze their attention distributions on 100 examples from the HotpotQA dataset. This analysis examines which input components these heads predominantly attend to, and whether their negative contributions correspond to known non-informative attention behaviors.

**Analysis Setup.** We analyze attention patterns on 100 HotpotQA examples by grouping input tokens into four categories: (i) sink tokens (e.g., special tokens and structurally dominant positions), (ii) supporting evidence sentences, (iii) question tokens, and (iv) remaining context. For each attention head, we compute the average proportion of attention mass assigned to each category.

**Results.** As shown in Table 10, attention heads assigned strong negative weights by Sentinel predominantly attend to sink tokens or question tokens, while allocating little to no attention to supporting

Methods	LongBench-En (2000-token constraint)				LongBench-Zh (2000-token constraint)			Overall AVG
	SingleDoc	MultiDoc	Summ.	En-AVG	SingleDoc	MultiDoc	Zh-AVG	
Qwen-2.5-0.5B-Instruct (500)	37.29	46.94	23.25	35.83	62.04	18.42	40.23	38.03
Qwen-2.5-0.5B-Instruct (1000)	38.35	47.43	23.66	36.48	61.43	18.57	40.00	38.24
Qwen-2.5-0.5B-Instruct (2000)	36.70	47.48	22.89	35.69	61.57	18.76	40.16	37.92
Qwen-2.5-0.5B-Instruct (3000)	37.73	46.16	23.03	35.64	62.24	18.57	40.41	38.02

Table 8: Performance of 0.5B models with different probing sizes (500, 1000, 2000, 3000) on LongBench.

Methods	LongBench-En (2000-token constraint)				LongBench-Zh (2000-token constraint)			Overall AVG
	SingleDoc	MultiDoc	Summ.	En-AVG	SingleDoc	MultiDoc	Zh-AVG	
Raw Attention (ratio 0.1)	25.79	36.54	20.39	27.57	35.03	16.33	25.68	26.62
Raw Attention (ratio 0.2)	33.19	41.09	21.63	31.97	48.45	17.23	32.84	32.41
Raw Attention (ratio 0.3)	34.91	43.74	22.39	33.68	55.09	18.14	36.62	35.15
Raw Attention (ratio 0.4)	37.63	45.95	22.88	35.49	58.78	17.82	38.30	36.89
Raw Attention (ratio 0.5)	37.47	44.70	23.25	35.14	60.63	17.42	39.03	37.09
Sentinel (ratio 0.1)	37.72	41.47	22.58	33.93	58.96	19.36	39.16	36.55
Sentinel (ratio 0.2)	39.90	45.97	23.37	36.42	59.50	17.92	38.71	37.56
Sentinel (ratio 0.3)	39.45	46.51	23.86	36.61	60.98	18.68	39.83	38.22
Sentinel (ratio 0.4)	39.93	46.62	23.38	36.65	59.51	18.77	39.14	37.89
Sentinel (ratio 0.5)	38.60	46.77	23.54	36.30	61.41	18.44	39.92	38.11

Table 9: Performance across compression ratios (chunk size = 1024).

Layer	Head	Probe Weight	Sink	Supporting	Question	Others
11	1	-13.16	0.89	0.01	0.05	0.04
3	0	-12.83	0.74	0.01	0.18	0.03
3	10	-10.22	0.08	0.00	0.84	0.02
21	9	-9.95	0.01	0.00	0.98	0.01
14	5	-9.47	0.00	0.03	0.85	0.06
3	5	-9.11	0.74	0.04	0.03	0.18
9	11	-8.15	0.96	0.00	0.03	0.01

Table 10: Examples of attention heads assigned strong negative weights by Sentinel, showing attention mass concentrated on sink or question tokens rather than supporting evidence.

evidence. In contrast, positively weighted heads focus primarily on evidence-bearing context.

**Implications.** This analysis shows that negatively weighted heads capture structurally dominant but semantically uninformative behaviors, such as attention sinks or self-focused query attention. Explicitly down-weighting these heads allows Sentinel to suppress spurious attention patterns and decode context utilization more robustly than methods that rely on raw attention or positively identified retrieval heads alone.

## H Latency and Inference Efficiency

We evaluate end-to-end inference latency across different Sentinel configurations, focusing on the effects of chunk size and attention feature design. Table 11 reports average and median latency per sample on the English LongBench dataset, mea-

sured on a single A100 GPU.<sup>1</sup>

With a chunk size of 1024 and All Layers attention features, Sentinel achieves  $1.13\times$  speedup over LLMingua-2 while reaching 38.02 F1.

To further improve runtime, we evaluate SENTINEL (SELECTED), which uses compact mRMR-selected features. At chunk size 1024, this variant reduces latency to 0.60s ( $1.30\times$ ) with only minor performance degradation (37.24 F1), offering an efficient alternative for low-latency scenarios.

## I LLM Evaluation Settings

For LLM-based evaluation, we adopt the official prompt templates and decoding settings from LongBench (Bai et al., 2024) to ensure consistency and comparability across methods. Unless otherwise specified, all decoding parameters are fixed for all datasets: the temperature is set to 0.0, the nucleus sampling parameter  $top\_p$  is 1.0, the random seed is fixed to 42, only a single generation is sampled ( $n = 1$ ), and streaming is disabled.

<sup>1</sup>We monkey-patch the model to extract only the final-token attention used by our method, replacing other activations with None to reduce overhead.

Method	Chunk Size	Proxy Model	Avg. Time (s) ↓	Med. Time (s) ↓	Speedup vs. LLMingua-2 ↑	Overall-AVG
LLMLingua-2 (trained)	512	XLNet-RoBERTa-Large (561M)	0.78	0.70	1.00×	28.35
Raw Attention	512	Qwen-2.5-0.5B-Instruct (494M)	1.01	0.84	0.77×	35.18
Raw Attention	1024	Qwen-2.5-0.5B-Instruct (494M)	0.65	0.54	1.20×	35.03
Sentinel (ours)	512	Qwen-2.5-0.5B-Instruct (494M)	1.02	0.84	0.76×	37.05
Sentinel (ours)	1024	Qwen-2.5-0.5B-Instruct (494M)	0.69	0.57	1.13×	38.02
Sentinel (ours) selected	1024	Qwen-2.5-0.5B-Instruct (494M)	0.60	0.49	1.30×	37.24

Table 11: Inference latency per QA sample on the full LongBench dataset (lower is better). LLMingua-2 is trained for token-level compression and limited to 512-token chunks. Sentinel uses a smaller, untrained decoder-only proxy and supports larger chunk sizes for improved efficiency. Speedup is relative to LLMingua-2 (chunk size 512). **Overall-AVG** denotes average accuracy across all chunk settings per method.