

Concept than Document: Context Compression via AMR-based Conceptual Entropy

Anonymous ACL submission

Abstract

Large Language Models (LLMs) face information overload when handling long contexts, particularly in Retrieval-Augmented Generation (RAG) where extensive supporting documents introduce redundant content that interferes with reasoning. Context engineering has emerged to address these challenges, yet existing methods rely on lexical or token-level features that fragment semantic units and fail to capture conceptually essential content. We propose an unsupervised context compression framework leveraging Abstract Meaning Representation (AMR) to preserve semantically essential information while filtering irrelevant text. By quantifying node-level entropy within AMR graphs, our method estimates the conceptual importance of each node, enabling retention of core semantics. Specifically, we construct AMR graphs from retrieved contexts, compute the conceptual entropy of each node, and identify statistically significant concepts to form a condensed, semantically focused context. Experiments on the PopQA and EntityQuestions datasets demonstrate that our method outperforms vanilla RAG and existing baselines, achieving superior accuracy while substantially reducing context length. To the best of our knowledge, this is the first work introducing AMR-based conceptual entropy for context compression, demonstrating the potential of structured linguistic representations in context engineering.

1 Introduction

Large Language Models (LLMs) are increasingly equipped with mechanisms to incorporate long contexts, allowing them to leverage external information beyond their training data (Lewis et al., 2020; Karpukhin et al., 2020). However, as the context length grows, LLMs often struggle to effectively identify and utilize truly relevant information, leading to performance degradation and inefficiency. This challenge, reflecting the trade-off between retrieval recall and precision, becomes particularly

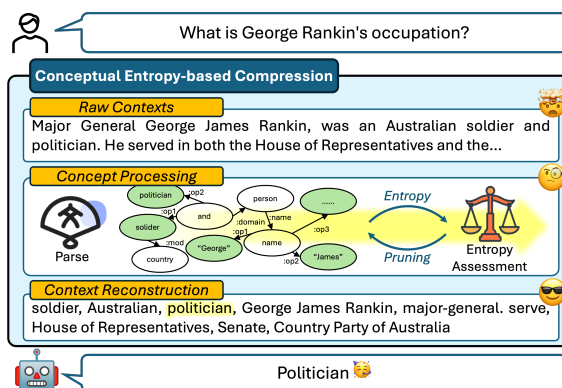


Figure 1: Long retrieved documents contain much irrelevant content; our method keeps only key AMR-based concepts to form a semantically focused context.

acute in scenarios such as Retrieval-Augmented Generation (RAG), where the inclusion of more retrieved documents raises the chance of accessing useful knowledge but simultaneously introduces overwhelming amounts of irrelevant text obscure the key facts (Shi et al., 2023; Jin et al., 2025a).

Context engineering has therefore become an effective strategy for enhancing the quality and efficiency of long-context utilization, aiming to distill essential information while reducing noise and redundancy (Mei et al., 2025). Existing approaches primarily focus on lexical or surface-level features for information filtering (Xu et al., 2024; Cheng et al., 2024). While these methods work well for certain queries, they may struggle with capturing complex semantic concepts and preserving factually important information. Moreover, traditional compression techniques may inadvertently remove crucial supporting evidence while retaining superficially relevant but semantically vacuous content.

To address the aforementioned limitations, we propose a novel context compression method that leverages Abstract Meaning Representation (AMR) (Banarescu et al., 2013) to identify and preserve semantically essential information. AMR

069	graphs provide a structured representation that ab-	identify and preserve core semantic informa-	120
070	stracts away from surface syntactic variations while	tion while filtering redundant content.	121
071	retaining core semantic content (Chen et al., 2025).		
072	Concepts assuming diverse semantic roles across	• Extensive experiments demonstrate that the	122
073	contexts naturally carry more informative value	proposed method outperforms vanilla and	123
074	for inference (Kuhn et al., 2023), which can be	other compression baselines by maintaining	124
075	quantified as higher information entropy in the role	robust semantic core preservation.	125
076	distribution of concept nodes (Nguyen et al., 2025).		
077	Moreover, cognitive studies suggest that the human	• The method achieves reductions in context	126
078	brain can automatically reconstruct scenarios im-	length and latency while preserving semantic	127
079	plied by essential concepts through pre-learned se-	integrity, offering a linguistically empowered	128
080	semantic knowledge (Binder et al., 2009; Horikawa,	framework for context engineering.	129
081	2025), and LLMs exhibit a similar capacity for		
082	concept-based scene understanding, providing the-	2 Related Work	130
083	oretical support for prioritizing semantically funda-		
084	mental concepts during reasoning (Du et al., 2025).	2.1 Context Engineering	131
085	Building on this foundation, our method con-	Context engineering has become a key strategy	132
086	structs AMR graphs from retrieved contexts to rep-	for managing and structuring information in LLM	133
087	resent entities and key semantic concepts in a struc-	workflows (Mei et al., 2025; Verma, 2024; Shi	134
088	tured form. For each concept node, we calculate	et al., 2024). Early approaches selected relevant	135
089	the information entropy to estimate its semantic	sentences or passages based on lexical similar-	136
090	contribution based on its inherent contextual uncer-	ity (Hwang et al., 2024), while other methods used	137
091	tainty. We then apply significance testing to iden-	neural models to reorganize retrieved contexts (Xu	138
092	tify concept nodes that exhibit statistically reliable	et al., 2024; Liu et al., 2024). Recent work exam-	139
093	informational salience, which constitute the struc-	ines learned context engineering techniques that op-	140
094	tural basis for reconstructing a compressed con-	optimize representations for downstream tasks. Jiang	141
095	text that preserves essential semantic content while	et al. (2024) uses instruction tuning to refine con-	142
096	suppressing redundant information. To mitigate	texts while preserving task-relevant information.	143
097	potential distortion caused by AMR’s abstraction	Selective-Context (Li et al., 2023b) applies atten-	144
098	from surface realization, the distilled concepts are	tion mechanisms to highlight critical segments. Jin	145
099	restored to their original textual expressions in the	et al. (2025b) emphasizes semantic integrity in	146
100	source contexts, ensuring factual consistency and	engineered contexts, integrating natural language	147
101	maintaining semantic clarity in the reconstructed	spans and semantic vectors to support dynamic evi-	148
102	compressed context for reasoning.	dence selection and improve answer quality.	149
103	We evaluate our method on two challeng-	2.2 AMR-enhanced Large Language Models	150
104	ing knowledge-intensive Q&A benchmarks,	Abstract Meaning Representation provides a struc-	151
105	PopQA (Mallen et al., 2023) and EntityQues-	tured formalism that abstracts away from syntac-	152
106	tions (Sciavolino et al., 2021), which require	tic variations, making it suitable for cross-lingual	153
107	reasoning over long-context factual evidence	and cross-domain applications (Wein and Opitz,	154
108	retrieved from external sources. Experimental	2024). Recent AMR parsing advances have made	155
109	results show substantial performance gains over	it practical to construct high-quality graphs from	156
110	vanilla RAG and other context compression	context (Bevilacqua et al., 2021; Zhou et al., 2021),	157
111	baselines, with more pronounced improvements	enabling applications across NLP tasks(Li et al.,	158
112	on instances involving long supporting documents.	2021; Liu et al., 2015; Song et al., 2019). With	159
113	These findings support our hypothesis that	the rise of LLMs, researchers have explored us-	160
114	AMR-based entropy filtering effectively isolates	ing AMR for semantic enhancement. Recent stud-	161
115	core semantic content while removing redundant	ies have examined AMR-driven chain-of-thought	162
116	information. The main contributions of this work	prompting, showing that structured semantic rep-	163
117	can be summarized as follows:	resentations can improve LLM performance across	164
		tasks (Jin et al., 2024). Zhang et al. (2025) inte-	165
118	• We propose a novel unsupervised context com-	grates AMR into LLM frameworks through struc-	166
119	pression framework that leverages AMR to	tured representation methods, although aligning	167

AMR’s graph structure with sequential processing remains challenging. These observations suggest that while full graph utilization is non-trivial, AMR nodes constitute stable semantic units that encode informative conceptual content, making them suitable for structured context engineering.

2.3 Information Theory in LLMs

Information-theoretic measures have become increasingly important in the era of LLMs, providing principled tools to understand and improve model behavior (Wang et al., 2025). LLMs have leveraged such analyses for interpretation and optimization (Nikitin et al., 2024). For instance, entropy-based selection of demonstration examples has been shown to enhance the performance of CoT prompting (Zhou et al., 2023). Beyond prompting, information-theoretic approaches have been applied to model compression, knowledge distillation, and efficient fine-tuning (Yin et al., 2024; Mao et al., 2024). These studies illustrate an emerging trend in which information theory provides both theoretical insights and practical tools for working with LLMs (Agarwal et al., 2025). In this work, we integrate graphical information-theoretic principles of AMR, leveraging high-entropy nodes as concise and informative representations of long contexts.

3 Methodology

3.1 Problem Formulation

The framework for transferring the context in raw documents to condensed concepts is as Figure 2. Given a query Q and a set of retrieved documents $D = \{d_1, d_2, \dots, d_n\}$ with corresponding correct answers $A = \{a_1, a_2, \dots, a_m\}$, our objective is to generate a compressed context C' that preserves the most semantically informative concepts essential for answering Q to yield $a_j \in A$, while substantially reducing the overall context length.

To create a controlled experiment that focuses exclusively on the impact of core concepts within the context on answer accuracy, we retain only documents that contain correct answers. This controlled setting enables us to isolate how our compression method affects the preservation of essential contextual information by eliminating interference from irrelevant documents. The hypothesis can be formalized as: $\forall d_i \in D, \exists a_j \in A$ such that $a_j \in d_i$.

Formally, we aim to learn a compression function $f(D) \rightarrow C'$ such that:

$$Acc(q, C') \gtrsim Acc(q, D) \text{ and } |C'| \ll |D| \quad (1)$$

where $C' \subseteq D$, $|C'|$ and $|D|$ are the lengths of the compressed and original contexts, respectively.

3.2 AMR Graph Construction

For each document $d_i \in D$, we construct it to the sentence-level AMR graphs with an mBart-based parser¹ trained in the AMR 3.0 corpus² to address potential multilingual concerns. Let $G_i = (V_i, E_i)$ denote the AMR graph for document d_i , where V_i represents the set of concept nodes and E_i represents the semantic relations between concepts. Each concept node $v \in V_i$ corresponds to a semantic concept (e.g., entities, predicates, or modifiers) and is associated with its textual realization in the raw document. The edges in E_i represent semantic relationships such as agent-of (ARG0), patient-of (ARG1), and various semantic roles.

Our approach is grounded in the cognitive hypothesis that both human comprehension and LLM inference can effectively reconstruct semantic scenarios from discrete informative concepts without explicit relational encoding (Xu et al., 2025; Fedorenko et al., 2024; Rogers et al., 2004; Wit and Gillette, 1999). This principle suggests that intelligent systems possess inherent capabilities to infer implicit relationships between concepts based on their learned background knowledge and contextual co-occurrence patterns (Brown et al., 2020; Cao et al., 2023; Suresh et al., 2023). Building on these foundations, we keep the concept nodes V_i and discard the explicit E_i in each G_i . This design ensures that the compressed context consists of discrete semantic concepts, avoiding the introduction of artificial relational symbols that may interfere with the LLM’s pre-trained language understanding capabilities while leveraging the model’s intrinsic ability in concept-based scenario reconstruction.

3.3 Information Entropy Computation

To identify the most informative concepts within each AMR graph, we employ an information-theoretic approach based on token-level perplexity measurements. For each concept node $v \in V_i$, we calculate its information entropy by leveraging the AMR generation model’s uncertainty when predicting the concept token sequence.

Given the AMR parsing model M with parameters θ , we obtain the probability distribution over the vocabulary for each token position in the

¹<https://github.com/BramVanroy/multilingual-text-to-amr>

²<https://catalog.ldc.upenn.edu/LDC2020T02>

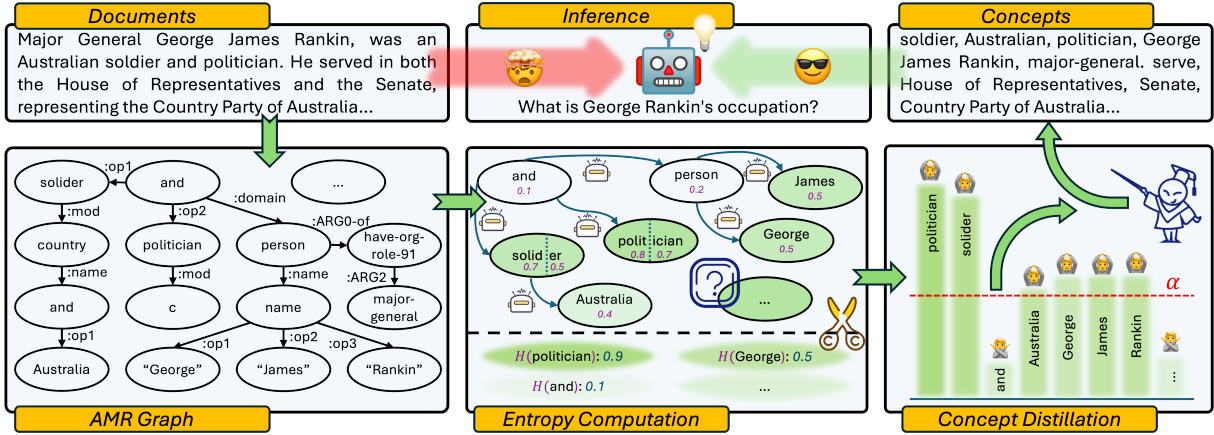


Figure 2: The conceptual entropy-based workflow converts the sparse context in raw supporting documents into condensed AMR-based concepts, forming a compact semantic representation for LLMs inference.

AMR linearization. However, modern tokenizers decompose words into subword units, requiring the aggregation to obtain concept-level entropy scores. For a concept v that corresponds to a complete word-level representation in d_i , the tokenizer may decompose it into the subword tokens $v = [s_1, s_2, \dots, s_m]$, $m \geq 0$. We compute the token-level entropy for each subword as:

$$E(s_j) = \exp(-\log P_\theta(s_j | s_{<j}, G_i)) \quad (2)$$

where $s_{<j}$ denotes the preceding tokens within the same concept. We aggregate token-level entropies into a concept-level entropy score as Eq. 3. Specifically, we identify concept boundaries by tracking tokens that begin with the special prefix "G" and accumulate entropy scores for all s_j belonging to the same conceptual unit. This aggregation strategy ensures that concepts composed of multiple subword tokens are not artificially penalized relative to single-token concepts. This alignment provides a balanced representation of the model's uncertainty across all subword components of a concept.

$$H(v) = \frac{1}{m} \sum_{j=1}^m E(s_j) \quad (3)$$

Compared to token-level entropy in linear text, computing entropy over AMR concept nodes leverages semantic structure to more precisely estimate informational content. High-entropy nodes often represent content-specific, less redundant meanings, thus providing more discriminative signals for downstream reasoning. This enables the compression process to highlight semantically rich units that may be obscured in the surface text.

3.4 Concept Distillation

The supporting document set D can be conceptualized as a coherent descriptive scenario corresponding to query Q , within which genuinely informative concepts can be identified through their statistically significant entropy deviations. Concepts exhibiting higher entropy relative to the general nodes carry more discriminative information and are thus more valuable for answering the query. For each $d_i \in D$ with concept entropy $\{H(v_1), H(v_2), \dots, H(v_{|V_i|})\}$, we perform a one-sample t-test to identify concepts with significantly higher information than the population mean:

$$t_{stat}(v_j) = \frac{H(v_j) - \bar{H}}{\frac{s}{\sqrt{n}}} \quad (4)$$

where \bar{H} is the sample mean entropy, s is the sample standard deviation, and $n = |V_i|$. We compute the corresponding p-value using the t-distribution with $n - 1$ degrees of freedom:

$$p(v_j) = 2 \times (1 - F_t(|t_{stat}(v_j)|, n - 1)) \quad (5)$$

where F_t is the cumulative distribution function. We then screen out concepts whose p-values satisfy $p(v_j) < \alpha$ as statistically significant high-information concepts. Our goal is not to identify only the most informative concepts, but rather to eliminate overly generic ones while preserving a relative conceptual basis for LLMs to infer the document's semantics. Considering the empirical validation of LLMs' inference, we adopt a relaxed threshold, $\alpha = 0.3$. This setting prevents the over-pruning of moderately informative concepts,

thereby ensuring that the retained set includes contextual signals. The ablation study to verify the different α settings is in Section C.

3.5 Context Compression and Reconstruction

The final compressed context C' is constructed by aggregating the concepts with significant entropy across all documents in D . For each document d_i , let $V_i = \{v \in V_i : p(v) < \alpha\}$ denote the set of statistically significant concepts. For each $c_i \in C'$, the compressed representation for document d_i is:

$$c_i = \bigodot_{v \in V_i} \phi(v) \quad (6)$$

where $\phi(\cdot)$ maps each concept v to its processed surface form through a sequence of linguistic post-processing steps designed to preserve semantic coherence and ensure linguistic fluency. These include *Temporal Expression Reconstruction*, where date and time expressions fragmented during AMR parsing are converted into natural language format, such as transforming "month 7 year 2025" into "July 2025"; *Redundancy Removal*, which eliminates consecutive duplicate concepts to reduce repetition while maintaining semantic diversity; and *Surface Realization*, which restores the processed concepts to their original textual forms in the raw document to mitigate potential distortions introduced by the AMR parsing process. This compressed form serves as the final input context, preserving the essential semantic signals while substantially reducing the original context length.

4 Experiments

4.1 Datasets and Implementation Details

We conduct comprehensive evaluations on two widely-adopted open-domain question-answering datasets that provide long-context supporting documents for RAG-based inference: **PopQA** (Mallen et al., 2023) and **EntityQuestions** (Sciavolino et al., 2021). For comprehensive evaluation, we use Contriever (Izacard et al., 2022) as the retriever for PopQA and BM25 (Robertson et al., 2009) for EntityQuestions, with retriever optimization beyond the scope of this work. Both datasets are equipped with ground-truth annotations indicating whether each supporting document contains the correct answer, denoted by the boolean indicator "hasanswer". To align the problem formulation in Eq. 1, we retain only documents where "hasanswer" = True, ensuring that performance variations stem from com-

pression effectiveness rather than irrelevant document interference. For each query Q , let K denote the number of answer-containing documents in the filtered D . The statistical characteristics of the curated $\langle Q, A, D \rangle$ triplets are summarized as follows:

Table 1: Statistical results of the amount of screened-out $\langle Q, A, D \rangle$ pairs from the datasets.

$K=$	1	2	3	4	5	6	7	8	9	10
PopQA	280	298	174	172	160	153	149	155	135	125
EQ	489	572	373	295	239	199	179	169	130	113

To mitigate reliance on parametric knowledge in LLM inference, we employ a structured prompting that prioritizes externally provided evidence over internal memory. We adopt the instruction as follows: "[Refer to the following facts to answer the question. Facts: C' . Question: Q]". Given that prompt intensity significantly influences inference behavior (Wu et al., 2024), we frame the supporting concepts C' as "facts" to establish a constrained knowledge boundary that minimizes interference from potentially conflicting parametric knowledge.

4.2 Baseline Methods

Our baseline evaluation examines two key dimensions: (1) diverse backbone LLM architectures, and (2) alternative context compression techniques. For backbone LLMs, we select mainstream publicly available LLMs, including GPT-Neo (1.3b and 2.7b) (Black et al., 2021), OPT (1.3b and 2.7b) (Zhang et al., 2022), BLOOM LM (560m, 7b1) (Le Scao et al., 2022), LLaMA-2-chat (13b) (Touvron et al., 2023), Llama-3.1-Instruct (8b) (Dubey et al., 2024), DeepSeek-V2-Lite (16b) (DeepSeek-AI, 2024), and Qwen3 (32b) (Team, 2025). The combination of backbone LLMs with contexts in raw supporting documents constitutes the *Vanilla* baseline.

For context compression, we implement five representative approaches that span different paradigms. We categorize these methods into three groups to answer the following questions: **Q1:** Can simple frequency-based measures suffice for identifying informative content? (*Statistical Method*). **Q2:** Can LLMs perform compression effectively through prompt-based reasoning? (*LLMs-driven Methods*). **Q3:** Can dedicated context compression models be more targeted and effective? (*Compression-specific Methods*).

The baselines corresponding to the above questions are as follows: (1) *Statistical Method*: TF-IDF, the statistical entropy-inspired method that

417 identifies salient terms using frequency–inverse
418 document frequency weighting to highlight in-
419 formative concepts. (2) *LLMs-driven Methods*:
420 prompt-based keyword extraction and summariza-
421 tion that leverage LLaMA-3.1-8B-Instruct with
422 prompts as Prompt A1 and Prompt A2 to generate
423 keywords and summarizations. (3) *Compression-*
424 *specific Methods*: Selective Context (SelCon) (Li
425 et al., 2023a) that employs trained models to iden-
426 tify relevant spans, and LLMLingua (Jiang et al.,
427 2023) uses budget-constrained token selection for
428 optimal compression. These baselines evaluate if
429 compressed contexts can preserve essential infor-
430 mation while reducing computational overhead.

431 4.3 Evaluation Metrics

432 We employ two evaluation metrics: accuracy (Acc)
433 and Area Under the Curve (AUC). The standard de-
434 viation (σ) of AUC is used as an auxiliary metric.
435 The Acc follows the exact match protocol of Mallen
436 et al. (2023), measuring if any generated answer
437 exactly matches any gold-standard $a_j \in A$ for a
438 given query Q . The σ assesses the stability of com-
439 pressed methods across different backbone LLMs.

440 The AUC provides a comprehensive assessment
441 across varying K . Specifically, AUC computes the
442 area under the Acc curve against K . Higher AUC
443 indicates superior overall performance across the
444 corresponding intervals. Given our focus on long-
445 context compression, we partition the AUC calcu-
446 lation into two intervals for the values of K : a
447 standard interval $I_s = [1, 10]$ that captures gen-
448 eral performance trends and a long-context interval
449 $I_l = [6, 10]$ that highlights performance under long
450 context. This decomposition provides clear insights
451 into both typical and challenging scenarios.

452 5 Results and Analysis

453 5.1 Overall Performance

454 The AUC results in I_s interval in Table 2 and Table 3
455 present the overall performance comparison across
456 both datasets. The full results in Acc are in Ta-
457 ble A1 and A2 respectively. In the PopQA dataset,
458 the proposed method achieves substantial gains
459 compared to the vanilla baseline. The most notable
460 improvements occur in larger models like Qwen3-
461 32B, Llama-2-chat-13b, and DeepSeek-V2-Lite.
462 In contrast, smaller models like Bloom-560m/7b1
463 show relatively modest improvements. On the
464 EntityQuestions dataset, the results exhibit sim-
465 ilar trends with some variations. The proposed

466 method achieves the best or second-best perfor-
467 mance across most configurations, with particularly
468 strong results on larger models like Qwen3-32B.
469 However, we observe slight performance degrada-
470 tion compared to vanilla on smaller models like
471 GPT-Neo-1.3B and Bloom-560m/7b1. Consider-
472 ing the previous observation, this phenomenon in-
473 dicates that compact LLMs may benefit from more
474 contextual information that retains rich linguistic
475 elements to reconstruct scenarios rather than ag-
476 gressive compression. This suggests a trade-off
477 between compression ratio and model capacity that
478 warrants consideration in practical deployments.
479 In addition, our method achieves a competitive σ
480 across diverse backbone LLMs, indicating it pre-
481 serves universally shared semantic cores rather than
482 model-specific preferences, forming a robust se-
483 mantic compression that maintains coherent rea-
484 soning chains across different architectures.

485 Compared to compression baselines, our method
486 demonstrates substantial advantages across differ-
487 ent paradigms. Against the statistical TF-IDF ap-
488 proach, we achieve overwhelming superiority on
489 both datasets, outperforming all backbone LLMs.
490 Although TF-IDF outperforms the vanilla setting
491 on certain backbone models, this improvement is
492 not consistent when examined across different ar-
493 chitectures, as indicated by the unstable results
494 with the highest σ . Its performance depends on
495 surface-level lexical patterns, which may occasion-
496 ally align with answer-bearing spans in simple con-
497 texts. However, TF-IDF lacks semantic structure
498 awareness and does not model how LLMs recon-
499 struct contextual meaning. As a result, it may either
500 discard essential cues or retain redundant tokens
501 that vary across models. The fluctuating perfor-
502 mance across backbones indicates answers of Q1
503 that frequency-based signals are insufficient for
504 reliably identifying informative content.

505 The LLM-driven baselines, Keywords and Sum-
506 mary, show limited performance in most settings.
507 Unlike statistical measures, these baselines depend
508 on generative rewriting, which makes them sen-
509 sitive to semantic integrity and prompts. These
510 factors lead to unreliable results across different
511 backbones. In addition, the generative paradigm
512 can introduce hallucinations into the rewritten
513 content, further increasing the uncertainty of the
514 compressed context. A notable trend is that the
515 summary-based compression achieves the lowest σ .
516 The reason is summary-compressed context is re-
517 main natural language, forming a continuous repre-

Table 2: The quantitative results of AUC \uparrow for the PopQA dataset, where the full name order of the LLMs is: GPT-Neo-1.3B, GPT-Neo-2.7B, OPT-1.3b, OPT-2.7b, Bloom-560m, Bloom-7b1, Llama-2-chat-13b, Llama-3.1-8B-Instruct, DeepSeek-V2-Lite, Qwen3-32B. The standard division is as $\sigma \downarrow$. The best results are in **bold**, and the second-best results are in underlined. The **increased** and **decreased** Δ are marked differently.

D	K	G-1.3	G-2.7	O-1.3	O-2.7	b-560	b-7b1	L-13	L3.1-8	DS-V2	Q3-32	$\sigma \downarrow$
Vanilla	I_s	553.32	550.79	585.12	596.31	<u>575.04</u>	664.92	583.57	701.36	575.00	251.99	119.63
	I_l	262.07	252.04	278.86	282.63	<u>284.04</u>	<u>318.37</u>	293.42	337.14	303.30	101.33	64.77
TF-IDF	I_s	354.04	508.48	486.22	523.84	<u>417.67</u>	608.85	623.00	650.98	179.28	210.62	165.39
	I_l	169.82	251.12	244.02	269.09	217.52	307.70	311.47	316.14	106.47	113.97	78.00
Keywords	I_s	423.52	449.40	532.66	547.01	497.93	588.64	552.55	606.34	295.62	271.88	116.40
	I_l	193.41	211.08	264.65	274.44	252.10	294.34	278.92	302.44	173.88	141.73	55.10
Summary	I_s	433.24	459.55	540.52	504.34	527.49	577.91	482.79	551.42	491.56	285.17	82.74
	I_l	206.04	223.84	267.55	242.91	268.18	294.93	252.27	270.41	269.50	138.74	44.81
SelCon	I_s	453.31	490.44	580.08	581.62	443.08	634.40	637.20	717.74	557.43	293.10	121.98
	I_l	209.18	228.22	286.68	284.62	216.25	307.80	309.70	339.02	295.48	156.93	57.34
Lingua	I_s	<u>554.94</u>	<u>553.15</u>	<u>607.40</u>	<u>617.07</u>	567.67	<u>665.73</u>	<u>645.21</u>	<u>743.76</u>	<u>643.01</u>	<u>325.39</u>	110.21
	I_l	263.89	<u>258.09</u>	<u>292.36</u>	<u>286.70</u>	280.85	317.55	<u>312.28</u>	<u>346.24</u>	318.18	163.83	50.08
Ours	I_s	600.62	611.43	625.14	648.91	587.98	677.77	678.51	756.44	648.90	356.55	<u>104.32</u>
	I_l	283.54	296.09	298.73	308.92	292.74	332.16	326.67	357.74	<u>318.06</u>	191.09	44.33
Δ	I_s	+47.30	+60.64	+40.02	+52.60	+12.94	+12.85	+94.94	+55.08	+73.90	+104.56	30.32
	I_l	+21.47	+44.05	+19.87	+26.29	+8.70	+13.79	+33.25	+20.60	+14.76	+89.76	23.57

Table 3: The AUC \uparrow results for the EntityQuestions dataset. The symbol definitions are same as Table 2.

D	K	G-1.3	G-2.7	O-1.3	O-2.7	b-560	b-7b1	L-13	L3.1-8	DS-V2	Q3-32	$\sigma \downarrow$
Vanilla	I_s	550.08	<u>608.54</u>	<u>618.05</u>	677.63	511.98	705.35	657.06	743.99	572.72	235.42	142.98
	I_l	259.35	<u>283.86</u>	<u>284.91</u>	318.26	236.82	329.58	296.63	338.60	313.36	87.65	72.88
TF-IDF	I_s	302.59	459.72	419.50	517.23	314.45	552.43	666.08	627.44	180.75	235.64	165.91
	I_l	146.52	239.16	188.60	259.91	155.99	273.13	<u>323.23</u>	276.02	107.46	112.64	75.92
Keywords	I_s	358.34	458.67	495.48	545.41	392.71	572.18	614.18	674.23	284.15	287.12	135.78
	I_l	171.09	229.08	245.89	276.19	190.40	282.74	310.78	323.42	175.65	128.99	65.40
Summary	I_s	336.92	366.90	450.84	437.40	396.18	498.25	435.01	511.30	448.16	210.08	88.12
	I_l	161.38	180.04	221.94	202.50	196.11	254.38	209.77	242.42	247.76	77.62	52.17
SelCon	I_s	278.08	329.18	359.08	391.45	251.39	401.26	531.96	545.13	395.29	226.52	<u>107.42</u>
	I_l	136.32	163.02	177.21	187.91	137.72	195.78	268.26	259.44	208.08	103.98	<u>52.52</u>
Lingua	I_s	541.93	598.45	592.69	644.01	<u>496.46</u>	670.92	<u>698.64</u>	<u>792.93</u>	<u>648.58</u>	<u>374.74</u>	115.86
	I_l	244.38	275.40	274.64	283.11	223.05	308.36	322.57	<u>357.82</u>	<u>307.12</u>	152.43	57.73
Ours	I_s	<u>546.46</u>	627.41	632.79	<u>662.16</u>	494.45	<u>688.73</u>	738.82	813.86	652.14	406.00	118.33
	I_l	<u>248.82</u>	294.48	298.31	<u>295.18</u>	<u>229.06</u>	<u>323.26</u>	343.58	371.30	307.05	181.50	55.95
Δ	I_s	-3.62	+18.87	+14.74	-15.47	-17.53	-16.62	+81.76	+69.87	+79.42	+170.58	61.27
	I_l	-10.53	+10.62	+13.40	-23.08	-7.76	-6.32	+46.95	+32.70	-6.31	+93.85	35.11

sensation showing lower sensitivity to surface-level changes. In contrast, the discrete keywords-based compression shows notable performance swings. These observations answer **Q2** by showing that LLM-driven baselines are not a reliable choice due to the uncertainty in inference.

Compared with the SelCon baseline, our method achieves higher AUC across configurations. We hypothesize that this gap stems from fundamental differences in our approaches: while both methods utilize information theory, SelCon operates at the phrase/sentence level through token-based self-information aggregation for content filtering, whereas our method uses AMR’s structured semantic representation to compute concept-level entropy based on semantic roles and connections in comprehensive contexts. The AMR-based entropy better preserves the conceptual coherence for complex

reasoning, as it captures semantic structures and dependencies that are crucial for maintaining clear inferential chains for reconstructing scenarios.

LLMLingua represents competitive baseline as a token-level compression technique. The advantage of our method relative to LLMLingua comes from the complementary strengths of semantic-level versus token-level compression: while LLMLingua selects tokens through iterative perplexity-based filtering and budget control, our AMR-based approach identifies coherent concept units that match the information structure. Both methods preserve essential information, but our semantic abstraction excels when maintaining conceptual relationships matters more than surface-level linguistic continuity. Moreover, our method enhances the interpretability and readability by preserving complete conceptual units as atomic elements and maintain-

ing lexical integrity, whereas token-level compression can fragment words that disrupt local linguistic structures. This property facilitates human understanding and debugging. Compared with other baselines, both SelCon and LLMingua achieve competitive AUC and σ , addressing **Q3** on the necessity of dedicated context compression methods.

5.2 Performance on Long Contexts

To further validate our method and highlight its characteristics, we analyze performance in the long-context interval I_l in Table 2 and Table 3, emphasizing behaviors that emerge specifically under long-context conditions. The proposed method achieves the competitive performance that keeps the same trend as in the I_s , but the gains are reduced. The reduction is expected since the I_l interval typically encompasses longer contexts or higher complexity scenarios, where the marginal benefit of improvements tends to diminish. However, a notable phenomenon is that σ is significantly lower for this interval, which contains longer but more concentrated concepts compared with the massive but dispersed interval, indicating the benefit of macro-level semantic constraints in capturing informative concepts within complex contexts in specific scenarios. Moreover, the low σ of Δ indicates consistent performance variance across backbones.

5.3 Compression Efficiency

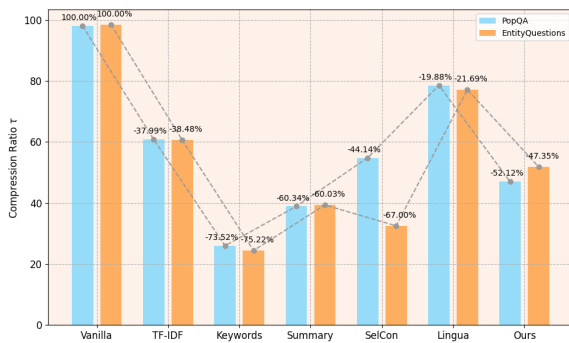


Figure 3: Comparison of token-level compression ratios across different context compression methods.

We examine the compression efficiency in terms of token-level reduction (τ) and inference latency (ms per instance). As shown in Figure 3, our method reduces the length to about 50% of the vanilla on average, while keeping the Acc stable in both datasets. Baselines such as Keywords and Summary yield lower token counts, but they often remove meaningful factual cues, leading to perfor-

mance drops. In contrast, operating at the concept level through AMR allows the compressed context to retain the core semantic units needed for reasoning, rather than relying on surface lexical signals.

Table 4: Inference time comparison (ms per instance)

LLMs	Vanilla	TF-IDF	Keywords	Summary	SelCon	Lingua	Ours
PopQA							
G-1.3	402.89	468.01	366.23	410.52	470.27	429.32	380.38
G-2.7	672.12	622.81	548.13	578.64	634.32	640.69	548.51
O-1.3	322.68	314.18	281.84	316.45	305.92	356.40	306.23
O-2.7	517.73	499.23	484.59	487.71	524.98	526.04	461.01
b-560	261.43	265.57	235.13	237.32	275.10	274.51	249.55
b-7b1	1130.13	1152.86	1006.23	1006.60	1150.33	1139.21	1058.83
L-13	1886.29	1405.44	1329.71	1364.58	1476.61	1507.88	1409.22
L3.1-8	1032.17	1091.39	688.09	644.89	1109.62	1089.12	888.58
DS-V2	1233.80	166.51	150.13	165.69	293.25	171.14	164.05
Q3-32	5283.06	5029.57	4795.76	4879.41	5094.38	5040.01	4783.34
EntityQuestions							
G-1.3	605.82	587.49	546.79	547.65	724.10	761.94	585.63
G-2.7	866.79	811.66	749.83	746.46	867.93	932.25	779.43
O-1.3	528.14	486.72	481.75	496.73	557.98	648.45	499.03
O-2.7	703.28	684.91	647.43	671.47	761.57	827.20	702.26
b-560	445.16	468.14	421.71	416.70	527.12	582.26	439.89
b-7b1	1319.82	1338.23	1196.88	1176.98	1190.92	1456.20	1279.33
L-13	1805.86	1786.85	1672.17	1743.26	1717.31	1881.69	1590.65
L3.1-8	1233.70	1282.96	871.17	836.18	1016.64	1398.92	1083.90
DS-V2	358.03	326.23	333.82	326.80	431.81	444.87	330.10
Q3-32	5239.69	5313.41	4996.61	5012.55	5120.52	5409.85	5168.47

The reduction in context length leads directly to faster inference, and the latency decreases in line with the length reduction. Table 4 shows that the proposed method lowers the average inference time compared to the vanilla setting. Baselines reducing latency via token pruning may fragment expressions and weaken local coherence, especially in long contexts. By retaining intact conceptual units, our compressed contexts remain stable for reasoning, enabling both shorter inference time and reliable answering, even under high compression.

6 Conclusion

This paper presents a compression method for context engineering that leverages conceptual information entropy of AMR to identify semantically crucial concepts. Our method shows improvements over baselines while achieving substantial compression ratios. The experiments demonstrate that AMR-based semantic analysis guides context compression effectively. The integration of structured linguistic representation with information-theoretic concept selection offers a paradigm to balance information retention with computational efficiency.

Future research includes extending our approach to multi-modal contexts, modeling cross-document concept relationships, and exploring adaptive compression strategies based on query complexity. Incorporating other stable linguistic representations is also a valuable direction to improve the efficiency and effectiveness in context engineering.

624 Limitations

625 Although the proposed method shows clear gains
626 in long-context settings, some limitations remain.
627 First, the current approach relies on the stability
628 of AMR parsers, and the performance may de-
629 cline when the parser produces incomplete or noisy
630 graphs. The parsing processing is based on the
631 sentence-level graph, so complex document-level
632 structures are easily ignored. These dependency
633 introduces upper bounds on covered conceptual
634 information in compression. Developing reliable
635 AMR parsers is a continuously valuable direction.

636 Second, the current setup evaluates compres-
637 sion under a controlled testing environment where
638 answer-containing documents are considered. This
639 design isolates the effect of compression but does
640 not fully reflect real-world retrieval pipelines,
641 where irrelevant or conflicting documents are com-
642 mon. Experimenting with the setting in a full re-
643 trieval stack and examining different retrievers' in-
644 fluence will be conducted in future work.

645 Finally, computing AMR graphs and entropy
646 scores introduces extra cost during preprocessing.
647 Although this cost occurs offline, it may restrict
648 the method in latency-sensitive systems or in large-
649 scale applications where many documents must
650 be processed. A crucial future work is exploring
651 high-efficiency solutions for these stages.

652 References

653 Shivam Agarwal, Zimin Zhang, Lifan Yuan, Jiawei Han,
654 and Hao Peng. 2025. The unreasonable effectiveness
655 of entropy minimization in llm reasoning. [arXiv preprint arXiv:2505.15134](#).

657 Laura Banarescu, Claire Bonial, Shu Cai, Madalina
658 Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin
659 Knight, Philipp Koehn, Martha Palmer, and Nathan
660 Schneider. 2013. [Abstract Meaning Representation for sembanking](#). In [Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse](#), pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.

665 Michele Bevilacqua, Rexhina Blloshmi, and Roberto
666 Navigli. 2021. One spring to rule them both:
667 Symmetric amr semantic parsing and generation
668 without a complex pipeline. In [Proceedings of the AAAI conference on artificial intelligence](#), volume 35, pages 12564–12573.

671 Jeffrey R Binder, Rutvik H Desai, William W Graves,
672 and Lisa L Conant. 2009. Where is the semantic
673 system? a critical review and meta-analysis of 120
674 functional neuroimaging studies. [Cerebral cortex](#), 19(12):2767–2796.

Sid Black, Gao Leo, Phil Wang, Connor Leahy,
and Stella Biderman. 2021. [GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow](#). If you use this software, please cite it using these metadata.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. [Language models are few-shot learners](#). In [Advances in Neural Information Processing Systems](#), volume 33, pages 1877–1901. Curran Associates, Inc.

Qi Cao, Takeshi Kojima, Yutaka Matsuo, and Yusuke Iwasawa. 2023. [Unnatural error correction: GPT-4 can almost perfectly handle unnatural scrambled text](#). In [Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing](#), pages 8898–8913, Singapore. Association for Computational Linguistics.

Huiyao Chen, Meishan Zhang, Jing Li, Min Zhang, Lilja Øvrelid, Jan Hajič, and Hao Fei. 2025. [Semantic role labeling: A systematical survey](#). [arXiv preprint arXiv:2502.08660](#).

Xin Cheng, Xun Wang, Xingxing Zhang, Tao Ge, Si-Qing Chen, Furu Wei, Huishuai Zhang, and Dongyan Zhao. 2024. [xrag: Extreme context compression for retrieval-augmented generation with one token](#). [Advances in Neural Information Processing Systems](#), 37:109487–109516.

DeepSeek-AI. 2024. [Deepseek-v2: A strong, economical, and efficient mixture-of-experts language model](#). [Preprint](#), arXiv:2405.04434.

Changde Du, Kaicheng Fu, Bincheng Wen, Yi Sun, Jie Peng, Wei Wei, Ying Gao, Shengpei Wang, Chuncheng Zhang, Jinpeng Li, and 1 others. 2025. [Human-like object concept representations emerge naturally in multimodal large language models](#). [Nature Machine Intelligence](#), pages 1–16.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 82 others. 2024. [The llama 3 herd of models](#). [CoRR](#), abs/2407.21783.

Evelina Fedorenko, Steven T Piantadosi, and Edward AF Gibson. 2024. [Language is primarily a tool for communication rather than thought](#). [Nature](#), 630(8017):575–586.

Tomoyasu Horikawa. 2025. [Mind captioning: Evolving descriptive text of mental content from human brain activity](#). [Science Advances](#), 11(45):eadw1464.

732	Taeho Hwang, Sukmin Cho, Soyeong Jeong, Hoyun Song, SeungYoon Han, and Jong C Park. 2024. Exit: Context-aware extractive compression for enhancing retrieval-augmented generation. arXiv preprint arXiv:2412.12559 .	788
733		789
734		790
735		791
736		792
737	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised dense information retrieval with contrastive learning . Transactions on Machine Learning Research .	793
738		794
739		795
740		796
741		797
742	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, Singapore. Association for Computational Linguistics.	798
743		799
744		
745		800
746		801
747		802
748		803
749	Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. 2024. Longllmlingua: Accelerating and enhancing llms in long context scenarios via prompt compression . In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL) , pages 1658–1677.	804
750		805
751		806
752		807
753		808
754		809
755		810
756	Jiajie Jin, Xiaoxi Li, Guanting Dong, Yuyao Zhang, Yutao Zhu, Yongkang Wu, Zhonghua Li, Qi Ye, and Zhicheng Dou. 2025a. Hierarchical document refinement for long-context retrieval-augmented generation. arXiv preprint arXiv:2505.10413 .	811
757		812
758		813
759		814
760		815
761	Yiqiao Jin, Kartik Sharma, Vineeth Rakesh, Yingtong Dou, Menghai Pan, Mahashweta Das, and Srijan Kumar. 2025b. Sara: Selective and adaptive retrieval-augmented generation with context compression. arXiv preprint arXiv:2507.05633 .	816
762		817
763		818
764		819
765		820
766	Zhijing Jin, Yuen Chen, Fernando Gonzalez Adauto, Jiarui Liu, Jiayi Zhang, Julian Michael, Bernhard Schölkopf, and Mona Diab. 2024. Analyzing the role of semantic representations in the era of large language models . In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers) , pages 3781–3798, Mexico City, Mexico. Association for Computational Linguistics.	821
767		822
768		823
769		824
770		825
771		826
772		827
773		828
774		829
775		830
776	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering . In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) , pages 6769–6781. Association for Computational Linguistics.	831
777		832
778		833
779		834
780		835
781		836
782		837
783	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2023. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation . In The Eleventh International Conference on Learning Representations .	838
784		839
785		840
786		841
787		842
		843
		844
		845
	Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, and 1 others. 2022. Bloom: A 176b-parameter open-access multilingual language model.	
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Igor Kulikov, Vishrav Chaudhary, Sebastian Wang, Wen-tau Yih Barta, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks . In Advances in Neural Information Processing Systems , volume 33, pages 9459–9474.	
	Chenliang Li, Bin Bi, Ming Yan, Wei Wang, and Song-fang Huang. 2021. Addressing semantic drift in generative question answering with auxiliary extraction . In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers) , pages 942–947, Online. Association for Computational Linguistics.	
	Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. 2023a. Compressing context to enhance inference efficiency of large language models . In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing , pages 6342–6353, Singapore. Association for Computational Linguistics.	
	Yucheng Li, Bo Dong, Chenghua Lin, and Frank Guerin. 2023b. Compressing context to enhance inference efficiency of large language models . In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP) .	
	Fei Liu, Jeffrey Flanigan, Sam Thomson, Norman Sadeh, and Noah A. Smith. 2015. Toward abstractive summarization using semantic representations . In Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies , pages 1077–1086, Denver, Colorado. Association for Computational Linguistics.	
	Shengjie Liu, Jing Wu, Jingyuan Bao, Wenyi Wang, Naira Hovakimyan, and Christopher G Healey. 2024. Towards a robust retrieval-based summarization system . arXiv preprint arXiv:2403.19889 .	
	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories . In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) , pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.	
	Wenhao Mao, Chengbin Hou, Tianyu Zhang, Xinyu Lin, Ke Tang, and Hairong Lv. 2024. Parse trees guided llm prompt compression . arXiv preprint arXiv:2409.15395 .	

846	Lingrui Mei, Jiayu Yao, Yuyao Ge, Yiwei Wang, Bao-	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann	902
847	long Bi, Yujun Cai, Jiazhi Liu, Mingyu Li, Zhong-Zhi	Dubois, Xuechen Li, Carlos Guestrin, Percy Liang,	903
848	Li, Duzhen Zhang, and 1 others. 2025. A survey of	and Tatsunori B. Hashimoto. 2023. Stanford alpaca:	904
849	context engineering for large language models. arXiv	An instruction-following llama model. https://	905
850	preprint arXiv:2507.13334 .	github.com/tatsu-lab/stanford_alpaca .	906
851	Dang Nguyen, Ali Payani, and Baharan Mirzasoleiman.	Qwen Team. 2025. Qwen3 technical report . Preprint ,	907
852	2025. Beyond semantic entropy: Boosting LLM	arXiv:2505.09388 .	908
853	uncertainty quantification with pairwise semantic	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	909
854	similarity . In Findings of the Association for	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	910
855	Computational Linguistics: ACL 2025 , pages 4530–	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	911
856	4540, Vienna, Austria. Association for Computa-	Bhosale, and 1 others. 2023. Llama 2: Open founda-	912
857	tional Linguistics.	tion and fine-tuned chat models. arXiv preprint	913
858	Alexander Nikitin, Jannik Kossen, Yarin Gal, and	arXiv:2307.09288 .	914
859	Pekka Marttinen. 2024. Kernel language entropy:	Sourav Verma. 2024. Contextual compression in	915
860	Fine-grained uncertainty quantification for llms	retrieval-augmented generation for large language	916
861	from semantic similarities . In Advances in Neural	models: A survey. arXiv preprint arXiv:2409.13385 .	917
862	Information Processing Systems , volume 37, pages	Jingyao Wang, Wenwen Qiang, Zeen Song, Changwen	918
863	8901–8929. Curran Associates, Inc.	Zheng, and Hui Xiong. 2025. Learning to think:	919
864	Stephen Robertson, Hugo Zaragoza, and 1 others. 2009.	Information-theoretic reinforcement fine-tuning for	920
865	The probabilistic relevance framework: Bm25 and	llms. arXiv preprint arXiv:2505.10425 .	921
866	beyond. Foundations and Trends® in Information	Shira Wein and Juri Opitz. 2024. A survey	922
867	Retrieval , 3(4):333–389.	of AMR applications . In Proceedings of the	923
868	Timothy T Rogers, Matthew A Lambon Ralph, Peter	2024 Conference on Empirical Methods in Natural	924
869	Garrard, Sasha Bozeat, James L McClelland, John R	Language Processing , pages 6856–6875, Miami,	925
870	Hodges, and Karalyn Patterson. 2004. Structure and	Florida, USA. Association for Computational Lin-	926
871	deterioration of semantic memory: a neuropsycholog-	guistics.	927
872	ical and computational investigation. Psychological	EC Wit and Marie Gillette. 1999. What is linguistic	928
873	review , 111(1):205.	redundancy. University of Chicago .	929
874	Christopher Sciavolino, Zexuan Zhong, Jinhyuk Lee,	Kevin Wu, Eric Wu, and James Zou. 2024. Clashe-	930
875	and Danqi Chen. 2021. Simple entity-centric ques-	sheval: Quantifying the tug-of-war between an	931
876	tions challenge dense retrievers . In Proceedings	LLM’s internal prior and external evidence . In	932
877	of the 2021 Conference on Empirical Methods in	The Thirty-eight Conference on Neural Information	933
878	Natural Language Processing , pages 6138–6148, On-	Processing Systems Datasets and Benchmarks	934
879	line and Punta Cana, Dominican Republic. Associa-	Track .	935
880	tion for Computational Linguistics.	Fangyuan Xu, Weijia Shi, and Eunsol Choi. 2024. RE-	936
881	Chunting Shi, Michihiro Yasunaga, Isabelle Augenstein,	COMP: Improving retrieval-augmented LMs with	937
882	Nikos Voskarides, Mikel Artetxe, Xiang Ren, Xi-	context compression and selective augmentation . In	938
883	aozhong Wan, Antoine Bosselut, Dragomir Radev,	The Twelfth International Conference on Learning	939
884	Wenpeng Yin, and 1 others. 2023. Replug: Retrieval-	Representations .	940
885	augmented black-box language models . arXiv	Qihui Xu, Yingying Peng, Samuel A Nastase, Martin	941
886	preprint arXiv:2301.12652 .	Chodorow, Minghua Wu, and Ping Li. 2025. Large	942
887	Kaize Shi, Xueyao Sun, Qing Li, and Guandong Xu.	language models without grounding recover non-	943
888	2024. Compressing long context for enhancing rag	sensorimotor but not sensorimotor features of human	944
889	with amr-based concept distillation. arXiv preprint	concepts. Nature human behaviour , pages 1–16.	945
890	arXiv:2405.03085 .	Mingjia Yin, Chuhan Wu, Yufei Wang, Hao Wang, Wei	946
891	Linfeng Song, Daniel Gildea, Yue Zhang, Zhiguo	Guo, Yasheng Wang, Yong Liu, Ruiming Tang, Defu	947
892	Wang, and Jinsong Su. 2019. Semantic neural ma-	Lian, and Enhong Chen. 2024. Entropy law: The	948
893	chine translation using amr . Transactions of the	story behind data compression and llm performance.	949
894	Association for Computational Linguistics , 7:19–31.	arXiv preprint arXiv:2407.06645 .	950
895	Siddharth Suresh, Kushin Mukherjee, Xizheng Yu, Wei-	Jiahuan Zhang, Tianheng Wang, Hanqing Wu, Ziyi	951
896	Chun Huang, Lisa Padua, and Timothy Rogers. 2023.	Huang, Yulong Wu, Dongbai Chen, Linfeng Song,	952
897	Conceptual structure coheres in human cognition	Yue Zhang, Guozheng Rao, and Kaicheng Yu.	953
898	but not in large language models . In Proceedings	2025. Sr-llm: Rethinking the structured repre-	954
899	of the 2023 Conference on Empirical Methods in	sentation in large language model. arXiv preprint	955
900	Natural Language Processing , pages 722–738, Sin-	arXiv:2502.14352 .	956
901	gapore. Association for Computational Linguistics.		

- 957 Susan Zhang, Stephen Roller, Naman Goyal, Mikel
958 Artetxe, Moya Chen, Shuohui Chen, Christopher De-
959 wan, Mona Diab, Xian Li, Xi Victoria Lin, and 1
960 others. 2022. Opt: Open pre-trained transformer
961 language models. [arXiv preprint arXiv:2205.01068](https://arxiv.org/abs/2205.01068).
- 962 Chuyue Zhou, Wangjie You, Juntao Li, Jing Ye, Ke-
963 hai Chen, and Min Zhang. 2023. [INFORM : In-
964 formation eNtropy based multi-step reasoning FOR
965 large language models](#). In [Proceedings of the
966 2023 Conference on Empirical Methods in Natural
967 Language Processing](#), pages 3565–3576, Singapore.
968 Association for Computational Linguistics.
- 969 Jiawei Zhou, Tahira Naseem, Ramón Fernandez As-
970 tudillo, and Radu Florian. 2021. [AMR parsing with
971 action-pointer transformer](#). In [Proceedings of the
972 2021 Conference of the North American Chapter
973 of the Association for Computational Linguistics:
974 Human Language Technologies](#), pages 5585–5598,
975 Online. Association for Computational Linguistics.

A Prompts for Baselines

Following the instruction-tuning framework of Taori et al. (2023), we design prompt templates for keyword extraction and summarization baselines, as detailed in Prompt A1 and Prompt A2.

Prompt A1: Keywords Extraction

```
[INST] «SYS»
Extract a few keywords from the
following content.
«/SYS»
Prompt = """Below is an instruction that
describes a task, paired with an input that
provides content.
### Instruction: {"" + Instruction +
"""}
### Input: {"" + D + """}
### Response: ""
[/INST]
```

Prompt A2: Summary Generation

```
[INST] «SYS»
Generate a short summary of the
following content.
«/SYS»
Prompt = """Below is an instruction that
describes a task, paired with an input that
provides content.
### Instruction: {"" + Instruction +
"""}
### Input: {"" + D + """}
### Response: ""
[/INST]
```

Based on the aforementioned observation, we set $\alpha = 0.3$ in our method, which represents the optimal trade-off, maximizing the discriminatory power of retained concepts while maintaining a compact and informative context for downstream inference. This tuning contributes significantly to the robustness and effectiveness of our context compression approach in context engineering.

B Accuracy Details

C Ablation Study

We perform an ablation study to analyze the impact of the hyper-parameter α , which controls the significance threshold in the concept pruning process, on the overall performance of our method and the results are shown in Table A3. This parameter determines which concepts are retained from the AMR graphs based on their entropy values to construct the compressed context. Table A3 shows that the lower values of α overly restrict the retained information, pruning out useful concepts and leading to degraded performance. Conversely, higher α values retain too many concepts, which may introduce noise and reduce compression efficiency.

Table A1: Accuracy (Acc \uparrow) comparison on the PopQA dataset. The best results for each LLM with setting K are in **bold**, and the next best results are in underlined. Δ here represents the difference between Ours and Vanilla, and the **increased** and **decreased** Δ are marked differently. The best results for each of K are **marked**.

LLMs	$C \setminus K$	1	2	3	4	5	6	7	8	9	10
		Vanilla	48.57	64.77	54.60	54.07	63.13	60.78	62.42	63.23	69.63
GPT-Neo-1.3B	TF-IDF	22.00	37.89	38.21	38.95	43.90	39.87	42.28	39.71	43.70	59.40
	Keywords	30.00	41.61	43.68	30.58	33.75	30.98	48.32	47.10	43.70	37.60
	Summary	23.93	41.61	47.13	48.84	52.50	50.33	46.98	48.39	56.30	58.40
	SelCon	39.29	54.03	42.53	48.26	51.88	55.56	46.98	49.03	52.59	65.60
	LLMLingua	53.93	53.69	56.90	58.72	64.38	60.78	63.76	65.81	65.93	76.00
	Ours	53.93	68.45	57.47	59.88	70.00	68.63	69.13	71.61	68.89	79.20
	$\alpha = 0.01$	19.64	32.21	42.53	48.84	53.13	54.90	65.10	67.10	62.96	76.80
	$\alpha = 0.05$	28.57	41.28	48.28	56.98	58.75	60.78	67.11	68.39	63.70	77.60
	$\alpha = 0.1$	28.57	41.28	48.28	56.98	58.75	60.78	67.11	68.39	63.70	77.60
	$\alpha = 0.5$	35.00	51.01	56.90	58.14	61.88	69.28	69.13	62.58	65.19	76.80
	Δ	+5.36	+3.68	+2.87	+5.81	+6.87	+7.85	+6.71	+8.38	-0.74	+6.40
	GPT-Neo-2.7B	Vanilla	51.07	69.46	59.77	52.91	59.38	63.40	59.73	57.42	65.19
TF-IDF		33.57	46.31	50.00	53.49	58.75	64.05	59.06	61.29	60.74	76.00
Keywords		30.71	43.60	48.28	31.16	33.88	33.88	33.88	47.10	51.88	59.20
Summary		22.86	39.93	50.00	50.00	56.25	56.21	52.35	57.42	55.56	60.80
SelCon		45.00	56.71	50.00	46.51	59.38	54.25	57.72	54.84	53.33	70.40
LLMLingua		52.14	70.13	59.20	50.54	59.38	59.48	63.76	61.29	63.70	79.20
Ours		51.07	70.45	60.34	61.63	64.38	66.01	75.17	69.68	77.03	82.40
$\alpha = 0.01$		20.71	32.21	40.23	55.23	53.13	60.12	59.06	64.52	62.96	77.60
$\alpha = 0.05$		29.64	41.61	53.45	56.98	60.63	63.40	68.46	64.52	65.93	76.00
$\alpha = 0.1$		30.36	47.99	50.57	59.30	62.50	64.05	72.48	71.61	71.11	81.60
$\alpha = 0.5$		40.36	53.02	61.49	64.53	64.38	66.01	72.48	67.74	69.63	81.60
Δ		0.00	+0.99	+0.57	+8.72	+5.00	+2.61	+15.44	+12.26	+11.84	+6.40
OPT-1.3b	Vanilla	52.14	67.11	63.22	57.56	61.25	62.09	69.13	71.62	66.67	80.80
	TF-IDF	36.79	47.32	46.55	47.09	53.75	58.17	59.06	60.00	61.48	68.80
	Keywords	31.79	45.30	55.17	35.23	64.38	64.05	65.10	68.39	60.74	76.80
	Summary	26.79	47.65	56.90	59.30	65.00	61.44	65.10	67.74	65.19	77.60
	SelCon	49.29	61.07	56.90	56.98	63.75	60.13	67.79	73.55	74.07	82.40
	LLMLingua	55.00	71.14	61.49	57.56	65.00	64.71	75.17	73.54	68.89	84.80
	Ours	54.29	69.13	68.39	59.30	68.13	68.62	74.50	74.19	73.33	84.80
	$\alpha = 0.01$	23.57	34.56	44.25	48.84	58.13	60.13	69.80	73.55	70.37	84.80
	$\alpha = 0.05$	30.36	44.30	55.17	59.88	60.00	62.09	73.83	73.55	70.37	82.40
	$\alpha = 0.1$	32.86	46.98	60.92	62.79	66.25	69.93	73.15	75.48	75.56	86.40
	$\alpha = 0.5$	42.14	57.72	60.34	59.88	65.63	68.28	75.17	77.42	71.85	84.00
	Δ	+2.15	+2.02	+5.17	+1.74	+6.88	+6.53	+5.37	+2.57	+6.66	+4.00
OPT-2.7b	Vanilla	49.64	66.78	62.64	63.72	65.00	61.44	64.43	70.32	75.56	83.20
	TF-IDF	33.21	48.32	49.43	51.16	56.88	64.71	64.43	70.32	65.19	73.60
	Keywords	35.36	43.96	58.05	58.14	65.00	59.48	66.44	70.97	68.89	76.80
	Summary	29.64	49.33	54.60	55.23	60.00	54.90	60.40	65.16	56.30	67.20
	SelCon	48.21	64.09	54.02	58.14	66.25	60.78	67.11	74.19	73.33	79.20
	LLMLingua	55.71	73.15	62.64	62.79	71.25	65.36	63.09	77.42	71.11	84.80
	Ours	55.36	70.13	64.37	70.35	72.50	69.93	76.51	78.06	77.78	83.20
	$\alpha = 0.01$	22.86	35.91	45.98	54.65	61.88	62.75	66.44	72.26	68.15	83.20
	$\alpha = 0.05$	33.57	45.30	59.77	61.05	66.25	66.67	73.15	78.06	77.78	80.80
	$\alpha = 0.1$	35.00	53.02	60.34	66.86	69.38	67.32	75.17	76.13	79.26	80.00
	$\alpha = 0.5$	45.36	63.42	62.64	68.02	66.25	73.20	73.83	71.61	77.78	80.00
	Δ	+5.72	+3.35	+1.73	+0.58	+7.50	+6.49	+12.08	+7.74	+2.39	0.00
Bloom-560m	Vanilla	51.07	62.42	54.02	56.60	61.25	62.75	66.44	72.90	73.33	80.00
	TF-IDF	27.14	34.90	36.78	43.60	48.75	45.10	57.72	52.26	52.59	64.80
	Keywords	26.43	44.30	56.32	48.26	55.63	56.21	62.42	63.23	60.74	75.20
	Summary	27.50	47.65	52.87	54.07	61.88	58.17	67.11	70.97	62.22	77.60
	SelCon	34.29	49.66	45.40	43.60	47.50	47.06	51.68	59.35	48.89	65.60
	LLMLingua	53.57	65.10	52.87	52.33	60.00	59.48	68.46	74.84	67.41	80.80
	Ours	52.86	66.44	55.74	55.23	59.38	64.05	71.81	76.77	73.33	77.60
	$\alpha = 0.01$	18.57	26.51	33.91	36.63	46.25	50.33	61.07	60.65	66.67	72.80
	$\alpha = 0.05$	26.07	39.93	41.95	49.42	51.88	54.25	70.47	61.29	68.89	71.20
	$\alpha = 0.1$	30.00	40.94	46.55	47.09	50.63	61.44	71.14	65.16	71.85	76.80
	$\alpha = 0.5$	33.21	44.63	54.02	53.49	63.13	62.75	71.14	69.03	63.70	72.00
	Δ	+1.79	+4.02	+1.72	+0.77	+1.87	+1.30	+5.37	+3.87	0.00	+4.40
Bloom-7bl	Vanilla	56.43	73.49	72.41	65.12	68.75	74.12	78.52	80.65	77.04	82.00
	TF-IDF	40.36	56.04	62.64	61.04	65.63	71.24	76.51	76.13	77.04	84.80
	Keywords	38.93	53.02	62.64	59.88	65.63	67.32	75.84	74.19	71.85	77.60
	Summary	32.50	49.66	60.34	56.40	64.38	71.90	74.50	71.61	74.07	77.60
	SelCon	53.21	66.11	67.24	65.70	65.00	71.90	74.50	78.71	77.04	83.20
	LLMLingua	57.86	72.82	72.99	67.44	68.75	74.50	75.84	81.94	78.52	88.00
	Ours	54.64	69.80	73.56	65.12	71.25	77.12	82.55	83.23	82.22	91.20
	$\alpha = 0.01$	22.14	35.57	43.68	50.58	64.38	61.44	71.14	70.97	71.85	82.40
	$\alpha = 0.05$	31.07	48.66	59.77	59.88	68.13	68.63	73.83	74.84	80.74	86.40
	$\alpha = 0.1$	36.43	51.68	59.20	59.88	70.63	71.90	75.17	77.42	82.96	89.60
	$\alpha = 0.5$	42.86	61.74	66.67	62.79	73.75	75.82	79.19	82.58	81.48	85.60
	Δ	+1.79	+6.69	+1.15	0.00	+2.50	0.00	+4.03	+7.58	+5.18	+4.00
Llama-2-chat-13b	Vanilla	51.78	60.40	56.32	61.94	59.38	64.24	69.80	74.84	79.26	84.80
	TF-IDF	48.57	59.01	63.79	62.11	65.63	73.20	77.18	78.71	77.78	82.40
	Keywords	36.43	52.01	55.17	58.72	58.13	62.75	68.46	72.26	69.63	74.40
	Summary	27.86	44.63	46.55	54.07	48.13	46.41	56.38	65.16	65.93	83.20
	SelCon	55.00	69.13	63.79	67.44	65.00	69.28	73.15	78.71	80.00	86.40
	LLMLingua	58.93	68.79	63.22	71.51	65.63	68.63	72.48	80.00	81.48	88.00
	Ours	59.64	69.46	69.54	72.67	73.75	73.20	81.21	82.58	81.48	89.60
	$\alpha = 0.01$	30.00	41.61	44.83	54.65	56.25	66.01	71.14	68.39	73.33	85.60
	$\alpha = 0.05$	36.79	53.69	55.74	56.40	62.50	64.71	73.83	72.26	76.30	83.20
	$\alpha = 0.1$	43.93	61.41	62.64	63.37	65.63	70.59	75.17	75.48	77.03	84.80
	$\alpha = 0.5$	55.00	68.46	68.97	69.19	70.00	77.12	80.54	78.06	77.04	82.40
	Δ	+7.86	+9.06	+13.22	+11.63	+14.37	+18.96	+11.31	+7.04	+2.22	+4.80
Llama-3.1-8B-Instruct	Vanilla	61.43	76.17	74.41	74.42	70.63	70.74	83.22	87.38	86.67	89.40
	TF-IDF	51.07	66.18	64.37	70.93	70.63	73.20	78.52	81.94	81.48	75.20
	Keywords	51.79	56.04	62.07	63.37	61.88	69.28	73.83	76.13	77.04	81.60
	Summary	38.93	56.38	55.75	65.12	50.63	67.32	68.46	67.10	65.19	72.00
	SelCon	68.21	76.85	76.44	76.16	75.63	79.08	84.56	87.10	82.22	91.20</

Table A2: Accuracy (Acc \uparrow) comparison the EntityQuestions dataset. The symbols' definitions are same as Table A1.

LLMs	C / K	K									
		1	2	3	4	5	6	7	8	9	10
GPT-Neo-1.3B	Vanilla	47.24	60.31	58.45	56.95	60.25	62.31	60.34	65.09	66.92	71.68
	TF-IDF	21.27	28.32	32.71	29.15	35.15	40.20	33.52	37.28	36.15	38.94
	Keywords	22.50	35.66	37.27	36.61	43.10	46.73	46.93	44.38	33.85	45.13
	Summary	22.29	34.09	36.46	35.93	41.84	32.16	36.87	40.83	43.85	47.49
	SelCon	21.06	25.70	29.76	29.15	31.80	29.65	30.17	35.51	40.77	30.09
	LLMLingua	51.53	64.16	60.86	56.95	56.90	65.83	62.01	58.58	57.69	66.37
	Ours	50.92	61.36	59.79	57.29	64.85	57.79	60.35	63.31	63.08	66.37
	$\alpha = 0.01$	19.02	30.59	40.48	42.37	40.17	44.72	55.31	53.25	53.08	58.41
	$\alpha = 0.05$	25.97	39.69	44.77	49.49	52.72	53.27	59.78	63.91	55.38	59.29
	$\alpha = 0.1$	28.63	46.33	50.94	53.56	58.58	59.30	63.69	66.27	58.46	60.18
	$\alpha = 0.5$	37.01	51.05	54.69	57.63	55.23	56.78	55.87	56.21	56.15	59.29
	Δ	+3.68	+1.05	+1.34	+0.34	+4.60	-4.52	+0.01	-1.78	-3.84	-5.31
	GPT-Neo-2.7B	Vanilla	54.40	64.86	65.42	64.75	67.78	69.35	71.51	68.64	72.31
TF-IDF		30.88	33.39	49.33	43.73	51.04	55.28	60.34	57.99	60.00	66.37
Keywords		29.65	41.78	46.92	48.14	49.79	56.28	55.87	59.17	58.46	54.87
Summary		21.06	35.14	39.14	35.93	42.26	47.74	39.66	44.38	50.00	44.25
SelCon		24.74	28.67	35.92	32.20	37.66	38.69	44.13	41.42	43.08	30.09
LLMLingua		54.19	65.56	63.81	62.71	70.71	66.33	70.53	68.44	67.69	69.91
Ours		54.21	62.94	69.71	66.78	70.71	71.36	74.86	71.60	73.85	76.99
$\alpha = 0.01$		20.45	33.22	48.53	48.47	50.21	50.75	69.27	62.13	60.77	69.91
$\alpha = 0.05$		30.67	44.76	57.64	60.00	58.16	62.81	68.72	71.60	63.85	72.57
$\alpha = 0.1$		36.20	52.27	59.59	62.71	63.60	63.82	72.07	71.01	67.69	68.14
$\alpha = 0.5$		48.46	60.14	68.10	68.14	69.04	70.35	71.51	73.96	71.54	73.45
Δ		-0.19	-1.92	+3.29	+2.03	+2.93	+2.01	+3.35	+2.96	+1.54	+3.54
OPT-1.3b		Vanilla	56.24	66.78	65.68	65.42	73.22	67.84	70.39	69.82	75.38
	TF-IDF	32.92	41.26	47.99	45.42	56.90	45.73	46.37	45.56	46.82	53.98
	Keywords	31.29	37.59	52.82	53.22	60.67	59.30	60.34	59.76	63.84	64.60
	Summary	27.40	40.03	46.65	50.17	51.46	53.77	54.75	50.89	61.54	55.75
	SelCon	27.61	31.64	37.53	37.29	42.26	38.69	43.58	49.11	46.15	38.05
	LLMLingua	54.81	66.26	66.49	57.29	68.20	64.82	73.18	71.01	63.08	69.91
	Ours	53.41	62.76	69.71	63.73	76.15	70.85	74.86	75.15	76.15	73.45
	$\alpha = 0.01$	21.47	35.66	49.06	51.86	50.63	56.78	63.69	62.72	63.85	61.95
	$\alpha = 0.05$	30.06	42.66	54.42	60.00	61.09	60.80	67.04	65.09	65.38	64.60
	$\alpha = 0.1$	34.97	51.57	61.13	59.66	64.44	64.32	69.16	67.46	64.60	64.60
	$\alpha = 0.5$	46.63	58.74	68.15	65.08	64.83	66.33	67.04	68.64	68.46	68.14
	Δ	-2.83	-4.02	+4.03	-1.69	+2.93	+3.01	+4.47	+5.33	+0.77	+2.65
	OPT-2.7b	Vanilla	57.46	70.80	72.12	71.53	76.99	78.39	77.65	79.29	82.31
TF-IDF		34.97	41.96	50.67	51.19	63.60	64.82	63.13	64.50	68.46	62.83
Keywords		33.95	43.01	52.82	59.32	65.69	62.81	69.83	67.46	70.77	73.45
Summary		27.20	42.13	48.79	52.20	53.56	49.25	49.16	52.07	52.31	48.67
SelCon		30.06	35.14	43.70	42.37	45.19	44.22	46.37	48.52	49.23	43.36
LLMLingua		57.87	70.80	75.34	74.24	74.90	73.37	69.83	68.64	67.69	80.53
Ours		56.62	71.50	74.80	76.61	76.57	78.39	71.51	71.01	72.31	82.30
$\alpha = 0.01$		25.15	36.19	47.99	56.27	55.23	60.80	68.16	68.64	72.31	69.91
$\alpha = 0.05$		34.76	45.45	57.10	66.10	64.02	70.85	69.27	71.60	74.62	72.57
$\alpha = 0.1$		39.47	54.72	58.71	66.10	69.87	69.35	73.18	71.01	74.62	77.88
$\alpha = 0.5$		53.37	65.21	69.44	74.24	74.48	71.36	69.27	72.78	72.31	74.34
Δ		-0.84	+0.70	+2.68	+5.08	-0.42	0.00	-6.14	-8.28	-10.00	+2.66
Bloom-560m		Vanilla	48.26	56.47	53.62	53.22	57.32	60.80	54.75	59.17	61.53
	TF-IDF	26.18	27.27	31.37	28.14	41.00	35.18	38.55	43.79	36.15	39.82
	Keywords	24.74	35.31	41.29	41.69	48.54	46.23	54.19	47.33	42.31	46.90
	Summary	21.68	34.97	43.70	40.00	45.19	50.75	51.40	46.75	46.92	51.33
	SelCon	16.77	19.23	25.20	20.68	25.10	30.15	34.08	36.69	34.62	34.51
	LLMLingua	44.38	54.90	56.03	48.47	59.41	64.82	50.84	57.40	52.31	60.18
	Ours	42.97	52.27	53.35	47.12	59.00	64.32	50.28	56.21	60.77	59.29
	$\alpha = 0.01$	17.59	23.25	30.56	31.19	35.56	37.69	42.46	46.15	37.69	55.75
	$\alpha = 0.05$	23.31	30.42	34.32	36.27	43.10	39.70	52.51	55.02	47.69	56.64
	$\alpha = 0.1$	29.24	36.19	39.14	39.66	47.70	43.22	53.63	55.03	48.96	61.06
	$\alpha = 0.5$	36.40	41.08	44.77	47.80	49.37	48.74	52.21	53.25	50.00	60.18
	Δ	-5.29	+4.20	-0.27	-6.10	+1.68	+3.52	-4.47	-2.96	-0.76	-2.66
	Bloom-7b1	Vanilla	58.28	74.65	74.26	76.61	79.91	82.41	75.98	84.62	83.08
TF-IDF		37.63	47.03	53.08	61.36	67.36	63.32	67.04	68.05	72.31	68.14
Keywords		34.56	50.17	56.57	62.71	69.04	67.34	68.72	73.96	69.23	74.33
Summary		28.63	40.73	49.33	51.86	55.23	64.82	57.54	65.09	66.15	66.37
SelCon		27.81	36.36	42.36	43.05	46.44	46.73	43.58	53.85	51.54	46.90
LLMLingua		52.15	71.85	71.58	70.84	82.01	80.40	72.63	81.66	75.38	76.99
Ours		51.12	71.50	71.31	72.88	83.26	81.91	73.74	82.25	83.84	84.96
$\alpha = 0.01$		23.72	34.44	46.38	49.49	48.12	55.78	62.57	66.27	67.69	71.68
$\alpha = 0.05$		31.29	44.76	54.96	58.64	60.25	63.82	70.95	76.92	71.54	74.34
$\alpha = 0.1$		36.81	53.67	60.45	67.78	72.75	72.75	72.07	78.70	75.58	75.58
$\alpha = 0.5$		50.31	60.31	62.47	70.51	74.48	77.39	74.86	79.88	72.31	72.57
Δ		-7.16	-3.15	-2.95	-3.73	+3.35	-0.50	-7.47	-3.37	+0.76	+4.42
Llama-2-chat-13b		Vanilla	54.40	71.69	71.05	73.22	79.08	76.38	74.30	72.78	71.54
	TF-IDF	49.28	55.24	69.17	70.51	84.10	78.39	80.45	81.66	80.77	82.30
	Keywords	39.47	53.85	60.86	64.41	67.36	74.37	73.18	80.47	79.23	81.42
	Summary	31.29	43.71	45.30	44.07	50.63	51.76	41.90	52.66	56.15	66.37
	SelCon	37.63	43.18	53.35	56.27	59.41	65.33	64.25	68.05	69.23	68.14
	LLMLingua	59.30	73.25	72.39	80.00	81.59	78.39	79.33	80.47	81.54	84.07
	Ours	64.83	76.40	79.36	80.34	84.52	84.42	85.47	86.39	86.15	86.72
	$\alpha = 0.01$	36.40	41.61	57.37	61.36	68.20	74.37	73.18	80.47	75.38	72.57
	$\alpha = 0.05$	45.40	55.77	68.90	71.53	80.33	77.39	81.56	83.43	79.23	80.53
	$\alpha = 0.1$	49.69	65.03	74.26	76.61	81.59	82.41	81.56	86.39	79.23	81.42
	$\alpha = 0.5$	65.64	76.40	77.75	82.03	84.10	82.91	81.01	81.07	81.54	81.42
	Δ	+10.43	+4.71	+8.31	+7.12	+5.44	+8.04	+11.17	+13.61	+14.61	+7.08
	Llama-3.1-8B-Instruct	Vanilla	66.67	81.29	84.18	82.03	82.85	83.42	86.03	85.21	85.38
TF-IDF		56.65	61.19	72.39	75.93	78.66	69.85	72.07	69.82	70.00	58.41
Keywords		55.62	63.81	71.31	71.53	76.15	80.40	86.03	81.66	76.15	78.76
Summary		39.26	48.08	54.42	55.59	60.25	61.81	55.87	60.36	60.77	69.03
SelCon		41.41	50.35	61.39	59.66	59.41	68.34	63.12	62.72	68.46	61.95
LLMLingua		74.44	83.31	86.60	93.86	87.45	89.9				

Table A3: The ablation study results of AUC \uparrow . The LLMs' order and symbol definitions are the same as Table 2.

Datasets	α	K	G-1.3	G-2.7	O-1.3	O-2.7	b-560	b-7bl	L-13	L3.1-8	DS-V2	Q3-32	
PopQA	Ours	I_s	600.62	611.43	625.14	648.91	587.98	677.77	678.51	756.44	648.90	356.55	
		I_l	283.54	296.09	298.73	308.92	292.74	332.16	326.67	357.74	318.06	191.09	
	0.01	I_s	474.99	476.62	513.82	521.05	427.70	521.88	534.01	646.48	430.18	275.57	
		I_l	261.01	255.40	286.18	279.83	249.96	285.88	288.66	339.13	231.95	151.76	
		ΔI_s	-125.63	-134.81	-111.32	-127.86	-160.28	-155.89	-144.50	-109.96	-218.72	-80.98	
	0.05	ΔI_l	-22.53	-40.69	-12.55	-29.09	-42.78	-46.28	-38.01	-18.61	-86.11	-39.33	
		I_s	518.36	527.80	555.57	585.21	486.72	593.21	575.42	692.99	444.91	328.01	
		I_l	268.39	268.61	290.00	302.73	263.38	306.92	296.35	344.75	228.82	190.62	
	0.1	ΔI_s	-82.26	-83.63	-69.57	-63.70	-101.26	-84.56	-103.09	-63.45	-203.99	-28.54	
		ΔI_l	-15.15	-27.48	-8.73	-6.19	-29.36	-25.24	-30.32	-12.99	-89.24	-0.47	
		I_s	518.36	555.59	590.69	604.98	508.20	611.86	615.68	721.41	484.82	330.48	
	0.5	I_l	268.39	288.02	302.36	<u>304.22</u>	<u>277.27</u>	316.30	305.38	<u>351.58</u>	249.50	182.36	
		ΔI_s	-82.26	-55.84	-34.45	-43.93	-79.78	-65.91	-62.83	-35.03	-164.08	-26.07	
		ΔI_l	-15.15	-8.07	+3.63	-4.70	-15.47	-15.86	-21.29	-6.16	-68.56	-8.73	
	EntityQuestions	Ours	I_s	546.46	627.41	632.79	662.16	494.45	688.73	738.82	813.86	652.14	406.00
			I_l	248.82	294.48	298.31	295.18	229.06	323.26	343.58	371.30	307.05	181.50
		0.01	I_s	398.68	468.53	475.96	513.12	321.22	478.44	586.43	693.08	490.18	319.13
			I_l	213.20	252.50	249.62	274.46	173.02	260.26	302.50	345.54	252.38	148.58
			ΔI_s	-147.78	-158.88	-156.83	-149.04	-173.23	-210.29	-152.39	-120.78	-161.96	-86.87
		0.05	ΔI_l	-35.62	-41.98	-48.69	-20.72	-56.04	-63.00	-41.08	-25.76	-54.67	-32.92
			I_s	461.64	539.16	523.81	572.68	379.61	554.66	661.10	743.58	497.84	348.16
			I_l	235.35	271.86	260.21	<u>287.21</u>	203.99	288.49	323.18	355.98	243.08	161.74
		0.1	ΔI_s	-84.82	-88.25	-108.98	-89.48	-114.84	-134.07	-77.72	-70.28	-154.30	-57.84
			ΔI_l	-13.47	-22.62	-38.10	-7.98	-25.07	-34.77	-20.40	-15.32	-63.97	-19.76
I_s			<u>501.54</u>	564.93	554.98	596.23	408.18	601.44	692.64	766.50	534.69	<u>364.75</u>	
0.5		I_l	<u>248.16</u>	276.75	268.54	292.42	209.26	299.94	<u>329.10</u>	<u>362.84</u>	263.05	161.30	
		ΔI_s	-44.92	-62.48	-77.81	-65.93	-86.27	-87.29	-46.18	-47.36	-117.45	-41.25	
		ΔI_l	-0.66	-17.73	-29.77	-2.76	-19.80	-23.32	-14.48	-8.46	-44.00	-20.20	
0.01		I_s	491.76	<u>613.74</u>	<u>581.67</u>	<u>632.94</u>	<u>435.51</u>	<u>633.65</u>	<u>720.34</u>	<u>798.94</u>	<u>592.58</u>	355.74	
		I_l	226.26	<u>288.91</u>	<u>271.38</u>	<u>287.21</u>	<u>209.92</u>	<u>302.03</u>	325.79	360.77	<u>292.97</u>	138.40	
		ΔI_s	-54.70	-13.67	-51.12	-29.22	-58.94	-55.08	-18.48	-14.92	-59.56	-50.26	
0.05		ΔI_l	-22.56	-5.57	-26.93	-7.97	-19.14	-21.23	-17.79	-10.53	-14.08	-43.10	