

# Efficient Lossless Text Compression with Large Language Models: Enhancing Cross- Lingual and Cross-Domain Applications

Anonymous ACL submission

## Abstract

In the era of information explosion, the rapid growth of multilingual and multi-domain textual data poses unprecedented challenges for efficient storage and transmission. Traditional lossless compression methods such as Huffman coding, LZ77, and zlib perform well in certain scenarios but often rely on fixed statistical rules. This limits their ability to capture deeper linguistic structures, especially in complex or domain-specific texts. To address these limitations, we propose two large language model-based lossless text compression methods: **DeepSeekZip** and **LlamaZip**, which respectively integrate DeepSeek-8B and Llama3-8B as predictive models with conventional zlib compression. By leveraging the models' capabilities in modeling complex language patterns, our approach significantly enhances compression performance. Extensive experiments across various languages and text domains demonstrate that DeepSeekZip and LlamaZip consistently achieve over 10% higher compression rates than zlib alone. Notably, DeepSeekZip performs better in Chinese text compression, while both models show comparable results in English. Furthermore, compression effectiveness varies across domains: news and medical texts are compressed more efficiently than legal and technical ones. This highlights the impact of structure, terminology, and contextual dependencies on compression outcomes.

## 1 Introduction

With the rapid advancement of the information age, the exponential growth of textual data has posed unprecedented challenges to storage and transmission efficiency. Developing effective text compression methods has thus become a pressing need, as it not only reduces storage costs but also significantly improves data transfer performance. Traditional compression algorithms, such as Huffman coding (Huffman, 1952) and LZ77 (Ziv and Lempel, 1977),

have performed well in some settings, but their reliance on fixed statistical patterns and rules restricts their ability to capture complex semantic structures. This is especially evident in domain-specific or multilingual texts, where conventional techniques often yield lower compression rates.

Classical lossless text compression methods are generally classified into three categories. (1) Dictionary-based approaches, such as LZ77 and LZW, replace repetitive substrings to improve compression ratio but fail to capture long-range dependencies. (2) Statistical coding methods, such as Huffman and arithmetic coding, model character-level frequency but lack the ability to adapt to complex contextual patterns. (3) Grammar-based methods attempt to infer context-free grammar rules for structural compression. However, the grammar-based methods suffer from high computational complexity and the absence of efficient random access (Shannon, 1948), making them less practical in real-world applications with dynamic or multilingual data.

To overcome the representational limitations of traditional methods, recent work has explored the use of deep neural networks to capture complex contextual structures. Goyal et al. (2018) proposed DeepZip, which leverages RNN-based conditional probability modeling combined with arithmetic coding. Transformer-based architectures have further improved the modeling of long-range dependencies. RWKV (Peng et al., 2023) introduces linear attention mechanisms to retain contextual information while reducing inference overhead. Similarly Perceiver (Jaegle et al., 2021) extends input scalability via iterative attention.

Building on this trend, researchers have integrated pre-trained language models into compression pipelines. Li et al. (2021) proposed two Transformer-based strategies for enhanced text compression: Explicit Text Compression (ETC) and Implicit Text Compression (ITC). ETC em-

085 ploys a separate sequence-to-sequence compression model with attention mechanisms to extract  
086 key semantic components, which are concatenated with the original input to enrich the encoding  
087 process. In contrast, ITC integrates the compression module directly into a non-autoregressive decoder,  
088 enabling end-to-end training and supporting seamless integration with downstream tasks such as machine  
089 translation. Both approaches effectively reduce input redundancy while maintaining task performance  
090 and demonstrate strong transferability across various NLP tasks. However, challenges persist. ETC  
091 necessitates a considerable amount of supervised data to define "key semantics," resulting in significant  
092 annotation costs. Although ITC offers greater flexibility, it may neglect low-frequency yet critical  
093 information, which could lead to semantic loss. Furthermore, ETC-generated compressed sequences  
094 often demonstrate low interpretability, and ITC's architecture demands extensive customization  
095 to integrate with existing models. While these methods show promise for semantic compression,  
096 they still fall short of achieving robust lossless compression performance and broad applicability.

097 Recent efforts have also explored leveraging BERT for lossless compression. Öztürk and Mesut  
098 (2024) proposed MLMCompress, which utilizes BERT's bidirectional contextual prediction to estimate  
099 token distributions and integrates it with arithmetic coding. The model achieves up to 38% higher  
100 compression on English datasets compared to NNCP and a 42% improvement in multilingual tasks.  
101 Furthermore, MLMCompress operates up to 35 times faster than GPTZip (Nishi et al., 2023)  
102 in certain settings, with 20% faster compression and up to 180% faster decompression. This  
103 demonstrates the efficiency and practicality of contextual modeling for compression. However, the  
104 use of masked language modeling inherently imposes limitations. BERT demonstrates continuity in  
105 prediction and effectively models long-range semantics. Furthermore, MLMCompress lacks optimization  
106 for Chinese and other non-Latin languages.

107 As large language models (LLMs) continue to evolve, their exceptional language modeling capabilities  
108 have made them increasingly attractive for compression tasks. Pre-trained models such as LLaMA  
109 and GPT learn high-quality token distributions over massive corpora, enabling accurate next-token  
110 prediction for entropy coding. Valmeekam et al. (2023) introduced LLMZip, which combines  
111 LLaMA with arithmetic coding and achieves state-

112 of-the-art performance on English text. However, most LLM-based methods to date primarily focus  
113 on English, and their performance on Chinese or domain-specific content remains underexplored.

114 To address these limitations, we propose a novel lossless text compression framework based on two  
115 competitive open-source LLMs: **LlamaZip** (based on LLaMA3-8B (AI, 2024b)) and **DeepSeekZip**  
116 (based on DeepSeek-8B (AI, 2024a)). Our method utilizes the LLM to predict token-level probabilities  
117 and subsequently applies standard zlib compression to the residual errors. This hybrid approach  
118 balances expressiveness and efficiency. **Figure 1** illustrates the overall architecture of our compression  
119 framework, highlighting the two-stage process from semantic modeling to entropy coding. We  
120 evaluate the proposed methods on both Chinese and English corpora, encompassing diverse domains  
121 such as news, law, medicine, and technology. Experiments results demonstrate that our approach  
122 outperforms traditional zlib-based compression by over 10% on average, with DeepSeekZip achieving  
123 superior performance on Chinese text, and both models perform comparably on English datasets.

## 124 Our contributions are threefold: 125

- 126 (1) We utilize a lossless text compression framework that integrates large language models  
127 with traditional entropy coding, thereby enhancing compression rate and generalizability. 128
- 129 (2) We conduct a comprehensive analysis of model performance across various languages and  
130 domains, emphasizing differences in adaptability and robustness. 131
- 132 (3) We demonstrate the practical benefits of LLM-based compression through extensive  
133 experiments, providing new insights into language-aware compression in multilingual contexts. 134

## 135 2 Related Work 136

137 In recent years, text compression has made significant strides, with researchers exploring diverse  
138 strategies to enhance compression efficiency and modeling capability. This section reviews five  
139 major methodological directions in the field. 140

### 141 2.1 Dictionary-based Compression. 142

143 Classic dictionary-based methods, such as LZ77 and LZ78 (Ziv and Lempel, 1977), replace repeated  
144

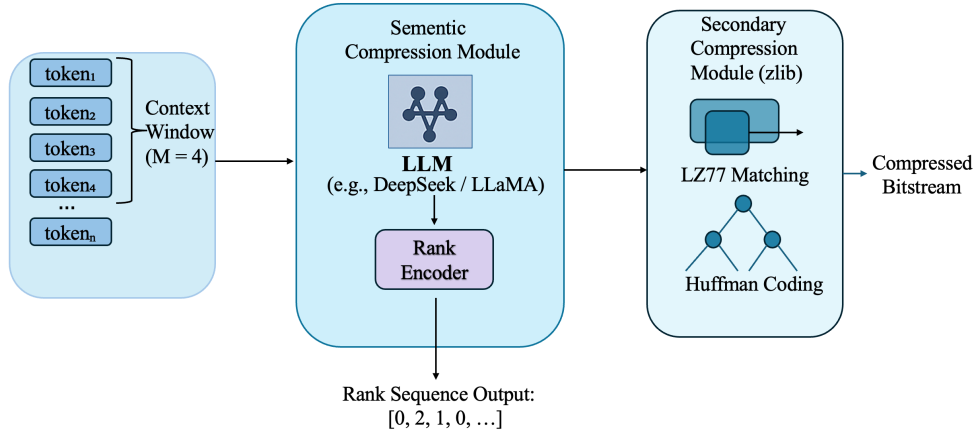


Figure 1: Semantic-to-Entropy Compression Pipeline. The LLM module generates token probability rankings, which are subsequently compressed using entropy-based methods such as zlib.

184 substrings with shorter references. More recently, 218  
 185 L3TC (Zhang et al., 2025) leverages a lightweight 219  
 186 Transformer variant (RWKV) combined with arith- 220  
 187 metic coding to achieve lower latency and com- 221  
 188 plexity. Chen et al. (Chen et al., 2003) further 222  
 189 proposed parallel dictionary lookup mechanisms 223  
 190 to enhance matching speed. However, dictionary- 224  
 191 based approaches face challenges in dynamic dic- 225  
 192 tionary construction, particularly in large-scale or 226  
 193 online compression settings where memory and  
 194 lookup overhead can become prohibitive. Multidic-  
 195 tionary systems can boost compression ratios but  
 196 significantly increase decoding complexity, which  
 197 limits their use in resource-constrained or real-time  
 198 devices.

## 199 2.2 Statistical Compression.

200 Statistical methods reduce redundancy by encod-  
 201 ing symbol frequencies. Huffman coding (Huff-  
 202 man, 1952) and arithmetic coding (Witten et al.,  
 203 1987) are foundational examples. Recently, k-th  
 204 order context models augmented with neural net-  
 205 works, such as RWKV (Peng et al., 2023), have im-  
 206 proved probabilistic estimation. Additionally, hy-  
 207 brid pipelines such as BWT+MTF+RLE (Adiego  
 208 et al., 2007) perform well in compressing struc-  
 209 tured data. However, high-order models require  
 210 large context buffers, resulting in increased mem-  
 211 ory usage and reduced adaptability to data streams  
 212 with changing distributions.

## 213 2.3 Grammar-based Compression.

214 Grammar-based methods aim to infer context-free  
 215 grammar rules that capture the underlying struc-  
 216 ture of text. Techniques based on smallest gram-  
 217 mar approximation and grammar induction un-

218 der higher-order entropy constraints have shown  
 219 promise (Gańczorz, 2018). However, these ap-  
 220 proaches suffer from computational intractability  
 221 (NP-hardness) and often relying on heuristics that  
 222 struggle with long sequences. Moreover, the re-  
 223 sulting compressed data structures lack efficient  
 224 random access, and their decoding speed is sub-  
 225 optimal, which limits their deployment in time-  
 226 sensitive scenarios.

## 227 2.4 Neural Arithmetic Coding.

228 Combining neural models with arithmetic coders  
 229 has recently emerged as a promising direction in  
 230 the field. L3TC (Zhang et al., 2025) employs  
 231 adaptive context modeling to dynamically adjust  
 232 token probabilities, approaching entropy bounds  
 233 while alleviating the bottlenecks of arithmetic cod-  
 234 ing through parallel block encoding. Nonetheless,  
 235 arithmetic coding requires high-precision float-  
 236 ing-point operations, rendering it unsuitable for low-  
 237 end devices or real-time compression tasks due to  
 238 its computational demands and sensitivity to nu-  
 239 merical instability.

## 240 2.5 LLM-based Compression.

241 Large language models (LLMs) such as LLaMA  
 242 and GPT have garnered significant attention in the  
 243 field of text compression due to their robust contex-  
 244 tual modeling capabilities. LLMZip (Valmeekam  
 245 et al., 2023) demonstrates the feasibility of com-  
 246 bining LLaMA3-8B with arithmetic coding to  
 247 achieve superior compression ratios on English cor-  
 248 pora. However, LLMZip suffers from extreme la-  
 249 tency—compressing just 10MB of text can take  
 250 up to 9.5 days—raising serious concerns about  
 251 its practicality. To address this challenge, Mittu

et al. (2024) introduced FineZip, a novel LLM-based compression framework incorporating on-line memory and dynamic context windows. By leveraging parameter-efficient fine-tuning (PEFT), FineZip reduces compression time from 9.5 days to 4 hours—a 54-fold improvement—while maintaining comparable compression ratios. It also surpasses traditional algorithms by achieving approximately 50% better compression rates on benchmark datasets. This work demonstrates that while LLM-based compression remains computationally intensive, its performance bottlenecks can be mitigated through architectural and optimization-level innovations.

## 2.6 Motivation and Gap.

Building upon these insights, we propose two LLM-based compression systems: **LlamaZip** and **DeepSeekZip**, which are based on LLaMA3-8B and DeepSeek-8B, respectively. Unlike prior works that primarily focus on English, we conduct a comprehensive evaluation across both English and Chinese corpora, covering multiple domains including news, medicine, law, and technology.

## 3 Method

We propose a two-stage compression framework that integrates the semantic modeling capabilities of large language models (LLMs) with traditional lossless compression techniques. Our approach comprises of a **Semantic Compression Module** and a **Secondary Compression Module**. This section provides a detailed overview of both modules and their mathematical formulations.

### 3.1 Semantic Compression Module

This module leverages LLMs to estimate the conditional probability distribution of each token based on the context, subsequently encoding the rank of the true token. The objective is to utilize the LLM’s language understanding to convert raw tokens into a more compressible rank sequence.

#### 3.1.1 Context Modeling and Temperature Scaling

Given a token context  $X = [x_{t-M}, \dots, x_{t-1}]$ , the LLM predicts the next token probability distribution:

$$P(x_t | X) = \text{LLM}(X) \quad (1)$$

To control sampling diversity, a temperature parameter  $T$  is introduced to rescale the output distribution:

When  $T \rightarrow 0$ , the model becomes deterministic:

$$Q_i = \frac{P_i^{1/T}}{\sum_j P_j^{1/T}} \quad (2)$$

$$Q_i = \begin{cases} 1 & \text{if } P_i = \max(P) \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

We set  $T = 0$  in all experiments to ensure deterministic outputs and enable reproducible, lossless decompression.

#### 3.1.2 Rank-Based Encoding

Instead of storing tokens directly, we record the rank of each ground-truth token under the LLM’s predicted distribution. The rank is computed as:

$$R_t = \text{rank}(P(w | X)) - 1 \quad (4)$$

where  $w$  is the true token at position  $t$ , and **the ranking is zero-based**—the most probable token has rank 0.

The final encoded sequence becomes:

$$\{R_1, R_2, \dots, R_n\} \quad (5)$$

**Example:** Consider the sentence “Artificial intelligence is rapidly advancing,” tokenized as “Artificial,” “intelligence,” “is,” “rapidly,” “advancing.” Using a context window of size  $M = 4$ , the LLM receives the first four tokens and predicts the fifth. Suppose the predicted token probabilities are:

Candidate Token	Probability
“advancing”	0.50
“evolving”	0.30
“expanding”	0.15
“adapting”	0.05

Table 1: Predicted token distribution.

Here, “advancing” is the correct token, ranked first in terms of probability. **Using zero-based indexing, its final rank is 0.** Thus, this token is encoded as 0. Repeating this procedure over the entire sequence results in a compact, rank-encoded representation.

This design utilizes the LLM’s comprehension of token dependencies to semantically compress text while maintaining reversibility.

## 3.2 Secondary Compression Module

Although the rank sequence is more compact than the original text, further compression is achievable through entropy coding. We utilize the Zlib library, which implements the DEFLATE algorithm by combining LZ77 and Huffman coding.

### 3.2.1 LZ77 Encoding

LZ77 minimizes redundancy by identifying repeated substrings within a sliding window. It encodes each match as a tuple  $\langle \text{length}, \text{distance} \rangle$  or outputs a literal character when no match is found.

Let  $W$  represent the sliding window size,  $S$  denote the search buffer, and  $p$  the current position. The match length  $L$  and distance  $D$  satisfy:

$$L = \max \{ l \mid S[p-l : p] = S[p-D-l : p-D] \}, \\ D \leq W \quad (6)$$

### 3.2.2 Huffman Coding

Huffman coding assigns shorter bit sequences to more frequent symbols. Let  $p(s)$  be the probability of symbol  $s$ , and  $l(s)$  its code length. Then:

$$H(S) \leq \mathbb{E}[l(s)] \leq H(S) + 1 \quad (7)$$

where the entropy  $H(S)$  is defined as:

$$H(S) = - \sum_s p(s) \log_2 p(s) \quad (8)$$

### 3.2.3 Bit Cost Estimation

To estimate the final compressed size, we compute the total number of bits in the output bitstream. After the rank sequences are passed through Huffman encoding, the bit length of each symbol  $s$  is denoted by  $l(s)$ , as introduced in Equation 7. The overall bit cost of the Zlib-compressed stream can be expressed as:

$$B_{\text{zlib}} = B_{\text{meta}} + \sum_{s \in \text{stream}} l(s) \quad (9)$$

Here,  $B_{\text{meta}}$  accounts for the overhead introduced by Zlib, including format headers, Huffman tree descriptors, and other control structures. The summation term represents the total number of bits consumed by symbol-level encoding in the stream.

Together, the semantic module reduces the entropy of the input by utilizing high-level language structures, while the Zlib module further minimizes bit-length through statistical coding. This hybrid approach achieves strong compression ratios while preserving exact reconstructability.

## 4 Experiments

### 4.1 Datasets

We use the text8 corpus from <http://mattmahoney.net/dc/text8.zip> as our base dataset. The text8 corpus is a cleaned, compressed excerpt of Wikipedia articles, which is widely employed in benchmarking text compression. To evaluate multilingual and domain-specific performance, we categorize text8 into four segments—**Medical**, **Legal**, **News**, and **Technical**—and use GPT-4o to generate corresponding Chinese translations. Note that we focus on small-scale data, with each text sample limited to around 10 KB, to examine the effectiveness of lightweight compression strategies.

The four domain subsets are described as follows:

- (1) **Medical Text:** Includes medical literature, case reports, and clinical descriptions, rich in domain-specific terminology and syntactic complexity.
- (2) **Legal Text:** Contains excerpts from legal codes and court decisions, featuring highly formal and nested grammatical structures.
- (3) **News Text:** Covers journalistic commentary and investigative reports, characterized by journalistic vocabulary and structured phrasing.
- (4) **Technical Text:** Consists of manuals and technical documentation, containing numerous domain-specific and instruction-oriented expressions.

All datasets undergo identical preprocessing steps, including text cleaning, tokenization, and normalization, to ensure fairness and consistency across all compression evaluations.

### 4.2 Experimental Setup

The experiments proceed in the following stages (see Figure 2):

- (1) **Segmentation and Tokenization:** Each document is split into chunks suitable for LLM input, followed by tokenization using each model’s native tokenizer.
- (2) **Prediction and Compression:** DeepSeek or LLaMA models are used to compute

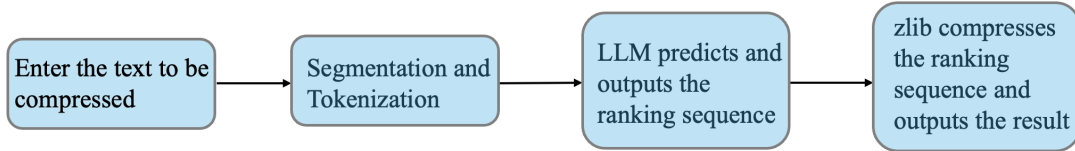


Figure 2: The experimental workflow of our compression framework.

next-token probability distributions and ranking sequences. These sequences are then compressed using `zlib` to produce final bit-streams.

- (3) **Evaluation:** Compression efficiency is measured using the *space-saving rate*  $\eta$ , computed as:

$$\eta = \frac{\text{Original Size} - \text{Compressed Size}}{\text{Original Size}} \times 100\%. \quad (10)$$

We conduct ablation experiments by varying the memory length parameter  $M \in \{128, 256, 512, 1024\}$  for each LLM to study the effect of context window size on compression rates. For each domain and language (Chinese and English), we select 5 representative 10 KB samples from `text8`, compress them using both `DeepSeek` and `LLaMA` models, and report the averaged results.

**Hardware and Software Environment.** All experiments are conducted on identical GPU servers. The compression pipeline is implemented in Python using `PyTorch` for model loading and inference.

### 4.3 Experimental Results

We conducted comparative experiments using both `LLaMA` and `DeepSeek` models. For each, we selected 5 Chinese and 5 English samples per domain (medical, legal, news, and technical) from the `text8` dataset. Each sample was approximately 10 KB in size. We then compressed these samples using each language model and computed the average space-saving rate. Additionally, we performed ablation studies by varying the memory length parameter  $M \in \{128, 256, 512, 1024, 2048\}$ , to evaluate the impact of context length on compression performance across different domains.

#### 4.3.1 Experimental Variables

- (1) **Text Domains and Language Types**

Model	M = 128	M = 256	M = 512	M = 1024
DeepSeekZip	42.09%	49.19%	56.19%	56.19%
LlamaZip	40.91%	48.72%	55.26%	55.26%
zlib	46.78%	46.78%	46.78%	46.78%

Table 2: Average space-saving rate (%) for different memory lengths  $M$  across all test samples.

We selected samples from four domains in the `text8` dataset. For each domain, we extracted **both Chinese and English** text samples of equal size to enable cross-linguistic comparison. Within each domain-language pair, we selected 5 thematically consistent samples and reported the average space-saving rate as the representative result.

- (2) **Comparison Models and Baseline**

We applied two large language models: **DeepSeekZip** and **LlamaZip**. Additionally, we used the traditional `zlib` compression algorithm as a baseline to compare against the LLM-based approaches.

- (3) **Main Variable: Memory Length  $M$**

We conducted experiments with memory lengths  $M = \{128, 256, 512, 1024, 2048\}$ . The objective was to observe how compression performance changes with increasing memory length, under a fixed input size constraint.

#### 4.3.2 Effect of Memory Length on space-saving rate

We compare the average space-saving rate  $\eta$  achieved by `DeepSeekZip`, `LlamaZip`, and the traditional `zlib` algorithm across all test samples. The results are summarized for memory lengths  $M = \{128, 256, 512, 1024, 2048\}$ . Since `zlib` is independent of model memory length, its space-saving rate remains constant across different settings.

As shown in Figure 3, the average space-saving rate of **DeepSeekZip** and **LlamaZip** increase significantly with the growth of memory length  $M$ ,

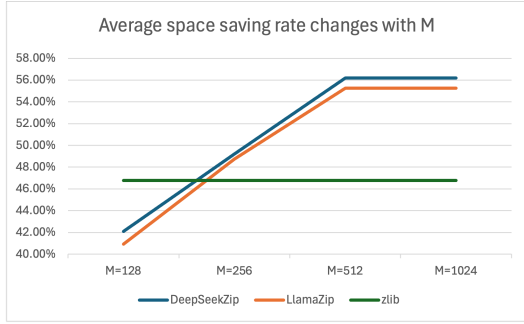


Figure 3: Average space-saving rate (%) as memory length  $M$  increases for DeepSeekZip, LlamaZip, and zlib.

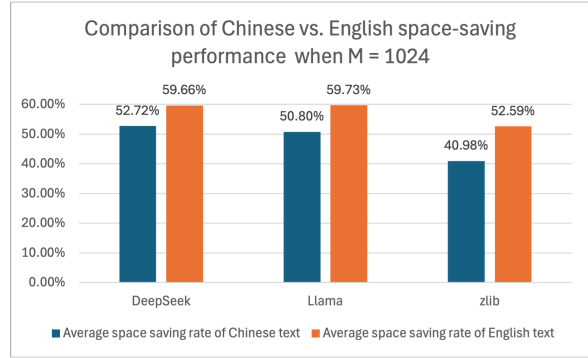


Figure 4: Comparison of Chinese vs. English space-saving performance when memory length  $M = 1024$ .

while the compression ratio of zlib remains constant. Detailed observations are as follows:

- (1) **From 128 to 256:** DeepSeekZip improves from 42.09% to 49.19%, and LlamaZip from 40.91% to 48.72%, while zlib remains at 46.78%. This demonstrates that longer memory significantly enhances LLMs’ contextual modeling, improving compression.
- (2) **From 256 to 512:** DeepSeekZip and LlamaZip continue to improve to 56.19% and 55.26%, respectively, while zlib stays flat. This highlights the value of long-range context in LLM-based compression.
- (3) **From 512 to 1024:** Space-saving rates plateau for both LLMs, with no further gains. This indicates diminishing returns beyond 512 tokens for fixed-length inputs (e.g., 10 KB).

**Overall Observations:** DeepSeekZip and LlamaZip consistently enhance space-saving rates as memory length grows, especially within  $M = 128 \sim 512$ , while zlib maintains constant due to its lack of contextual modeling. DeepSeekZip generally outperforms LlamaZip, though the gap narrows with larger  $M$ . Beyond  $M = 512$ , both models exhibit saturation, reflecting marginal gains.

#### Marginal Effect and Potential Causes:

- (1) **Emergence of Marginal Returns:** A significant improvement is observed when  $M$  increases from 128 or 256 to 512. However, the space-saving rate gains diminish beyond  $M = 512$  and become negligible up to  $M = 2048$ . This suggests that in moderate-length documents, once the model captures most of the redundant or structured information, further extending the memory yields limited additional benefits.

Model	Medical (zh)	Legal (zh)	Technical (zh)	News (zh)	Avg.
DeepSeekZip	56.03%	48.64%	48.36%	57.83%	52.72%
LlamaZip	51.83%	48.07%	46.88%	56.41%	50.80%
zlib	40.31%	38.76%	38.28%	46.55%	40.98%

Table 3: Space-saving rates of different models on Chinese texts in four domains ( $M = 1024$ ).

- (2) **Possible Causes:** There are two likely explanations. First, if the input document is relatively short or lacks recurring patterns, then increasing memory size may yield diminishing returns once redundancy is fully addressed. Second, Transformer-based models may experience decreasing efficiency in utilizing extended context, resulting in weaker marginal benefits.

#### 4.3.3 Comparison of Chinese vs. English Space-Saving Performance ( $M = 1024$ )

To evaluate the effectiveness of different models on Chinese and English texts, we compared the space-saving rate of DeepSeekZip, LlamaZip, and zlib under a fixed memory length  $M = 1024$ . The following tables present the average space-saving results across four domains (see Tables 3 and 4).

- (1) **Chinese Text Analysis** DeepSeekZip achieves the highest average space-saving rate for Chinese texts at 52.72%, outperforming LlamaZip (50.80%) and zlib (40.98%). Across all four Chinese domains, DeepSeekZip consistently outperforms the other methods, demonstrating superior effectiveness in addressing the unique characteristics of the Chinese language. In contrast, zlib yields the lowest space-saving rate and remains unaffected by memory length, highlighting its lack of semantic awareness and context modeling.

Model	Medical (en)	Legal (en)	Technical (en)	News (en)	Avg.
DeepSeekZip	59.67%	58.08%	60.63%	60.27%	59.66%
LlamaZip	59.85%	58.72%	60.45%	59.90%	59.73%
zlib	52.51%	51.64%	52.80%	53.41%	52.59%

Table 4: Space-saving rates of different models on **English** texts in four domains ( $M = 1024$ ).

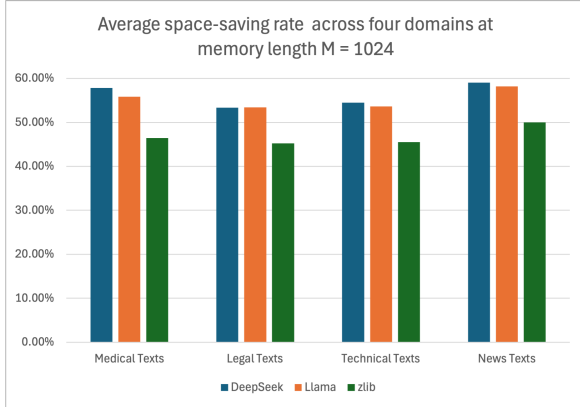


Figure 5: Performance differences of each model in different fields at memory length  $M = 1024$ .

**(2) English Text Analysis** DeepSeekZip and LlamaZip perform nearly identically on English texts, with average space-saving rates of 59.66% and 59.73% respectively. Both models significantly outperform zlib (52.59%), confirming the advantage of LLMs in modeling English semantics and capturing long-range dependencies.

**(3) Summary** DeepSeekZip shows a clear advantage over LlamaZip when processing Chinese texts, while both methods perform similarly on English texts. Both LLM-based approaches significantly outperform zlib in all instances, particularly with Chinese data. These results underscore the robustness of DeepSeekZip and LlamaZip across various languages and domains, showcasing notable improvements in space-saving efficiency.

#### 4.3.4 Domain-Specific Space-Saving Performance ( $M = 1024$ )

To compare the performance of **DeepSeekZip**, **LlamaZip**, and **zlib** across different domains. The average space-saving rate is computed for each model within each domain.

##### Analysis:

- Medical Texts:** DeepSeekZip achieves the highest space-saving rate at 57.85%, outperforming LlamaZip (55.84%) and zlib (46.41%). This indicates DeepSeekZip’s su-

Model	Medical	Legal	Technical	News
DeepSeekZip	57.85%	53.36%	54.50%	59.05%
LlamaZip	55.84%	53.40%	53.67%	58.16%
zlib	46.41%	45.20%	45.54%	49.98%

Table 5: Average space-saving rate (%) across four domains at memory length  $M = 1024$ .

perior capability in capturing medical expressions.

- Legal Texts:** DeepSeekZip and LlamaZip perform nearly equally with rates of 53.36% and 53.40%, respectively—both significantly better than zlib’s 45.20%. This suggests both LLM-based methods handle the formal and complex syntax of legal documents more effectively than traditional compression methods.
- Technical Texts:** DeepSeekZip slightly outperforms LlamaZip (54.50% vs. 53.67%), and both surpass zlib (45.54%), reflecting strong performance in compressing structured and terminology-rich technical content.
- News Texts:** DeepSeekZip shows superior performance with a 59.05% space-saving rate, ahead of LlamaZip (58.16%) and zlib (49.98%), highlighting its strength in modeling the variability and redundancy present in journalistic writing.

Overall, DeepSeekZip demonstrates consistent superiority across all domains, particularly in medical and news texts. LlamaZip follows closely and also significantly outperforms the traditional zlib approach, which demonstrates the weakest results across all categories. This reinforces the advantages of LLM-based methods in domain-sensitive semantic modeling for compression.

## 5 Conclusion

Text compression is evolving toward LLM-based systems that leverage deep semantics. We propose **DeepSeekZip** and **LlamaZip**, combining LLM prediction with entropy coding. Results show that longer memory lengths (e.g.,  $M = 512$ ) improve compression by capturing semantic redundancy. Our methods outperform traditional baselines while remaining lossless. Future work includes structure-aware modeling and efficient decoding like *FineZip* (Mittu et al., 2024).



## 621 Limitations

622 While our framework demonstrates strong perfor- 671  
623 mance across languages and domains, several lim- 672  
624 itations remain. First, LLM-based compression 673  
625 methods are computationally expensive compared 674  
626 to traditional algorithms like zlib, especially during 675  
627 inference. This restricts their deployment in real- 676  
628 time or resource-constrained environments. Sec- 677  
629 ond, although we evaluated multilingual settings, 678  
630 our non-English experiments relied on machine- 679  
631 translated corpora, which may not fully capture the 680  
632 complexities of native-language structures. Third, 681  
633 the current approach does not incorporate adap- 682  
634 tive memory or dynamic context resizing, which 683  
635 could be important for handling variable-length 684  
636 documents more efficiently. Lastly, our evalua- 685  
637 tion focuses on average space-saving rates, without 686  
638 measuring decoding latency or memory overhead, 687  
639 which are also critical in practical applications. Fu- 688  
640 ture work could address these issues through model 689  
641 optimization, native multilingual pretraining, and 690  
642 efficiency-aware benchmarks.

## 643 Ethics Statement

644 Our work builds upon publicly available large lan- 693  
645 guage models (DeepSeek and LLaMA) and stan- 694  
646 dard compression libraries (zlib). All datasets used, 695  
647 including the text8 corpus and its domain-specific 696  
648 subdivisions, are derived from open-access sources. 697  
649 To enable multilingual evaluation, we generate cor- 698  
650 responding Chinese texts from the English corpus 699  
651 using GPT-4o machine translation, followed by 700  
652 basic validation to ensure semantic consistency. 701  
653 These translations are used for research purposes 702  
654 only, and no personally identifiable or sensitive 703  
655 user data is involved. 704

656 We recognize that automatically translated data 705  
657 may not fully capture the linguistic richness and 706  
658 diversity of native Chinese corpora, which could 707  
659 limit generalizability and introduce subtle biases. 708  
660 Future work should explore more authentic and di- 709  
661 verse Chinese datasets to support broader language 710  
662 fairness. 711

663 Moreover, while LLM-based compression meth- 712  
664 ods offer significant gains in efficiency, they come 713  
665 with non-negligible environmental costs due to 714  
666 the computational demands of large-scale models. 715  
667 Our experiments are conducted using existing pre- 716  
668 trained models to minimize additional carbon foot- 717  
669 print. No human subjects were involved in this 718  
670 study, and no privacy or safety concerns arise from 719

our methodology. 671

We advocate for continued efforts toward sus- 672  
tainable, inclusive, and responsible NLP research. 673

## 674 References

- 675 Joaquín Adiego, Gonzalo Navarro, and Pablo de la 676  
677 Fuente. 2007. Using structural contexts to compress 678  
679 semistructured text collections. *Information Process- 680*  
681 *ing & Management*, 43(3):769–790. 682
- 683 DeepSeek AI. 2024a. Deepseek-r1-distill-llama- 684  
685 8b. [https://huggingface.co/deepseek-ai/ 686](https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B)  
687 [DeepSeek-R1-Distill-Llama-8B](https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B). Accessed: 688  
689 2024-05-19. 690
- 691 Meta AI. 2024b. Meta-llama-3-8b. [https:// 692](https://huggingface.co/meta-llama/Meta-Llama-3-8B)  
693 [huggingface.co/meta-llama/Meta-Llama-3-8B](https://huggingface.co/meta-llama/Meta-Llama-3-8B). 694  
695 Accessed: 2024-05-19. 696
- 697 David Chen, Enoch Peserico, and Larry Rudolph. 2003. 698  
699 A dynamically partitionable compressed cache. In 700  
701 *Proceedings of the Singapore-MIT Alliance Sympo- 702*  
703 *sium*. Technical Report, MIT Laboratory for Com- 704  
705 puter Science.
- 706 Michał Gańczorz. 2018. Entropy bounds for grammar 707  
708 compression. *arXiv preprint arXiv:1804.08547*. 709
- 710 Mohit Goyal, Kedar Tatwawadi, Shubham Chandak, and 711  
712 Idoia Ochoa. 2018. Deepzip: Lossless data compres- 713  
714 sion using recurrent neural networks. *arXiv preprint 715*  
716 *arXiv:1811.08162*. 717
- 718 David A Huffman. 1952. A method for the construction 719  
720 of minimum-redundancy codes. *Proceedings of the 721*  
722 *IRE*, 40(9):1098–1101.
- 723 Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol 724  
725 Vinyals, Andrew Zisserman, and Joao Carreira. 2021. 726  
727 Perceiver: General perception with iterative atten- 728  
729 tion. In *International conference on machine learn- 730*  
731 *ing*, pages 4651–4664. PMLR. 732
- 733 Zuchao Li, Zhuosheng Zhang, Hai Zhao, Rui Wang, 734  
735 Kehai Chen, Masao Utiyama, and Eiichiro Sumita. 736  
737 2021. Text compression-aided transformer encoding. 738  
739 *IEEE Transactions on Pattern Analysis and Machine 740*  
741 *Intelligence*, 44(7):3840–3857. 742
- 743 Fazal Mittu, Yihuan Bu, Akshat Gupta, Ashok De- 744  
745 vireddy, Alp Eren Ozdarendeli, Anant Singh, and 746  
747 Gopala Anumanchipalli. 2024. Finezip: Pushing the 748  
749 limits of large language models for practical lossless 750  
751 text compression. *arXiv preprint arXiv:2409.17141*. 752
- 753 Erika Nishi, Hayato Mizutani, and Takuya Hashimoto. 754  
755 2023. Gptzip: Lossless text compression using 756  
757 gpt. <https://github.com/erika-n/GPTzip>. Ac- 758  
759 cessed: 2025-05-08. 760
- 761 Emir Öztürk and Altan Mesut. 2024. Learning-based 762  
763 short text compression using bert models. *PeerJ 764*  
765 *Computer Science*, 10:e2423. 766

722 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak,  
723 Samuel Arcadinho, Stella Biderman, Huanqi Cao,  
724 Xin Cheng, Michael Chung, Matteo Grella, and 1  
725 others. 2023. Rwkv: Reinventing rns for the trans-  
726 former era. *arXiv preprint arXiv:2305.13048*.

727 Claude E Shannon. 1948. A mathematical theory of  
728 communication. *The Bell system technical journal*,  
729 27(3):379–423.

730 Chandra Shekhara Kaushik Valmeekam, Krishna  
731 Narayanan, Dileep Kalathil, Jean-Francois Chamber-  
732 land, and Srinivas Shakkottai. 2023. Llmzip: Loss-  
733 less text compression using large language models.  
734 *arXiv preprint arXiv:2306.04050*.

735 Ian H Witten, Radford M Neal, and John G Cleary. 1987.  
736 Arithmetic coding for data compression. *Communi-  
737 cations of the ACM*, 30(6):520–540.

738 Junxuan Zhang, Zhengxue Cheng, Yan Zhao, Shihao  
739 Wang, Dajiang Zhou, Guo Lu, and Li Song. 2025.  
740 L3tc: Leveraging rkwv for learned lossless low-  
741 complexity text compression. In *Proceedings of the  
742 AAI Conference on Artificial Intelligence*, pages  
743 13251–13259.

744 Jacob Ziv and Abraham Lempel. 1977. A universal  
745 algorithm for sequential data compression. *IEEE  
746 Transactions on information theory*, 23(3):337–343.

## 747 A Full Results and Compression Pipeline

### 748 A.1 Experimental Results by Domain and 749 Memory Length

750 Table 6 presents the full results of our compression  
751 experiments. We report the average space-saving  
752 rates (%) across four domains—Medical, Legal,  
753 Technical, and News—in both Chinese and English,  
754 using memory lengths  $M \in \{128, 256, 512, 1024\}$ .  
755 Each score is averaged over five samples per con-  
756 figuration.

### 757 A.2 Visualization

758 To visually complement the compression process,  
759 we include a diagram illustrating the semantic-to-  
760 entropy pipeline (Figure 1). The figure shows how  
761 the LLM module predicts token-level probability  
762 rankings, which are then transformed into rank  
763 sequences and further compressed using entropy-  
764 based algorithms such as `zlib`.

### 765 A.3 Experimental environment

766 We used DeepSeek-8B and LLaMA3-8B as frozen  
767 pre-trained models, each with approximately 8B  
768 parameters. Experiments were conducted on A100  
769 (80GB) GPUs. No additional fine-tuning or train-  
770 ing was performed.

## B Licenses of Used Artifacts

771 We summarize the licenses of all third-party arti-  
772 facts used in this work: 773

- 774 • **DeepSeek-R1-Distill-Llama-8B**: Released 774  
775 under the **Apache 2.0 License**, available at 775  
776 [https://huggingface.co/deepseek-ai/](https://huggingface.co/deepseek-ai/DeepSeek-R1-Distill-Llama-8B) 776  
777 `DeepSeek-R1-Distill-Llama-8B`. 777
- 778 • **Meta-Llama-3-8B**: Released under Meta’s 778  
779 custom license for research use, available 779  
780 at [https://huggingface.co/meta-llama/](https://huggingface.co/meta-llama/Meta-Llama-3-8B) 780  
781 `Meta-Llama-3-8B`. 781
- 782 • **zlib**: A classical entropy coding library re- 782  
783 leased under the **Zlib License**. 783

784 All licenses permit the use of these models and 784  
785 tools for non-commercial research purposes. No 785  
786 modifications were made to the original released 786  
787 artifacts. 787

Model	Medical		Legal		Technical		News		Avg.
	zh	en	zh	en	zh	en	zh	en	
<b>M = 128</b>									
DeepSeek	35.73	47.84	33.79	50.18	35.47	49.50	37.09	47.15	42.09
Llama	34.05	47.66	32.34	50.06	32.42	48.51	35.33	46.92	40.91
zlib	40.31	52.51	38.76	51.64	38.28	52.80	46.55	53.41	46.78
<b>M = 256</b>									
DeepSeek	47.56	54.11	41.04	53.86	41.09	54.10	47.12	54.67	49.19
Llama	44.35	54.17	43.18	54.27	39.77	54.42	45.47	54.15	48.72
zlib	40.31	52.51	38.76	51.64	38.28	52.80	46.55	53.41	46.78
<b>M = 512</b>									
DeepSeek	56.03	59.67	48.64	58.08	48.36	60.63	57.83	60.27	56.19
Llama	51.83	59.85	48.07	58.72	46.88	60.45	56.41	59.90	55.26
zlib	40.31	52.51	38.76	51.64	38.28	52.80	46.55	53.41	46.78
<b>M = 1024</b>									
DeepSeek	56.03	59.67	48.64	58.08	48.36	60.63	57.83	60.27	56.19
Llama	51.83	59.85	48.07	58.72	46.88	60.45	56.41	59.90	55.26
zlib	40.31	52.51	38.76	51.64	38.28	52.80	46.55	53.41	46.78

Table 6: Space-saving rates across domains, languages, and memory lengths.

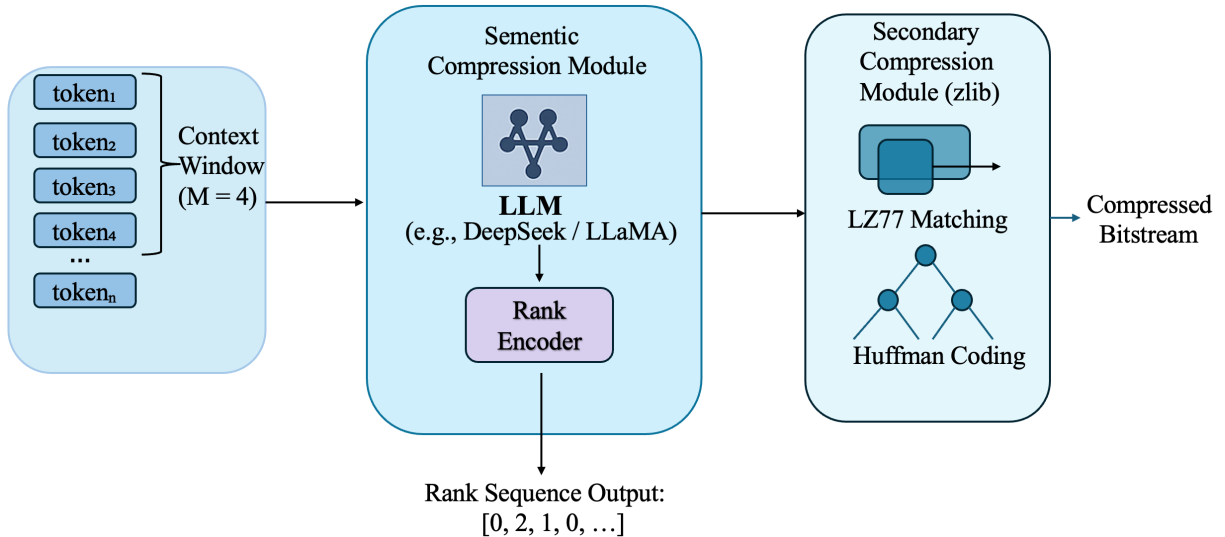


Figure 6: Semantic-to-Entropy Compression Pipeline. The LLM module outputs token rankings, which are subsequently encoded into compact sequences and compressed via entropy coding.