

# *ShortLang*: Compressed Text for efficient LLMs

Praneeth Vadlapati

*Independent researcher*

praneethv@arizona.edu

ORCID: 0009-0006-2592-2564

**Abstract:** The rapid increase in the size and complexity of modern language models has generated renewed interest in methods that can reduce computational requirements without compromising semantic fidelity. This paper introduces ShortLang, a minimal-length, semantically-preserving textual representation framework designed to optimize language model reasoning, training efficiency, and storage requirements. ShortLang compresses natural language into a concise symbolic form while retaining core meaning as measured by embedding similarity. It can serve as an intermediate representation for downstream tasks such as chunking, vector storage, retrieval-augmented generation, model training, and multi-step reasoning. This paper outlines key principles behind ShortLang, methods for generating and validating ShortLang representations, and its potential applications in large-scale machine learning systems. We further discuss strategies for building automatic ShortLang converters, including both rule-based systems and fine-tuned summarization models, and we examine theoretical and practical considerations for future research.

The source code is available at [github.com/Pro-GenAI/ShortLang](https://github.com/Pro-GenAI/ShortLang).

**Keywords:** Large Language Models, LLMs, Generative AI, Artificial Intelligence, AI

## I. INTRODUCTION

Recent advances in large language models (LLMs) have demonstrated remarkable capabilities in reasoning, summarization, planning, and knowledge retrieval. However, these systems often operate on lengthy and redundant natural-language inputs, resulting in high token usage, increased computational load, and elevated costs. As models scale to more complex tasks and larger corpora, the need for compact representations becomes increasingly urgent. Traditional summarization cannot always meet this need because it aims to preserve readability and narrative flow, which often leads to residual redundancy.

ShortLang proposes an alternative: a deliberately compressed textual encoding that prioritizes semantic sufficiency over linguistic naturalness. It transforms verbose natural language into a dense, symbolic, and minimally redundant form. Unlike classical shorthand or linguistic compression methods, ShortLang is designed specifically for machine consumption rather than human readability. Its primary goal is to retain all necessary information for machine inference with minimal token overhead.

### A. Disadvantages with current approaches

Existing techniques for optimizing language model performance often rely on summarization, token pruning, or heuristic chunking, yet these methods are limited by their dependence on natural-language structures. Traditional summarization aims to create human-readable condensations, which inevitably preserve syntactic redundancy and stylistic markers that are unnecessary for machine inference. Token pruning and compression algorithms, while effective for eliminating some forms of redundancy, frequently disrupt semantic coherence or introduce noise that degrades

downstream model accuracy. Furthermore, current chunking strategies operate on verbatim text, meaning that even trivial or semantically light components contribute to storage cost, embedding dimensionality, and retrieval inefficiency. As models scale, these limitations manifest in higher computational overhead, increased memory demands, and unnecessary expenditure in environments where tokenized operations are directly tied to financial cost.

### *B. Proposed system and its benefits*

The proposed system, ShortLang, introduces a structured methodology for transforming verbose input into a minimal-length, semantically preserved representation optimized for machine processing. Unlike standard summarizers, ShortLang does not attempt to maintain readability but instead focuses on retaining only the semantic core of text. Through a combination of linguistic compaction rules, abbreviation mapping, and model-driven semantic compression, ShortLang removes extraneous syntactic elements while preserving the conceptual information needed for reasoning. Its benefits are multifold: it significantly reduces token counts in both inference and training contexts; it increases the semantic density of stored knowledge units; and it improves computational throughput by allowing models to process meaning rather than linguistic scaffolding. By utilizing embedding-based similarity checks, the system can objectively measure semantic retention, thereby enabling more rigorous validation of compression quality than traditional human-readability metrics.

### *C. New use cases of the system*

ShortLang enables new capabilities that extend beyond the scope of traditional summarization and compression frameworks. Because it yields a semantically dense representation, ShortLang is well-suited for high-efficiency vector storage in large-scale retrieval systems, allowing organizations to store more conceptual information within limited database budgets. It also facilitates more efficient inter-agent communication in multi-agent LLM architectures, where concise semantic messages reduce latency and improve coordination. In training pipelines, ShortLang can compress large datasets into compact forms that retain the essential patterns models must learn, supporting low-resource training and fine-tuning within constrained computational environments. Moreover, ShortLang opens the possibility of creating domain-specific semantic codes for scientific literature, legal documents, and technical manuals, enabling rapid reasoning by models without processing the full linguistic complexity of the original sources. These new use cases highlight ShortLang’s potential to serve not only as a compression tool but also as a foundational representational format for machine cognition.

## *D. Applications of ShortLang*

### *1. Reasoning Optimization*

Reasoning chains often require models to process long passages. Converting these passages into ShortLang before feeding them into the model can reduce token usage and shift model capacity toward inference rather than redundant text processing.

### *2. Training Data Compression*

Massive corpora used for training LLMs contain substantial linguistic redundancy. ShortLang compression could reduce training data size while preserving the conceptual diversity necessary for effective model generalization. This can significantly lower training costs and reduce hardware requirements.

### *3. Efficient Chunking for Vector Embedding*

Before splitting text into chunks for embedding or retrieval-augmented generation, applying ShortLang compression reduces the number of tokens per chunk and may increase semantic density. This allows vector databases to store more information in fewer embeddings, improving retrieval performance and reducing storage requirements.

### *4. Vector Database Storage and Retrieval*

Vector databases often contain millions of stored text segments. By storing ShortLang representations instead of full natural-language text, organizations can reduce storage costs and accelerate nearest-neighbor search. Retrieval pipelines can reconstruct readable text on demand if necessary, but for machine-only workflows, ShortLang is sufficient.

### *5. Multi-Agent and Multi-Step Systems*

In systems where multiple LLM agents communicate, messages often contain redundant language. ShortLang provides a compact communication protocol that preserves meaning while reducing computational load across agents.

## *E. Related work*

ShortLang builds on longstanding research in text summarization, token compression, and semantic representation, while diverging from prior approaches in its prioritization of machine-oriented rather than human-oriented text transformation. Early work in extractive and abstractive summarization provided methods for condensing text but largely focused on readability and narrative coherence. Subsequent research in information retrieval introduced embeddings and vector representations that capture semantic meaning, laying the groundwork for embedding-based evaluation of textual similarity. Approaches such as byte-pair encoding and other subword compression techniques contributed to efficient tokenization but did not address semantic redundancy at the sentence or document level. More recently, work on system prompts, thought compression, and chain-of-thought distillation explored ways to minimize the reasoning steps or verbosity of LLMs, yet these methods operate at inference time rather than at the level of data preprocessing or representational transformation. ShortLang extends this body of research by offering a systematic means of encoding meaning into a compact textual form designed expressly for machine consumption, complementing and enhancing existing efforts in summarization, retrieval, and model optimization.

## **II. METHODS**

### *A. Rule-Based Compression Techniques*

The first methodological component of ShortLang relies on deterministic, rule-based compression strategies designed to eliminate linguistic redundancy while preserving semantic content. This approach begins with stopword removal, targeting function words such as articles, conjunctions, and auxiliary verbs that generally contribute syntactic structure rather than meaning. Beyond stopword pruning, the method incorporates systematic abbreviation of multi-word named entities into canonical identifiers, enabling significant length reduction without semantic distortion. Additional rules address redundant modifiers, repeated contextual information, and predictable syntactic constructions that can be safely collapsed. The rule-based system also includes domain-specific abbreviation dictionaries, permitting the compression of frequently occurring terminology in specialized corpora such as biomedical or legal texts. Although deterministic, these rules are crafted to avoid altering semantic embeddings significantly, and they

provide a baseline compression mechanism with transparent behavior and minimal computational cost.

#### *B. Model-Based Semantic Compression*

While rule-based methods provide consistency, they lack the capacity to interpret context or perform nuanced semantic abstraction. To address these limitations, the second methodological component employs model-based compression through fine-tuned language models. These models are trained on paired datasets consisting of full natural-language inputs and their ShortLang equivalents. Through supervised learning, the model acquires the ability to identify the minimal set of tokens required to preserve meaning, producing compressed outputs that often exceed the efficiency of rule-based systems. This approach leverages the model's internal representation of semantic relations, enabling it to collapse complex sentences, infer implicit relationships, and remove narrative filler. Model-based compression is especially effective in domains where meaning is distributed across long sequences, as the model can selectively retain only the concepts critical for downstream reasoning tasks. Over time, iterative refinement and reinforcement-based feedback can be used to improve compression quality by optimizing embedding similarity between original and compressed representations.

#### *C. Hybrid Compression Framework*

A third methodological pathway integrates the strengths of both rule-based and model-based strategies into a hybrid compression framework. In this approach, rule-based preprocessing is applied first to remove trivial redundancies, reduce token noise, and standardize entity representations. The output of this deterministic stage serves as input to a model-based compressor, which performs higher-order semantic shortening. By reducing linguistic variability before model inference, the hybrid system minimizes the model's burden and enhances its ability to produce stable, compact outputs. This layered design improves generalization across domains and reduces hallucinations by constraining the model to operate on cleaner input signals. Moreover, hybrid compression enables adjustable compression levels, where users can modulate the intensity of rule-based pruning and model-based abstraction depending on the requirements of a specific application, such as vector retrieval, inter-agent reasoning, or dataset minimization.

#### *D. Embedding-Based Validation and Quality Assessment*

Ensuring semantic retention is central to the effectiveness of ShortLang, and thus the final methodological component focuses on embedding-based evaluation techniques. Text embeddings derived from widely used embedding models serve as quantitative proxies for semantic meaning. After compressing text into ShortLang, the cosine similarity between the embeddings of the original and compressed texts is computed. High similarity scores indicate successful preservation of meaning, whereas lower scores suggest semantic drift or excessive compression. This evaluation method allows for objective, scalable, and domain-agnostic assessment without requiring human interpretation. In experimental workflows, embedding-based analysis can be complemented by downstream task evaluations, such as comparing model reasoning accuracy, retrieval performance, or classification outcomes when using original versus ShortLang-processed inputs. Together, these metrics form a robust validation pipeline that guides model fine-tuning, informs rule adjustments, and ensures the system maintains an optimal balance between brevity and semantic fidelity.

### **III. RESULTS**

The experiments are yet to be performed. Results will be displayed here after completion.

#### **IV. DISCUSSION**

Despite its promise, ShortLang poses several challenges. Over-compression may lead to semantic drift or loss of nuance. Embedding similarity does not perfectly capture meaning, particularly for long or conceptually complex texts. Some tasks require stylistic or syntactic detail that ShortLang intentionally strips away. Moreover, automatic generation of ShortLang may require careful tuning to avoid inconsistencies or hallucinations, especially when using model-based methods.

#### **V. CONCLUSION**

ShortLang represents a promising direction for enhancing efficiency in large language model workflows by reducing token overhead while retaining semantic content. By providing a minimal viable text representation, ShortLang enables faster reasoning, lower training costs, more efficient storage, and more scalable vector retrieval. Through embedding-based validation and a combination of rule-based and model-based generation techniques, ShortLang can serve as a robust intermediate representation for machine reasoning. As models continue to scale and demand more resources, ShortLang offers a compelling tool for optimizing the performance and cost-effectiveness of language-based AI systems.

#### **REFERENCES**

- [1] <To be added>