

Weekly Technical Report

Subject: Tokenization in the Context of Large Language Models (LLMs)

Prepared by: Anbhi Thakur

Week Covered: July 14 – July 19, 2025

Date of Submission: July 20, 2025

1. Purpose of the Study

The objective of this week's research was to develop a comprehensive understanding of tokenization as applied to NLP and large language models (LLMs). The study aimed to explore how tokenization works, why it is a necessary preprocessing step for LLMs, the types of tokenization strategies employed, and how token limits affect model performance, latency, and API cost. This research serves as a foundation for building and optimizing future NLP-based applications and prompt designs.

2. Scope of Work

The scope included:

- Understanding the theoretical foundation of tokenization in NLP
- Analyzing the role of tokenization in model architecture and training
- Evaluating different tokenization strategies (word, subword, character, BPE)
- Exploring the impact of token constraints on model inference and cost
- Presenting practical examples and industry usage patterns
- Documenting findings in a structured report format

3. Topics Covered

a. Definition and Importance of Tokenization

Tokenization is the process of dividing raw text into manageable and processable units called tokens. These tokens can represent full words, subwords, characters, or byte-pair encodings. This segmentation is essential because language models operate on numerical representations of text. Tokenization converts raw strings into token IDs, which are further mapped to embeddings.

b. Tokenization Strategies in LLMs

Modern LLMs employ subword-level tokenization using techniques like Byte Pair Encoding (BPE), WordPiece, and SentencePiece. These methods ensure efficient handling of rare, compound, or out-of-vocabulary words, better generalization across languages, and improved robustness for code and non-standard inputs.

c. Necessity of Tokenization for LLMs

Tokenization is critical for:

- Converting textual data into numerical token IDs
- Mapping to learned vector representations (embeddings)
- Efficient vocabulary management
- Handling language diversity and variant input structures

d. Impact of Token Limits

LLMs have predefined token limits which affect prompt capacity and response length. Exceeding token limits can lead to input truncation or model errors. Additionally, more tokens increase inference latency and cost. Proper token usage ensures more predictable and cost-effective model deployment.

Model Comparison Example:

GPT-4 Turbo: 128K tokens, \$0.01 per 1K tokens (input), \$0.03 per 1K tokens (output)

Claude 3 Opus: 200K tokens, \$0.015 per 1K tokens (input), \$0.075 per 1K tokens (output)

4. Key Insights

- **Foundational Role of Tokenization**

Tokenization forms the basis of text preprocessing in all transformer-based LLMs. Without proper tokenization, models cannot effectively encode, interpret, or generate natural language. It directly governs how meaning is represented and retained within the model's context window.

- **Influence on Model Generalization and Performance**

The choice of tokenizer architecture (e.g., BPE, WordPiece, SentencePiece) determines how well a model handles linguistic diversity, morphological complexity, and domain-specific text. Subword-level tokenization allows models to manage rare or novel words, reducing the need for retraining and improving generalization.

- **Implications for System Efficiency and Cost Management**

Tokenization directly impacts the length of input and output sequences. Since LLM inference and API usage are priced per token, optimizing token counts is essential for deploying scalable and cost-effective solutions. Efficient prompt engineering strategies—such as compression, instruction tuning, or content summarization—are necessary to reduce redundancy and stay within token limits.

- **Trade-offs Between Precision and Compression**

Developers must balance compact token usage with semantic fidelity. Over-compressing prompts can lead to context loss, while overly verbose prompts can cause truncation, latency, and higher API charges. Thus, tokenization must be managed with both performance and linguistic integrity in mind.

5. Knowledge and Skills Acquired

- **Advanced Understanding of Tokenization Algorithms**

Acquired an in-depth conceptual and practical understanding of major tokenization approaches:

- *Byte-Pair Encoding (BPE)*: Utilizes statistical merging of frequent character sequences to generate subword vocabularies.
- *WordPiece*: Originally developed for BERT, it prioritizes frequency and maximum likelihood estimates while maintaining lexical coherence.
- *SentencePiece*: Treats the input as a raw byte stream, making it suitable for multilingual tasks without language-specific preprocessing.

- **Token Optimization Techniques**

Gained the ability to:

- Predict token count based on input structure.
- Evaluate token utilization efficiency across different prompts.
- Design optimized, compressed prompts that meet token limitations while preserving intent.

- **Awareness of Operational Impacts**

Developed awareness of how tokenization affects:

- **Latency**: Longer token sequences increase compute time linearly or more.
 - **Accuracy**: Poor tokenization may misrepresent rare or compound words.
 - **Cost**: Every 1,000 tokens processed contributes to overall API billing, making token-aware development essential for sustainability.
-

6. Tasks Completed

- **Comprehensive Literature Review**

Reviewed foundational and recent academic papers, technical documentation, and model-specific tokenizer implementations to understand current best practices in tokenization.

- **Technical Report Drafting and Structuring**

Authored a detailed technical report covering theoretical background, tokenization mechanisms, practical examples, token limit implications, and comparative model pricing.

- **Behavioral and Cost Analysis of Tokenization Across Models**

Conducted an analytical comparison of token behaviors across various LLMs (e.g., GPT-4, Claude, PaLM), including their token limits, encoding methods, and token pricing per input/output. Mapped this data to cost-optimization recommendations.

- **Summarization of Tokenization Strategies and Best Practices**

Synthesized findings into actionable strategies for developers and researchers, focusing on:

- Prompt compression without semantic loss
 - Token-aware design for multilingual and specialized use cases
 - Trade-offs between prompt verbosity and context retention
-

7. Challenges Encountered

- **Deciphering Proprietary Tokenization Mechanisms**

Many LLM providers (e.g., OpenAI, Anthropic) use proprietary tokenizers with limited public documentation. This required indirect analysis through APIs, tokenizer simulators, and usage examples to reverse-engineer their behavior.

- **Complexity in Model Comparison and Token Billing**

Comparing different LLMs was complicated by:

- Varying token definitions (e.g., Unicode-based vs byte-level)
 - Non-uniform pricing models (input vs output token billing rates)
 - Differences in context window sizes and maximum generation lengths
- Careful alignment and normalization of this information was necessary for a meaningful comparative analysis.

- **Prompt Design Trade-offs**

Designing prompts that balance brevity and effectiveness proved challenging, especially when dealing with complex tasks (e.g., summarization, multi-turn dialogue). Fine-tuning prompts for both clarity and token efficiency required iterative experimentation.

8. Proposed Next Steps

- Develop a prompt optimization tool to reduce token usage
- Explore code-based tokenization and multilingual applications
- Apply token efficiency strategies in upcoming LLM integration projects