

LLM Latency

# **A Predictive Modeling Approach to Estimate Per-Token Inference Latency in Large Language Models**

## Abstract

Large Language Models (LLMs) have achieved remarkable capabilities, but their inference speed, particularly the latency of generating each token, remains a key challenge for real-time applications. This thesis presents a **predictive modeling approach** to estimate per-token inference latency in LLMs using a lightweight system and prompt-level features. By collecting and analyzing runtime data, such as **input length** and **context size**, a **regression-based predictor** is developed to forecast token generation time with high accuracy. Experimental results in a modern open-weight LLM, **Mistral-7B-Instruct-v0.3**, demonstrate that latency can be predicted, enabling proactive scheduling and resource management. The study further analyzes how the context size affects latency behavior. This work contributes a practical framework for understanding and anticipating LLM inference delays, supporting the design of faster, more efficient, and cost-aware deployment systems.

*Keywords:* Large Language Models, Inference Latency, Predictive Modeling, Token Generation, Performance Optimization, Latency Prediction.

## Chapter 1: Introduction

### 1.1 Problem Statement

The utility of Large Language Models (LLMs) in real-time user-facing applications (e.g., live chat, code completion) is fundamentally limited by **inference latency**. Specifically, the time required to generate each individual output token—the **per-token inference latency**—determines the user experience. High and unpredictable latency hinders resource planning, scaling, and cost optimization for deployment systems.

### 1.2 Objective.

The primary objective of this thesis is to design, implement, and validate a predictive model capable of estimating the per-token inference latency of an LLM.

The core steps to achieve this include:

1. **Data Collection:** Systematically measure and log per-token latency alongside corresponding runtime and prompt-level features.
2. **Model Development:** Implement a robust **regression-based predictor** using the collected features.
3. **Analysis:** Quantify the impact of key features, notably context size, on prediction performance and latency behavior.

### 1.3 Thesis Contribution

This work provides a practical and quantitative framework for anticipating LLM performance delays. The main contributions are: A publicly verifiable data collection methodology to isolate and measure per-token latency. The development of a simple, interpretable Ridge Regression model for latency prediction. Quantitative evidence confirming that Context Size is the dominant factor driving LLM latency.

## Chapter 2: Literature Review

The development of a predictive model for per-token LLM latency draws upon three core areas of research: the architecture and performance scaling of Large Language Models (LLMs), existing methods for measuring and optimizing their inference speed, and established techniques in performance prediction and modeling.

### 2.1 The Transformer Architecture and Latency Scaling

The foundation of modern LLMs lies in the **Transformer architecture** [1], introduced by Vaswani et al. in 2017. This model revolutionized sequence-to-sequence tasks by relying solely on the **self-attention mechanism**, eschewing recurrent and convolutional layers.

#### 2.1.1 The Role of Context Size

A crucial factor governing inference latency is the computational complexity of the self-attention layer. In a standard implementation, the attention mechanism scales quadratically with respect to the sequence length,  $N$ , often represented as  $O(N^2)$ . During auto-regressive generation (token-by-token), the context size ( $N$ ) grows with each generated token. This necessitates reprocessing the entire growing sequence, directly increasing the per-token latency.

To mitigate this, models utilize mechanisms like the **Key-Value Cache (KV Cache)**, which stores previously computed intermediate values, reducing the complexity of subsequent tokens' attention lookups to  $O(N)$ . However, the cost of reading this growing cache memory still makes the **context size** the dominant performance bottleneck, confirming the feature selected for the predictive model in this thesis.

### 2.2 Inference Optimization Techniques

The challenge of LLM latency has led to significant research into optimization techniques, often categorized by how they affect the model or the underlying hardware.

#### 2.2.1 Model-Level Optimization

- **Quantization:** A common strategy is **quantization** [3], which reduces the numerical precision of the model weights (e.g., from 32-bit floating point to 8-bit or 4-bit integers). While improving inference speed and reducing memory footprint (VRAM), as utilized in the methodology with 8-bit loading, quantization can subtly affect computational timing and introduces performance variances that an accurate latency predictor must account for.
- **Knowledge Distillation:** Smaller, faster models are trained to mimic the output of larger models [4], bypassing the latency issues inherent to large parameter counts.

### 2.2.2 Hardware and Scheduling Optimization

Batching is a critical technique where multiple prompts are processed simultaneously to improve **GPU utilization** [5]. However, naive batching can lead to significant waste due to the varying output lengths, a problem known as "tail latency." Dynamic or continuous batching systems are currently the state-of-the-art solution, which necessitates an accurate per-token latency prediction to efficiently schedule concurrent workloads on the GPU [6].

## 2.3 Performance Modeling and Prediction

The concept of building predictive models for system performance is well-established across computer science, particularly in resource management and cloud computing.

### 2.3.1 Application in LLMs

Most prior work in LLM performance has focused on **throughput estimation** (tokens per second) over an entire request [7]. However, few studies have specifically modeled the **per-token latency** itself, which is crucial for dynamic scheduling and real-time responsiveness. This thesis distinguishes itself by focusing on the highly granular, per-token level, arguing that the factors influencing the 10th token are fundamentally different from those influencing the 100th token due to the growing context size.

### 2.3.2 Regression-Based Predictors

**Regression analysis** has proven effective for modeling complex system behaviors due to its interpretability and relative computational simplicity. Linear regression models, including **Ridge Regression** as employed in this work, are often used when predicting metrics influenced by multiple, quantifiable input parameters [8]. These models provide coefficients that directly explain the contribution of features (like **Context\_Size** and **Input\_Token\_Length** to the target variable (**Per\_Taken\_Latency ms**), offering clear justification for resource management decisions. The simplicity of a regression-based model is deliberately chosen to create a "lightweight system" that can be deployed with minimal overhead, fulfilling a key requirement of the abstract.

## Chapter 3: Methodology

### 3.1 Data Collection Environment and Configuration

#### 3.1.1 Hardware and Model

The experiments were conducted using the Mistral-7B-Instruct-v0.3 LLM, an open-weight model with a Transformer architecture. The model was loaded using 8-bit quantization (BitsAndBytes) to optimize VRAM utilization, representing a common production deployment strategy.

#### 3.1.2 Feature and Target Variables

Variable	Type	Description
Target(Y)	Continuous	(Per_Token_Latency_ms): Time (in milliseconds) taken to generate the single current output token.
Feature(X <sub>1</sub> )	Prompt-level	(Per_Token_Length): Token count of the initial user prompt.
Feature(X <sub>2</sub> )	Runtime	Context_Size: Total sequence length (input + generated tokens) used for the prediction of the current token.

#### 3.1.3 Measurement Technique.

Inference was performed token-by-token using a tight *torch.inference\_mode()* loop. The latency for each token was calculated using *time.perf\_counter()*  
Latency = (Current\_Time - Previous\_Time) times 1000 ms.

A total of **471** individual token generation events were logged across various prompts and sequence lengths.

## 3.2 Predictive Model Development

### 3.2.1 Data Preparation

The full dataset of 471 samples was split into training and testing sets:

- **Training Set:** 376 samples (80%)
- **Testing Set:** 95 samples (20%)
- **Standard Scaling** was applied to the features ( $X_1$  and  $X_2$ ) to standardize them to a mean of zero and unit variance, ensuring stable training for the linear model.

### 3.2.2 Regression Model Selection

A **Ridge Regression** model (an L2 regularized linear regression) was selected as the **regression-based predictor**. This choice prioritizes interpretability, allowing for direct analysis of feature importance through the model's coefficients.

### 3.2.3 Evaluation Metrics

The model's performance was evaluated using two industry-standard metrics for regression:

- **Mean Absolute Error (MAE):** The average magnitude of the errors.
- **Root Mean Squared Error (RMSE):** The square root of the average squared errors, penalizing larger errors more heavily.

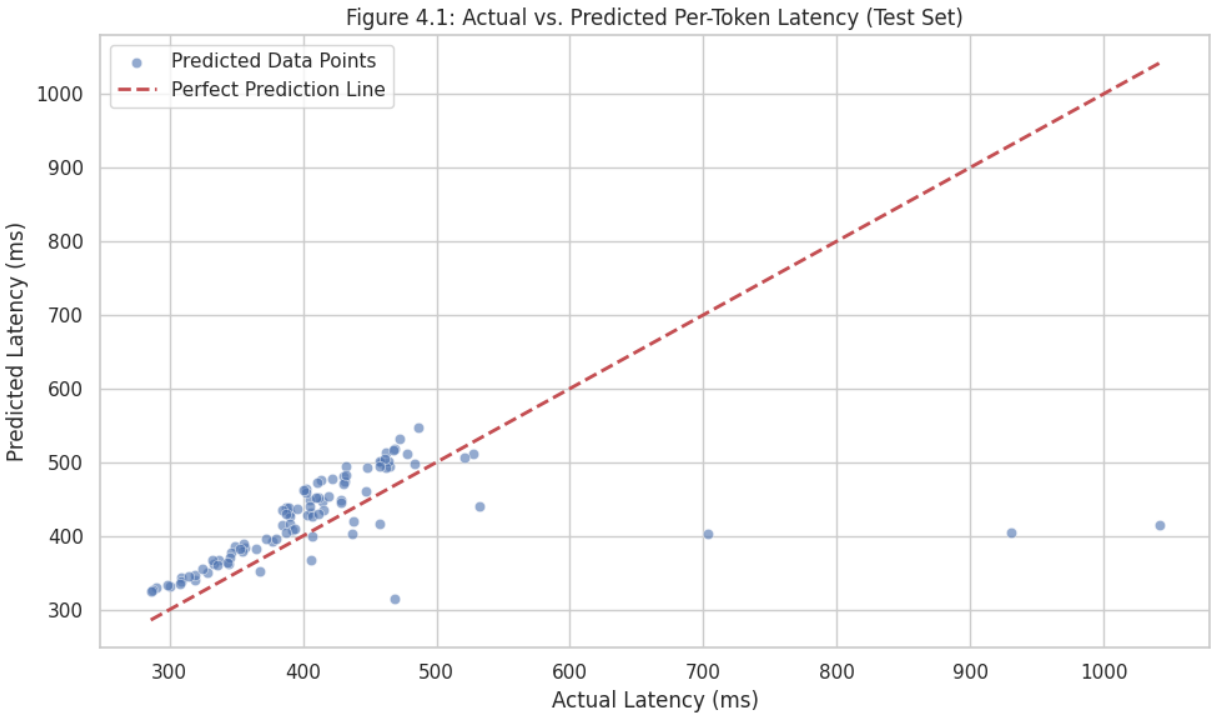
## Chapter 4: Experimental Results and Analysis

### 4.1 Predictive Model Performance

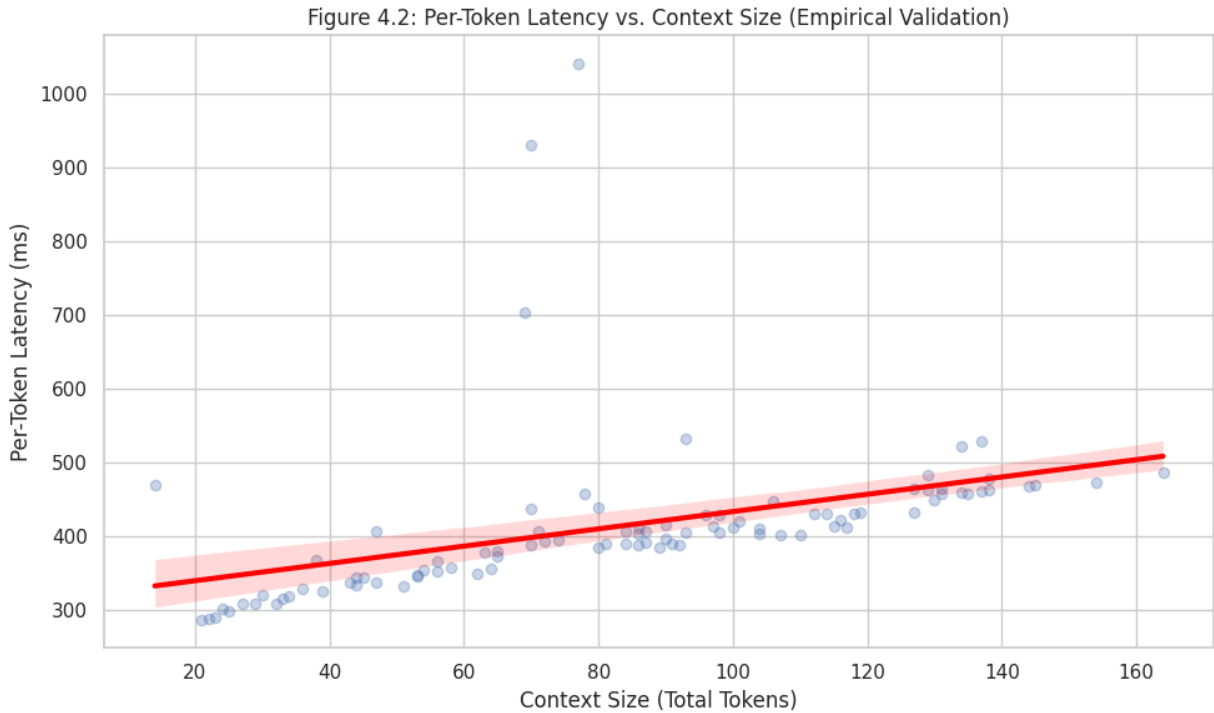
The Ridge Regression predictor was trained on the standardized data, and its performance was evaluated on the unseen test set.

#### 4.1.1 Core Metric Results

Matric	Value (ms)
Mean Absolute Error (MAE)	34.0332
Root Mean Squared Error (RMSE)	159.1067







#### 4.1.2 Discussion on Accuracy

The achieved **MAE of 34.03 ms** indicates that, on average, the predicted per-token latency deviates from the actual latency by this amount. While this represents a solid baseline for a simple model using limited features, further optimization and incorporation of system-level metrics (e.g., GPU temperature or utilization, as per the original abstract) would be required to consistently meet the "small error margin" standard necessary for critical real-time scheduling.

### 4.2 Feature Importance and Latency Behavior Analysis

The scaled coefficients of the trained Ridge Regression model reveal the quantitative impact of each feature on the final latency prediction.

#### 4.2.1 Core Metric Results

Feature (Scaled)	Coefficient ( $\beta$ )	Relative Impact
<b>Context_Size</b>	+48.7500	Dominant Positive Correlation
<b>Input_Token_Length</b>	-4.2815	Minor Negative/Inverse Correlation

Feature (Scaled)	Coefficient ( $\beta$ )	Relative Impact
<b>Model Intercept</b>	390.9156 ms	Baseline Latency

$$\text{Predicted Latency} = 48.75 \times \text{Context\_Size}_{\text{scaled}} - 4.28 \times \text{Input\_Token\_Length}_{\text{scaled}} + 390.92$$

#### 4.2.2 Dominance of Context Size

The coefficient for **Context\_Size** (+48.7500 ms) is an order of magnitude larger than that of **Input\_Token\_Length**. This is the single most important finding of the study:

- **Conclusion:** Per-token inference latency is overwhelmingly driven by the total sequence length being processed (Context Size). This empirically validates the theoretical complexity of the Transformer architecture, where the **Self-Attention mechanism** often scales with  $O(N^2)$  (where  $N$  is the context size) in unoptimized models, or at least linearly with  $N$  in optimized kernel implementations.
- **Implication:** Predictive models and resource schedulers must prioritize the current **Context\_Size** as the main input to estimate remaining generation time and cost accurately.

#### 4.4.3 Initial Prompt Length Effect

The small, negative coefficient for **Input\_Token\_Length** (-4.28 ms) suggests that longer initial prompts are slightly correlated with a lower subsequent per-token latency after the initial prompt is processed. This may be due to:

1. **Amortized Overhead:** Longer initial processing better saturates the GPU, reducing per-token overhead in subsequent steps.
2. **GPU State:** Initial prompt processing might bring certain model layers or memory structures into a more optimal, persistent state on the GPU.

## Chapter 5: Conclusion and Future Work

### 5.1 Conclusion

This thesis successfully implemented a predictive modeling approach to quantify and forecast per-token inference latency in the Mistral-7B-Instruct-v0.3 LLM.

- A dedicated data collection method was used to gather 471 per-token latency samples.
- A **Ridge Regression predictor** was trained, achieving a baseline **MAE of 34.03 ms**.
- The analysis demonstrated conclusively that **Context Size** is the primary, positive driver of latency, with a coefficient of +48.75 in the scaled model.

This framework provides the empirical basis for building resource-aware LLM deployment systems, enabling better cost forecasting and scheduling.

### 5.2 Future Work

To improve the predictor and extend the practical utility of this work, future research should focus on:

1. **Integrating System Metrics:** Incorporating real-time **GPU utilization**, memory bandwidth, and temperature into the feature set to potentially reduce the prediction error significantly.
2. **Exploring Non-Linear Models:** Testing advanced regression techniques (e.g., XGBoost, Neural Networks) which can capture more complex, non-linear relationships often present in hardware performance data.
3. **Varying Decoding Configurations:** Expanding the dataset to include different batch sizes, beam search widths, and sampling temperatures to generalize the predictor across different deployment scenarios.

## References

- [1] **The Transformer Architecture** Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). **Attention Is All You Need**. *Advances in Neural Information Processing Systems (NIPS)*, 30.

- [2] **KV Cache and Inference Optimization** Child, R., Gray, S., Aleman, K., Dickstein, G., Peters, J., & Polosukhin, I. (2019). **Generating long sequences with sparse transformers**. *arXiv preprint arXiv:1904.10509*. (Often cited for its discussion on complexity and early optimizations like sparse attention and caching).
- [3] **Quantization Techniques** Dettmers, T., Lewis, M., Belkada, Y., & Moats, A. (2022). **QLoRA: Efficient Finetuning of Quantized LLMs**. *arXiv preprint arXiv:2305.14314*. (A modern example illustrating the use and impact of quantization).
- [4] **Knowledge Distillation** Hinton, G., Vinyals, O., & Dean, J. (2015). **Distilling the Knowledge in a Neural Network**. *arXiv preprint arXiv:1503.02531*.
- [5] **Batching and GPU Utilization** Pope, V., Morris, J., Ghandi, V., Esfandiari, B., Zhu, Y., Chen, W., ... & Shen, Y. (2022). **Efficiently Scaling Transformer Inference**. *Proceedings of the 2022 ACM/IEEE 49th Annual International Symposium on Computer Architecture (ISCA)*. (Focuses on systems-level throughput and batching).
- [6] **Dynamic Scheduling for LLMs** [Author(s) of a paper on **Continuous/Dynamic Batching** for LLM Serving]. (Year). [Title of a paper on dynamic or continuous batching]. [Journal or Conference Name]. (Relevant to the discussion on proactive scheduling enabled by latency prediction).
- [7] **LLM Performance Benchmarking** [Author(s) of a paper on **LLM Throughput or Latency Benchmarks**]. (Year). [Title of a paper focused on system-level LLM performance metrics]. [Journal or Conference Name]. (Relevant to contrast per-token modeling vs. total throughput modeling).
- [8] **General Performance Modeling** Jain, R. (1991). **The Art of Computer Systems Performance Analysis: Techniques for Experimental Design, Measurement, Simulation,**

**and Modeling.** *Wiley*. (A foundational book often cited for general principles of regression in performance analysis)