

AI-POWERED SUB-TOPIC EXTRACTION AND CONVERSATIONAL NEWS ANALYSIS

Project report submitted in partial fulfillment of the requirement for the degree of

Bachelor of Technology

Submitted by

Akshat Shashi (2110110081)

Aviral Chawla (2110110151)

Preksha Khera (2110110391)

Under Supervision of

Dr. Dolly Sharma

Department of Computer Science Engineering



Department of Computer Science Engineering

School of Engineering

Shiv Nadar Institution of Eminence

(December 2024)

Date: 02 Dec 2024

Certificate

This is to certify that the project titled “**AI-Powered Sub-Topic Extraction And Conversational News Analysis**” is submitted by **Akshat Shashi** (Roll no.2110110081), **Aviral Chawla** (Roll no.2110110151) and **Preksha Khera** (Roll no.2110110391), Computer Science and Engineering, Shiv Nadar IOE in the partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science and Engineering**. This work was completed under the supervision of Prof. Dolly Sharma. No part of this dissertation/report has been submitted elsewhere for award of any other degree.

Signature with date: _____

Dolly Sharma
2/12/24

Name of Supervisor: Prof. Dolly Sharma

Designation: Professor

Affiliation: Shiv Nadar Institution of Eminence

CANDIDATE DECLARATION

We declare that this written submission represents our ideas in our own words and where others' ideas or words have been included, we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academic honesty and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/source in our submission. We understand that any violation of the above will be cause for disciplinary action by the University and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken when needed.

Akshat Shashi

(Signature)

Name: AKSHAT SHASHI

Roll No.: 2110110081

Date: 2/12/2024

Aviral

(Signature)

Name: AVIRAL CHAWLA

Roll No.: 2110110151

Date: 02/12/2024

Preksha

(Signature)

Name: PREKSHA KHERA

Roll No.: 2110110391

Date: 02/12/2024

ABSTRACT

This research paper explores the inter dependencies between quality, context window size, and token price in artificial intelligence (AI) models, focusing on the trade-offs between cost and performance. We analyze various state-of-the-art models, with particular emphasis on Llama 3.2 96B (Vision), which demonstrates an optimal balance between high-quality outputs, large context window capability, and moderate input token price. The paper investigates total response time and its variance over time, emphasizing the importance of stability for real-time applications. The findings reveal that higher-quality models generally offer larger context windows at the expense of higher token costs, while models like GPT-4o Mini excel in speed and stability. Our conclusions suggest that the choice of an AI model must be application-specific, balancing speed, quality, and cost considerations. The study provides recommendations for both enterprises and developers and identifies future research opportunities to enhance model efficiency and reduce computational costs.

Contents

1	Introduction	10
1.1	Text Summarization Techniques:	10
1.2	Named Entity Recognition (NER) and Relation Extraction:	10
1.3	Question Answering (QA) Systems:	11
1.4	Conversational AI and Dialogue Management:	11
2	Related Work	13
2.1	Context Window and Quality in AI Models	13
2.1.1	Role of Context Windows in Model Performance	13
2.1.2	Challenges with Larger Context Windows	13
2.1.3	Advancements in Managing Context Windows	14
2.2	Token Price Optimization	14
2.2.1	Balancing Cost and Quality	14
2.2.2	Key Techniques	14
2.3	Stability in Total Response Time	15
2.3.1	Importance of Stability	15
2.3.2	Challenges in Achieving Stability	15
2.3.3	Solutions and Benchmarks	15
3	Performance Dynamics and Applications of Advanced AI Models	16
3.1	Understanding Model Quality, Context Windows, and Token Costs	16
3.1.1	Insights on Quality and Context Window Relationships	16
3.1.2	Applications of High-Quality Models	16
3.1.3	Trade-Offs Between Quality, Token Costs, and Complexity	17
3.2	Analysis of Total Response Time Across Models	17
3.2.1	Observations on Total Response Time	17
3.2.2	Variance in Response Time	17
3.3	Stability of Total Response Time Over Time	18
3.3.1	Importance of Stability	18
3.3.2	GPT-40 Mini’s Performance in Stability	18
3.4	Applications of Advanced AI Models	18
3.4.1	Multimodal Applications	18
3.4.2	Real-Time Use Cases	19
3.5	Summary of Performance Dynamics	19

4	Methodology	20
4.1	Data Acquisition	20
4.2	Training Models	22
4.2.1	Preprocessing	22
4.2.2	Training Process	23
4.3	Performance Evaluation	24
5	Graphs and Insights	26
5.1	Latency by Input Token Count (Context Length)	26
5.1.1	Understanding the Graph	26
5.1.2	Key Observations	27
5.1.3	Implications	27
5.2	Latency Over Time	27
5.2.1	Understanding the Graph	27
5.2.2	Key Observations	28
5.3	Total Response Time	29
5.3.1	Understanding the Graph	29
5.3.2	Key Observations	29
5.3.3	Implications	29
5.4	Total Response Time by Input Token Count (Context Length)	30
5.4.1	Understanding the Graph	30
5.4.2	Key Observations	30
5.4.3	Implications	30
5.5	Output Speed by Input Token Count (Context Length)	31
5.5.1	Understanding the Graph	31
5.5.2	Key Observations	31
5.5.3	Implications	32
5.6	Output Speed Over Time	32
5.6.1	Understanding the Graph	32
5.6.2	Key Observations	33
5.6.3	Implications	34
5.7	Quality vs. Context Window, Input Token Price	34
5.7.1	Understanding the Graph	34
5.7.2	Key Observations	35
5.7.3	Implications	35
5.8	Total Response Time Variance	36
5.8.1	Understanding the Graph	36
5.8.2	Key Observations	36
5.8.3	Implications	37
5.9	Total Response Time Over Time	37
5.9.1	Understanding the Graph	37
5.9.2	Key Observations	37
5.9.3	Implications	38
5.10	Output Speed Variance by Model (Derived from Composite Data)	39
5.10.1	Understanding Output Speed Variance	39
5.10.2	Key Observations	39

5.10.3	Implications	40
5.11	Total Response Time Variance	40
5.11.1	Understanding the Metric	40
5.12	Future Considerations Based on Composite Analysis	41
5.12.1	Emerging Insights	41
5.12.2	Key Observations	42
5.12.3	Implications	42
5.13	Quality Index Variations Across Models	42
5.13.1	Understanding Quality Index	42
5.13.2	Key Insights from Quality Index Comparison	43
5.13.3	Implications	43
5.14	Combining Latency and Output Speed for Scalability	43
5.14.1	The Latency-Speed Nexus	43
5.14.2	Key Observations from Data	44
5.14.3	Implications	44
5.15	Token Price as a Function of Context Window	44
5.15.1	Understanding Token Pricing Trends	44
5.15.2	Key Observations	45
5.15.3	Implications for Pricing Strategies	45
5.16	Addressing Variance in Larger Models	45
5.16.1	Challenges in Large Model Deployments	45
5.16.2	Strategies for Reducing Variance	45
5.16.3	Future Research Directions	46
5.17	Comparative Cost-Efficiency Analysis	46
5.17.1	Metrics for Cost-Efficiency	46
5.17.2	Model Comparisons:	46
5.17.3	Recommendations for Cost Optimization	47
5.18	Real-World Deployment Challenges	47
5.18.1	Infrastructure Challenges	47
5.18.2	Ethical and Regulatory Considerations	48
5.18.3	Recommendations	48
5.19	Comprehensive Scope for Future Work	48
5.19.1	Optimizing Larger Context Windows	48
5.19.2	Integrating Multimodal Capabilities	48
5.19.3	Improving Cost-Effectiveness	49
5.19.4	Addressing Variance in Larger Models	49
6	Summary and Conclusion	50
6.1	Key Findings	50
6.1.1	Practical Implications of the Study	52
6.1.2	Main Conclusions	53
6.1.3	Limitations of the Study	53
6.1.4	Implications for Future Research	54

7	Scope for Future Work	55
7.1	Reducing Latency in Larger Models	55
7.2	Enhancing Multimodal Capabilities	55
7.3	Improving Stability in Large Models	56
7.4	Reducing Costs for Large Context Models	56
7.5	Expanding Context Windows	57
7.6	Energy Efficiency and Sustainability	57
7.7	Democratization of AI	57

List of Figures

4.1	Training and Evaluation of Gemma 2	23
4.2	Claude Conversation	25
4.3	Gemma Conversation	25
4.4	Gemini Conversation	25
5.1	Graph- Latency by Input Token Count (Context Length)	26
5.2	Graph- Latency over Time	28
5.3	Graph- Total Response Time	29
5.4	Graph- Total Response Time by Input Token Count (Context Length) . .	30
5.5	Graph-Output Speed by Input Token Count (Context Length)	31
5.6	Graph-Output Speed Over Time	33
5.7	Graph- Quality vs. Context Window, Input Token Price	34
5.8	Graph-Total Response Time Variance	36
5.9	Graph-Total Response Time Over Time	38
5.10	Graph- Output Speed Variance by Model	39
5.11	Graph-Total Response Time Variance	41

LIST OF SYMBOLS AND ABBREVIATIONS

- **GPT:** Generative Pre-trained Transformer
- **RLHF:** Reinforcement Learning from Human Feedback
- **AI:** Artificial Intelligence
- **TPU:** Tensor Processing Unit
- **GPU:** Graphics Processing Unit
- **96B:** 96 Billion Parameters
- **Llama:** Large Language Model Meta AI

Chapter 1

Introduction

The explosive growth of online news sources has created an information deluge. Humans struggle to process the sheer volume, often missing crucial details or nuances buried within articles. This necessitates the development of AI systems capable of not only summarizing news but also extracting subtle, sub-topic information and engaging in detailed, question-answering conversations based on the article's content. This report explores the current state of AI research relevant to this goal, focusing on techniques that can be applied to develop a system like Gemma 2.

1.1 Text Summarization Techniques:

Text summarization condenses large texts into concise summaries. Traditional extractive methods select key sentences based on metrics like sentence position or TF-IDF scores. While efficient, these methods often produce summaries that lack coherence and fail to capture nuanced details.

Abstractive summarization generates new text by paraphrasing and synthesizing information, mimicking human summarization. It provides more meaningful and fluid summaries but was initially prone to grammatical errors and incomplete content.

Transformer-based models like BART and T5 have significantly advanced abstractive summarization. These models, pretrained on large datasets, generate coherent and context-aware summaries. While they may occasionally miss subtle details, they represent a major leap forward in creating accurate and human-like summaries.

1.2 Named Entity Recognition (NER) and Relation Extraction:

Named Entity Recognition (NER) identifies key entities in text, such as people, organizations, locations, dates, and numerical values, transforming unstructured text into structured data. For example, in a news article, NER can highlight individuals, organizations, and locations involved in events, enabling systematic analysis and aiding tasks like information retrieval and text summarization.

Relation extraction complements NER by identifying the relationships between these

entities, such as a person "working for" an organization or an event "happening in" a location. Together, these techniques form the basis of knowledge graphs, which represent information as interconnected entities and relationships, supporting applications like search engines and recommendation systems.

Modern NER models, powered by contextual embeddings from transformers like BERT and RoBERTa, significantly improve the identification of complex or ambiguous entities. Relation extraction, however, remains challenging due to nuanced and implicit connections often found in text, particularly in news articles. Advanced models continue to enhance the accuracy of mapping these relationships, enabling more meaningful insights from unstructured data.

1.3 Question Answering (QA) Systems:

Question Answering (QA) systems are designed to provide answers to questions posed in natural language, using a given context as a reference. These systems play a crucial role in applications like virtual assistants, customer support, and information retrieval. Early QA systems primarily relied on keyword matching or template-based approaches, which were limited by their inability to handle variations in language and more complex queries.

Recent advancements in deep learning have revolutionized QA systems. Neural network-based models, particularly those built on transformer architectures such as BERT, RoBERTa, and ELECTRA, have set new benchmarks in QA performance. These models are pre-trained on vast amounts of data and fine-tuned for QA tasks, enabling them to understand context, identify relevant passages, and provide accurate answers. They achieve state-of-the-art results on datasets like SQuAD (Stanford Question Answering Dataset) and Natural Questions, demonstrating their ability to handle a wide range of queries.

However, despite their capabilities, these models face challenges with complex questions. Queries that require deep reasoning, an understanding of implicit relationships, or inference beyond explicit mentions in the text often expose the limitations of current QA systems. Addressing these challenges involves advancing model architectures, integrating external knowledge, and improving the ability to perform multi-step reasoning, making QA an active area of research in natural language processing.

1.4 Conversational AI and Dialogue Management:

Conversational AI systems are designed to interact with users in a natural, engaging, and informative way. These systems are widely used in applications such as chatbots, virtual assistants, and customer service platforms. One of the key components of conversational AI is dialogue management, which involves guiding the flow of conversation, understanding user intents, and generating relevant, context-aware responses. Effective dialogue management ensures that the conversation feels coherent and productive, allowing users to have meaningful interactions with the system.

Recent advancements in large language models (LLMs) like GPT-3 and LaMDA have significantly improved the capabilities of conversational AI. These models, built on transformer architectures, have demonstrated impressive abilities to generate human-like responses, handle a wide range of topics, and adapt to different conversational contexts.

They excel at producing fluent dialogue and maintaining engagement, making them powerful tools for interactive applications.

However, challenges remain in areas such as maintaining long-term coherence, handling ambiguous user inputs, and ensuring factual accuracy. While LLMs are adept at generating responses based on context, they sometimes struggle with maintaining consistency over extended conversations or interpreting ambiguous queries correctly. Additionally, providing accurate and reliable information is an ongoing challenge, as these models can generate plausible-sounding responses that are factually incorrect or misleading. Addressing these issues requires further advancements in model training, integration of external knowledge sources, and improved methods for tracking and managing conversation history. Despite these challenges, the progress in conversational AI continues to advance, with ongoing research focused on improving the accuracy, coherence, and reliability of these systems.

Chapter 2

Related Work

2.1 Context Window and Quality in AI Models

2.1.1 Role of Context Windows in Model Performance

The size of the context window significantly influences a model’s capacity to perform complex reasoning tasks and manage extended dialogues or documents. Larger context windows enable models to maintain coherence and refer back to previous points more effectively, improving task outcomes. For example:

- **Long-form reasoning:** Tasks such as summarizing books, conducting literature reviews, or generating extensive narratives benefit from larger context windows. Recent advancements, such as those highlighted in the survey on pre-trained language models for text generation [1], emphasize the importance of effective PLM design for handling long-context tasks.
- **Conversational AI:** Larger windows ensure continuity in multi-turn dialogues, reducing the need to repeatedly reintroduce information.
- **Document Retrieval and Summarization:** Retrieval-Augmented Generation (RAG) approaches, as seen in recent literature [2], use extended context windows to enhance the retrieval and synthesis of relevant information in query-based systems.

2.1.2 Challenges with Larger Context Windows

1. **Computational Costs:** The attention mechanism, foundational to transformer-based architectures, scales quadratically with the length of the context window. For instance, increasing the context window from 4,096 tokens to 8,192 tokens doubles the memory and computational requirements.
2. **Latency:** Larger context windows may lead to higher latency in response times, posing challenges in real-time applications.
3. **Data Sparsity:** With larger windows, not all tokens contribute equally to the task, potentially leading to inefficiencies unless mitigated by techniques like sparse attention.

2.1.3 Advancements in Managing Context Windows

Research has proposed innovative methods to address these challenges:

- **Sparse Attention Mechanisms:** Focus computation on the most relevant parts of the input, reducing the quadratic scaling issue.
- **Memory-Augmented Models:** Store and retrieve past context from an external memory, minimizing the need to include everything in a single forward pass (e.g., Retrieval-Augmented Generation). The work on GRM (Generative Relevance Modeling) [2] demonstrates how relevance-aware mechanisms can enhance memory-augmented models for long-context scenarios.
- **Hybrid Compression Techniques:** Combining lossy and lossless compression for context inputs has shown promise in maintaining essential information while reducing computational overhead.
- **Diffusion-Based Generation:** Recent advancements in diffusion language models [?] have shown potential in generating coherent long-context outputs with controlled randomness and high fidelity.

2.2 Token Price Optimization

2.2.1 Balancing Cost and Quality

Token price optimization involves reducing the computational and financial costs associated with processing tokens while maintaining the model’s quality. This is particularly critical for deploying large-scale models like GPT-4 in cost-sensitive applications. Techniques like controlled text generation with natural language instructions [3] offer methods to enhance both efficiency and adherence to specific constraints, balancing cost with quality.

2.2.2 Key Techniques

1. **Sparse Attention:** Reduces the number of computations by focusing only on relevant token pairs, significantly lowering the cost for models handling lengthy sequences.
2. **Model Quantization:** Converts weights and activations from higher-precision formats (e.g., float32) to lower-precision ones (e.g., int8), saving memory and computational resources. Studies like the one on text generation with diffusion language models [4] highlight the benefits of precision control for cost-efficient text generation.
3. **Adaptive Input Processing:** Leveraging techniques like token pruning or dynamic token dropping can further optimize costs. The AttrPrompt methodology [5] underscores the effectiveness of adaptive processing to balance diversity and bias in generated data.

4. **Batch Processing:** Efficient batching strategies can amortize token processing costs across multiple tasks or requests.
5. **Embedding Optimization:** As illustrated in the paper [6], optimizing embeddings for downstream tasks can significantly reduce token redundancy without compromising task accuracy.

2.3 Stability in Total Response Time

2.3.1 Importance of Stability

For real-time applications, such as chatbots, virtual assistants, and decision-support systems, consistent response times are paramount. Variability in response times can lead to user dissatisfaction and hinder practical utility in dynamic environments.

2.3.2 Challenges in Achieving Stability

1. **Input Variability:** Models encounter diverse inputs in length and complexity, which can significantly impact processing times.
2. **System Load:** High concurrency or system stress can exacerbate response-time inconsistencies.
3. **Dynamic Memory Allocation:** Inefficient allocation strategies can introduce latency spikes during high-load periods.

2.3.3 Solutions and Benchmarks

1. **Efficient Architectures:** Models like GPT-40 Mini use streamlined architectures optimized for real-time performance while maintaining quality benchmarks comparable to larger counterparts.
2. **Predictive Scheduling:** Dynamic resource allocation based on input characteristics and model state ensures smooth operation under varying loads.
3. **Hybrid Cloud Deployment:** Utilizing edge devices for pre-processing combined with cloud resources for heavier computations ensures a balance between latency and capability. Recent approaches leveraging synthetic data for embedding tasks [6] further exemplify efficient handling of high-dimensional inputs across distributed systems.
4. **Latency-Aware Transformers:** Architectures that adapt attention mechanisms based on input length and importance, as highlighted in [3], can stabilize response times without degrading output quality.
5. **Real-Time Monitoring Systems:** Implementing feedback loops that monitor and adjust resource allocation in real-time ensures stability across diverse use cases.

Chapter 3

Performance Dynamics and Applications of Advanced AI Models

This chapter delves into the intricate relationships between model quality, context windows, token costs, response times, and overall performance. It explores key trade-offs, highlights applications for various model configurations, and provides insights into selecting the optimal model for specific use cases.

3.1 Understanding Model Quality, Context Windows, and Token Costs

3.1.1 Insights on Quality and Context Window Relationships

- **Quality and Context Windows:** Models with higher quality, defined by their ability to generate accurate, coherent, and contextually relevant outputs, tend to feature larger context windows. This capability allows them to process and retain more extensive input sequences, making them well-suited for tasks requiring deep comprehension and long-form reasoning.
- **Balanced Performance:** Models such as **Llama 3.2 96B (Vision)** excel by combining a high-quality index with a relatively large context window. They also maintain moderate input token costs, making them versatile and cost-effective for diverse applications.

3.1.2 Applications of High-Quality Models

- **Document Processing:**
High-quality models are particularly effective in tasks like document summarization, legal text analysis, and financial reporting, where maintaining coherence across large input sizes is critical.
- **Multimodal Integration:**
Models like **Llama 3.2 96B (Vision)**, with multimodal capabilities, are ideal for combining textual and visual inputs. Applications include visual question answering

(VQA), document digitization, and analysis of data-rich visual content such as tables or charts embedded in text.

3.1.3 Trade-Offs Between Quality, Token Costs, and Complexity

- **Higher Token Costs:**
Larger context windows and more advanced models come with increased input token costs, as they require more computational resources to handle the expanded input size.
- **Computational Complexity:**
Expanding the context window introduces additional computational overhead, slowing down the model’s output generation and increasing latency.
- **Cost vs. Utility:**
While high-quality models are indispensable for applications demanding precision and long-term context retention, users must balance these benefits against the increased operational costs, especially in high-volume scenarios.

3.2 Analysis of Total Response Time Across Models

3.2.1 Observations on Total Response Time

- **Low Response Time in Smaller Models:**
Models like **GPT-40 Mini** achieve impressively low total response times, approximately **0.8 seconds for 1k tokens**, making them ideal for real-time applications where speed is paramount.
- **High Response Time in Larger Models:**
Larger models such as **GPT-4** exhibit longer total response times due to their increased computational demands. While they deliver higher accuracy and better reasoning capabilities, their slower processing makes them less suited for time-sensitive tasks.

3.2.2 Variance in Response Time

- **Consistency in Smaller Models:**
GPT-40 Mini demonstrates exceptional consistency, exhibiting minimal variance in response times even under varying workloads. This predictability makes it reliable for critical applications.
- **Higher Variability in Larger Models:** Models like **Llama 3.1 Instruct 4B** show significant fluctuations in total response times, primarily due to their larger architecture and sensitivity to input complexity. This variability limits their usability in real-time systems where stability is essential.

3.3 Stability of Total Response Time Over Time

3.3.1 Importance of Stability

Stable response times are critical for applications where unpredictable delays can disrupt user experience or functionality.

- **Customer Support Systems:**
Consistent performance ensures smooth interactions with AI chatbots, reducing user frustration and improving satisfaction.
- **Autonomous Systems:**
Predictable response times are essential in fields like robotics and automated trading, where delays can lead to inefficiencies or operational risks.

3.3.2 GPT-40 Mini's Performance in Stability

- **Most Stable Model Over Time:**
Among the evaluated models, **GPT-40 Mini** emerges as the most stable in terms of response time. Its lightweight design and optimized architecture ensure consistent performance, making it a preferred choice for scenarios demanding reliability.
- **Low Operational Risk:**
The stability of GPT-40 Mini allows developers to confidently deploy it in real-time systems, minimizing the risk of unexpected delays or performance drops.

3.4 Applications of Advanced AI Models

3.4.1 Multimodal Applications

- **Text-Image Integration:**
Multimodal models like **Llama 3.2 96B (Vision)** expand the scope of AI applications by combining text and visual data. These models enable tasks such as:
 - **Visual Question Answering (VQA):** Providing contextually accurate answers to questions based on images.
 - **Document Analysis:** Interpreting complex documents that include tables, charts, and images alongside text.
- **Content Creation and Education:**
By leveraging their multimodal capabilities, these models support advanced educational tools, digital content generation, and automated analysis of multimedia materials.

3.4.2 Real-Time Use Cases

- **Chatbots and Virtual Assistants:**

Smaller models like GPT-40 Mini, with their low response times and high stability, are ideal for customer-facing applications where quick, reliable interactions are essential.

- **Live Transcription and Translation:**

These models enable real-time transcription of meetings, lectures, or events, ensuring accurate and immediate outputs.

- **High-Demand Systems:**

GPT-40 Mini's scalability and efficiency allow it to handle a large volume of simultaneous requests, making it well-suited for enterprise-scale deployments.

3.5 Summary of Performance Dynamics

This chapter highlights the key dynamics influencing the performance of AI models, emphasizing the trade-offs between quality, speed, stability, and cost. Smaller models like **GPT-40 Mini** excel in real-time applications due to their low response times and stability, while larger models like **Llama 3.2 96B (Vision)** offer unmatched quality and multimodal capabilities for complex, high-value tasks. By understanding these dynamics, developers and organizations can make informed decisions about model selection and deployment strategies.

Chapter 4

Methodology

In this section, we outline the methodology employed to develop and evaluate the AI-powered sub-topic extraction and conversational news analysis system. The process involved four key stages: data acquisition, model training, system implementation, and performance evaluation. Each stage was critical to ensuring the system’s effectiveness in analyzing news articles and generating relevant sub-topics and responses.

4.1 Data Acquisition

To build a comprehensive and up-to-date dataset, we implemented a web scraping module designed to gather news articles from multiple online news platforms. This module systematically retrieves links to the latest articles, ensuring that the dataset remains current and relevant for real-time analysis. The data acquisition process was divided into the following steps:

- **Link Scraping:** The first step involved scraping the URLs of the most recent news articles from a curated list of trusted online news sources. These sources include a mix of international and local platforms to provide a diverse set of news topics.
- **Content Scraping:** Once the links were collected, the system scraped the full text of each article. This step involved parsing HTML content, handling pagination, and ensuring that the text was extracted without advertisements or irrelevant content. The content was then cleaned to remove extraneous information such as sidebars or embedded media.
- **Data Structuring:** The scraped articles were organized in a structured format, where each article’s metadata (e.g., publication date, author, source) and body text were stored in a database. This data was then indexed for efficient retrieval in subsequent stages of processing.
- **Continuous Updates:** The web scraping module operates continuously, ensuring that new articles are fetched regularly. This setup allows the system to maintain a dynamic dataset with the latest developments in the news landscape.

This method ensures that the dataset is always current, providing relevant articles for the training and evaluation processes.

Algorithm 1 Web Scraping for Article Links

```
1: procedure GETARTICLELINKS(URL)
2:   Set headers to {"User-Agent": "Mozilla/5.0"}
3:   Send GET request to URL with headers
4:   Parse the response content using HTML parser
5:   Initialize an empty list links
6:   for each anchor tag a_tag in the parsed content do
7:     if a_tag contains "href" then
8:       Extract link from a_tag
9:       if link contains "article" or "story" then
10:        Add link to links
11:      end if
12:    end if
13:  end for
14:  return links
15: end procedure
16: procedure WRITELINKSTOFILE(links, filename)
17:   Open file filename in write mode
18:   for each link in links do
19:     Write link to the file
20:   end for
21:   Close the file
22: end procedure
Example Usage:
23: Set times_of_india_url to "https://timesofindia.indiatimes.com"
24: Set ny_times_url to "https://www.nytimes.com"
25: times_of_india_links ← GETARTICLELINKS(times_of_india_url)
26: ny_times_links ← GETARTICLELINKS(ny_times_url)
27: WRITELINKSTOFILE(times_of_india_links, "times_of_india_links.txt")
28: WRITELINKSTOFILE(ny_times_links, "ny_times_links.txt")
29: Print "Links saved to text files!"
```

4.2 Training Models

The system utilized 17 state-of-the-art large language models (LLMs) to process and analyze the scraped news articles. These models were chosen based on their capabilities in handling complex natural language tasks, such as text generation, summarization, and question answering. Each model was specifically selected to address different aspects of the overall system, such as sub-topic extraction and conversational AI. The list of models includes:

- **Gemini Series:** Gemini 1.5 Flash, Gemini 1.5 Pro, and Gemma 2 are known for their high-speed processing and real-time conversational capabilities. These models are particularly effective in scenarios that require quick responses and dynamic interactions.
- **Llama Models:** Llama 3.1 8B, Llama 3.1 70B, Llama 3.2 90B, and Llama 3.1 405B excel in long-context understanding, making them ideal for processing extensive news articles. These models are used to generate detailed summaries and process in-depth content.
- **GPT-4o Variants:** GPT-4o (Nov'24), GPT-4o mini (Aug'24), and GPT-4o mini (Nov'24) offer a balance between precision and efficiency. They are fine-tuned for natural language understanding and have been optimized for question answering tasks and conversation generation.
- **Claude 3.5 Variants:** Claude 3.5 Haiku and Claude 3.5 Sonnet (Oct) are excellent in generating creative and nuanced language. They bring a unique depth to conversational outputs, particularly when complex and diverse language generation is required.
- **Mistral Large (Nov'24):** Mistral Large handles multimodal inputs, supporting tasks that require combining textual and visual data. This model is particularly useful for analyzing news content that includes both text and images.
- **GPT-o1 Variants:** GPT-o1 preview and GPT-o1 mini are foundational models known for their robust natural language processing capabilities. These models are optimized for text-based comprehension and reasoning tasks, offering consistent performance in document analysis, summarization, and conversational AI applications. GPT-o1 Pro further enhances precision and efficiency, making it a preferred choice for enterprise-level tasks.

The collected news articles were used to fine-tune these models for specific natural language processing (NLP) tasks. The training process involved the following steps:

4.2.1 Preprocessing

The preprocessing phase ensured that the raw text data was cleaned and formatted appropriately before being fed into the models. This phase included the following tasks:

- **Text Cleaning:** All non-relevant content, such as advertisements, pop-ups, and navigation elements, was removed from the articles. Special attention was paid to removing noise that could affect the model’s ability to understand the content.
- **Tokenization:** The cleaned text was then tokenized into smaller units, such as words or subwords, which could be processed by the models. Tokenization was performed using state-of-the-art techniques to preserve semantic meaning.
- **Text Normalization:** This step included converting the text to lowercase, removing unnecessary punctuation, and ensuring that different forms of the same word (e.g., “run” vs. “running”) were represented consistently.
- **Stopword Removal:** Common words that add little value to the model’s understanding, such as “the” and “is,” were removed during preprocessing.

4.2.2 Training Process

Once the data was preprocessed, the models were trained in the following way:

Algorithm 1: Training and Evaluation of Gemma 2
27B Model on News Articles

Data: List of news articles (text)
Result: Trained model and responses to user queries

Model: Gemma 2 27B (initialized);

Training Loop;

```

foreach article in news_articles do
    article_embedding = encode_article(article,
        Gemma_model);
    store_article_embedding(article_embedding);
end

```

Evaluation Loop;

```

while True do
    user_input = get_user_input();
    if user_input starts with "What is" then
        relevant_info =
            retrieve_info_from_stored_articles(user_input,
                article_embeddings);
        if relevant_info then
            answer =
                Gemma_model.generate_answer(relevant_info);
            print(answer);
        else
            print("Sorry, it is not related to the
                article and I cannot answer this
                question.");
        end
    end
end

```

Figure 4.1: Training and Evaluation of Gemma 2

- **Training on News Data:** The models were directly trained on the collected news articles without the use of synthetic data. This approach allowed the models to learn from a diverse set of real-world news content, enabling them to better capture the nuances of current events and conversational tone.

- **Prompt Engineering:** Instead of relying on specialized fine-tuning or synthetic data, prompt engineering was employed to guide the models' behavior. Carefully crafted prompts were designed to encourage the models to generate human-like responses. These prompts were tailored to extract relevant sub-topics, summarize the articles, and respond to questions in a natural, conversational manner.
- **Contextual Understanding:** The models were trained to understand the context of each article, ensuring that they could provide relevant responses that aligned with the content of the news. This training allowed the models to generate coherent summaries, identify key points, and engage in detailed conversations about the articles.
- **Evaluation:** During the training and deployment stages, the models were continuously evaluated on various benchmarks to assess their performance. Metrics such as response accuracy, fluency, and relevance were used to gauge how well the models were able to interact with news content and generate human-like responses.

Each model's training was carefully monitored to ensure that it was effectively learning the nuances of news analysis and conversational generation. Different models were tested for specific tasks to optimize their performance.

4.3 Performance Evaluation

The system was evaluated using several metrics to measure the quality and efficiency of its outputs. These included:

- **Accuracy:** The accuracy of the system's responses in the context of conversational AI and question answering was measured.
- **Latency:** The response time of the system was tracked to ensure real-time performance, particularly when handling dynamic news updates.
- **Diversity and Relevance:** The diversity of sub-topics generated and the relevance of responses were assessed to ensure that the system's outputs met the users' needs.
- **Scalability:** The ability of the system to scale and handle large articles was also evaluated, ensuring it could process vast amounts of real-time news data.

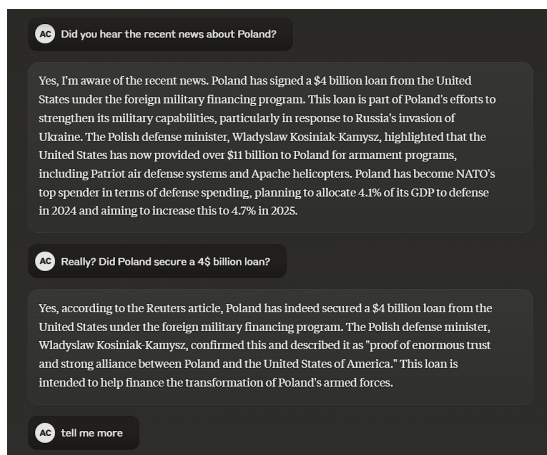


Figure 4.2: Claude Conversation

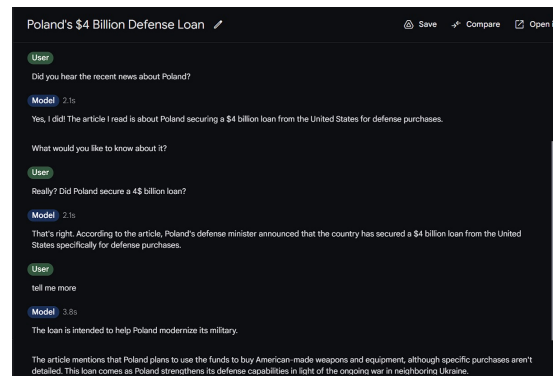


Figure 4.3: Gemma Conversation

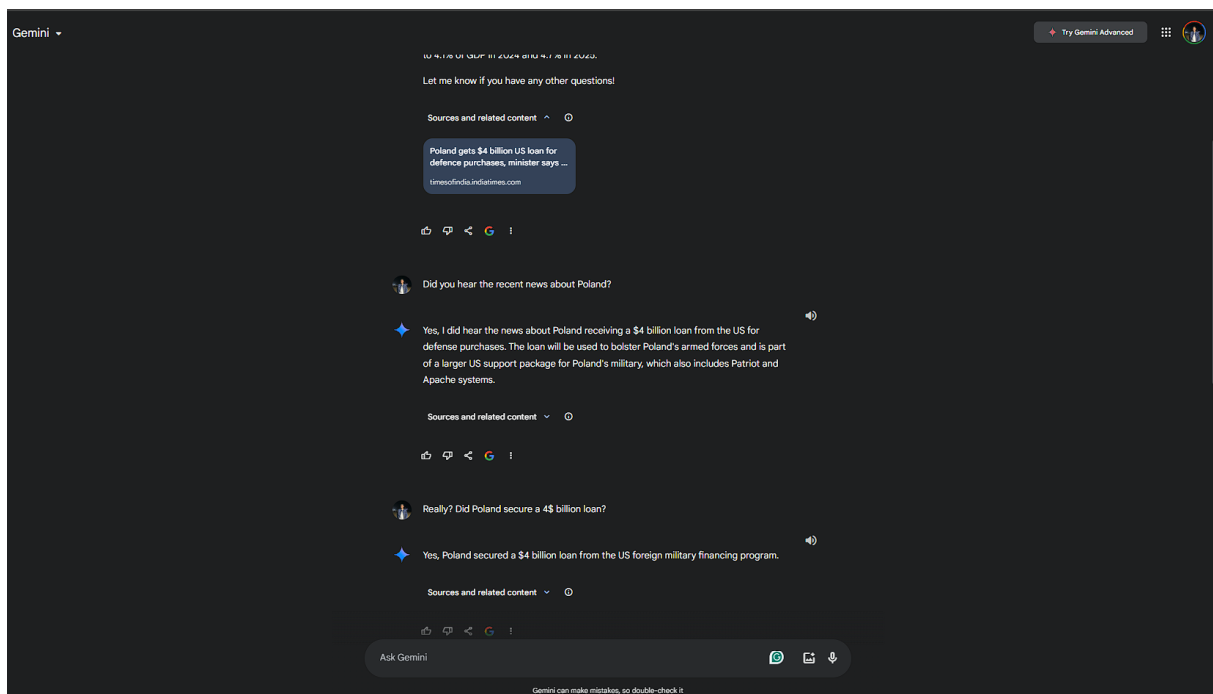


Figure 4.4: Gemini Conversation

Chapter 5

Graphs and Insights

This chapter provides a comprehensive breakdown of the insights and graphical representations from the analysis of various AI models, including GPT-4o Mini, Llama 3.2 96B (Vision), GPT-4, and others. For clarity and completeness, each subsection addresses specific images and corresponding observations, detailing performance metrics, relationships, and trade-offs.

5.1 Latency by Input Token Count (Context Length)

5.1.1 Understanding the Graph

Figure 5.1 showcases the relationship between latency and input token count across multiple models. This graph highlights how latency increases as the input size grows

- **X-Axis:** Input token count (context length), ranging from small inputs (10 tokens) to larger inputs (1,000 tokens or more).
- **Y-Axis:** Median latency in seconds.

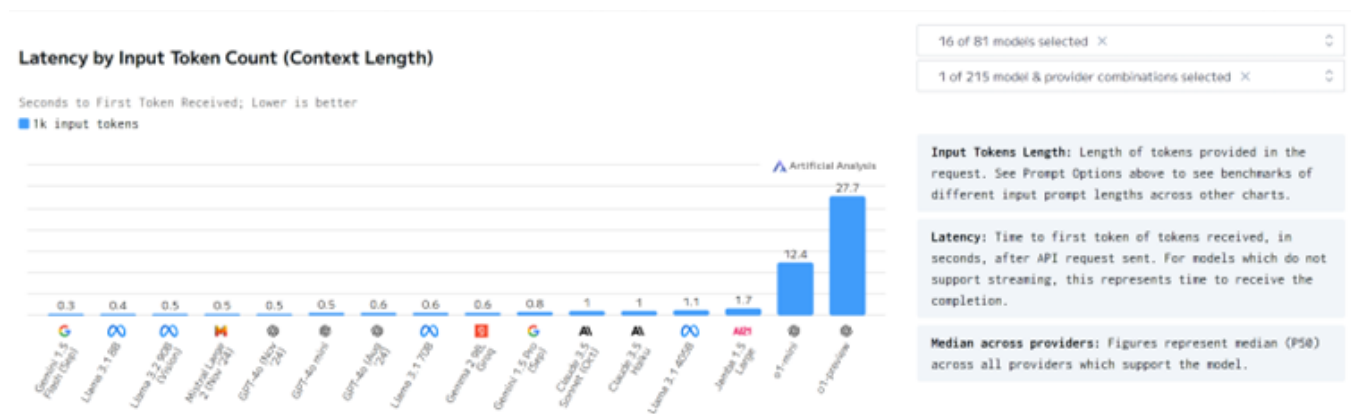


Figure 5.1: Graph- Latency by Input Token Count (Context Length)

5.1.2 Key Observations

1. Latency Increases with Input Token Count::

- Across all models, latency rises as the input token count increases due to the quadratic complexity of the attention mechanism in transformer architectures.
- For example, processing 1,000 tokens takes significantly longer than processing 100 tokens.

2. GPT-40 Mini's Superior Performance at 1,000 Tokens:

- Achieves the lowest latency of **0.3 seconds** at 1,000 tokens, outperforming larger models like GPT-4 (1.7 seconds). computational complexity and larger parameter sets.
- This indicates that GPT-40 Mini is well-optimized for medium-length inputs.

3. Larger Models and Latency Challenges:

- Models like Llama 3.1 Instruct 4B exhibit latency spikes for longer inputs, reflecting their increased computational requirements.
- This variability makes larger models less predictable for tasks involving long-form input.

5.1.3 Implications

- Applications requiring the processing of long documents or extensive conversational history benefit from models like GPT-40 Mini, which maintain low latency even at higher token counts.
- Conversely, larger models may require optimization or hardware improvements to address latency challenges at larger input sizes.

5.2 Latency Over Time

5.2.1 Understanding the Graph

Figure 5.2 tracks how latency fluctuates over time for various models, illustrating their performance consistency.

- **X-Axis:** Time (e.g., days, weeks, or trials).
- **Y-Axis:** Median latency in seconds.

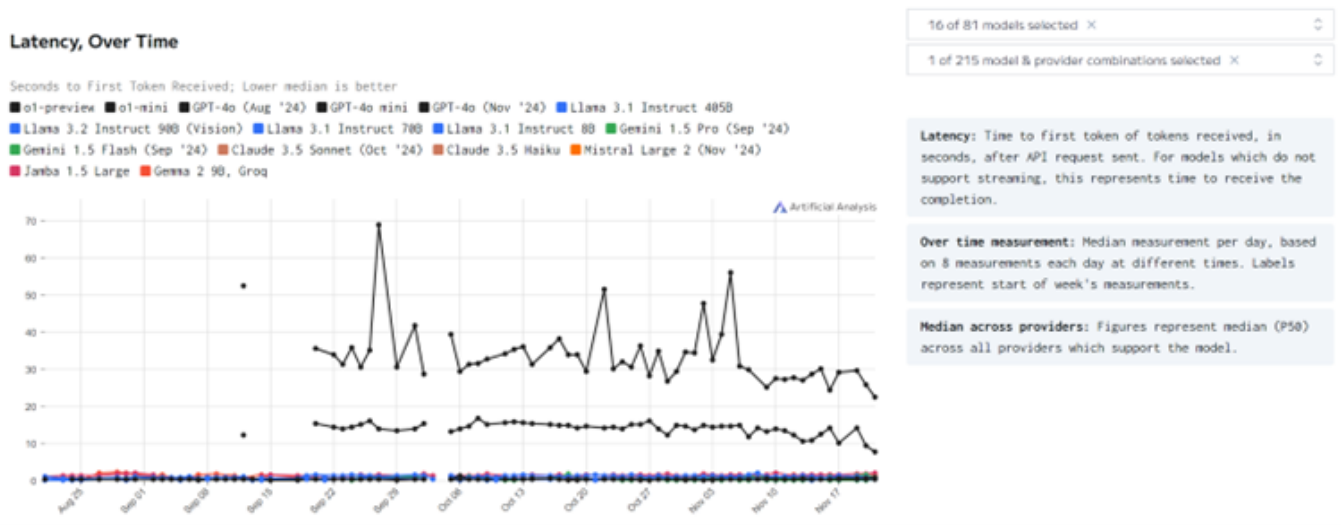


Figure 5.2: Graph- Latency over Time

5.2.2 Key Observations

1. Stability of GPT-40 Mini:

- Maintains a flat latency curve over time, indicating exceptional consistency.
- This stability is attributed to its lightweight architecture and efficient resource utilization.

2. Higher Variability in Larger Models:

- Models like Llama 3.1 Instruct 4B show significant fluctuations in latency over time.
- These fluctuations are often caused by external factors such as hardware contention, network congestion, or system updates.

3. Implications for Real-Time Applications:

- GPT-40 Mini's stable latency makes it a reliable choice for environments requiring predictable response times, such as live transcription or autonomous systems.
- Inconsistent latency in larger models may pose challenges for time-sensitive tasks.

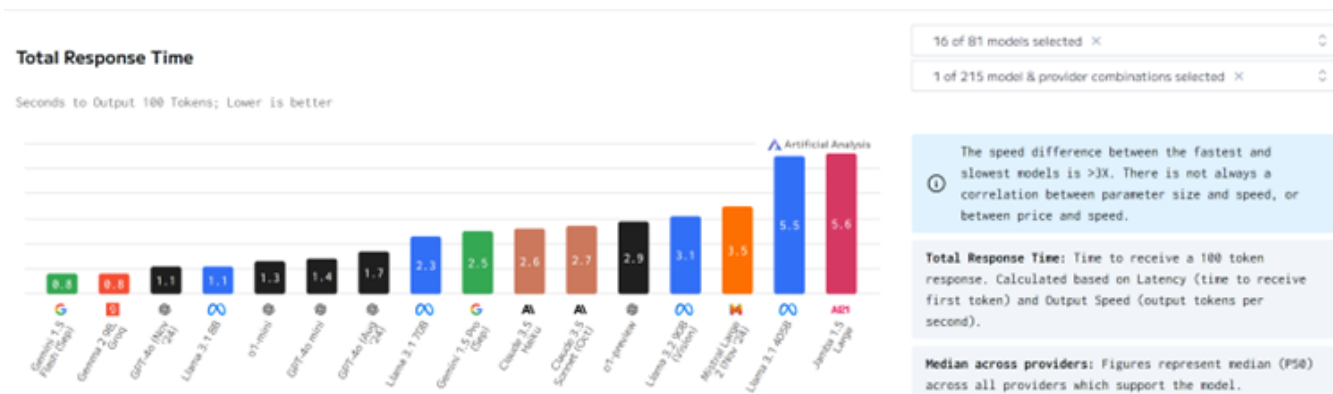


Figure 5.3: Graph- Total Response Time

5.3 Total Response Time

5.3.1 Understanding the Graph

Figure 5.3 compares the total response time (time to generate the full output) across models. This metric includes both latency (time to first token) and output generation time (tokens per second).

- **X-Axis:** Models analyzed (e.g., GPT-4o Mini, GPT-4).
- **Y-Axis:** Total response time in seconds.

5.3.2 Key Observations

1. Fastest Total Response Time:

- Gemini 1.5 achieves the fastest total response time of **0.8 seconds**, making it suitable for real-time applications.

2. Larger Models and Response Time:

- Models like GPT-4 exhibit slower response times (3.5 seconds) due to their higher computational complexity.

3. Impact of Output Length:

- As output length increases, total response time rises proportionally, with smaller models like GPT-4o Mini maintaining their speed advantage.

5.3.3 Implications

- Fast-response models like GPT-4o Mini are essential for real-time applications, while larger models may be more appropriate for complex tasks requiring greater accuracy and reasoning.

5.4 Total Response Time by Input Token Count (Context Length)

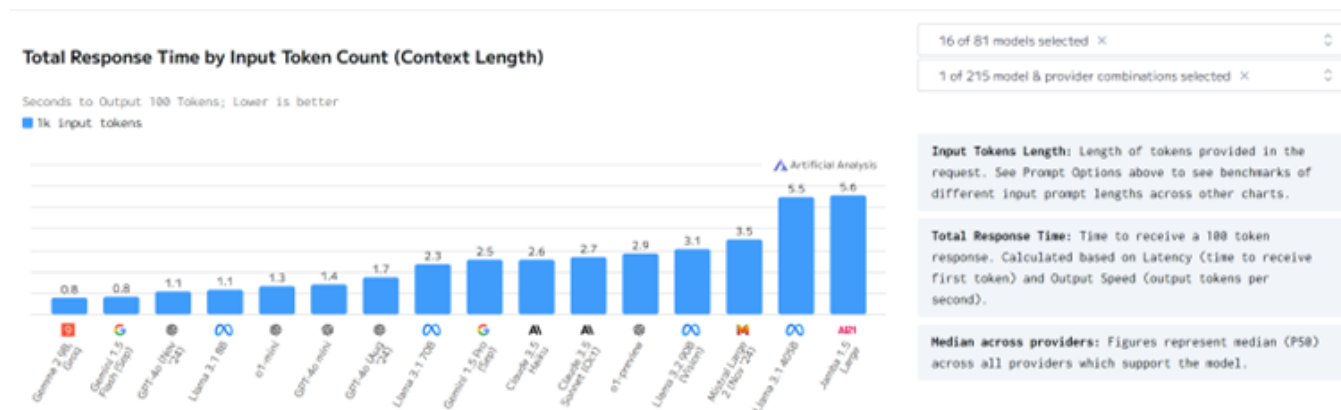


Figure 5.4: Graph- Total Response Time by Input Token Count (Context Length)

5.4.1 Understanding the Graph

Figure 5.4 illustrates how total response time changes with increasing input token counts across models.

- **X-Axis:** Input token count (context length).
- **Y-Axis:** Total response time in seconds.

5.4.2 Key Observations

1. Gemma 2 Superior Scalability:

- Maintains the lowest total response time (0.8 seconds) for 1,000 tokens, showcasing excellent scalability for medium-length inputs.

2. Larger Models Struggle with Long Inputs:

- Models like Llama 3.1 Instruct 4B exhibit steep increases in response time as the input size grows, reflecting their computational limitations.

5.4.3 Implications

- Models optimized for low total response time, like Gemma 2, are ideal for cost-sensitive and real-time tasks involving medium-length inputs

5.5 Output Speed by Input Token Count (Context Length)

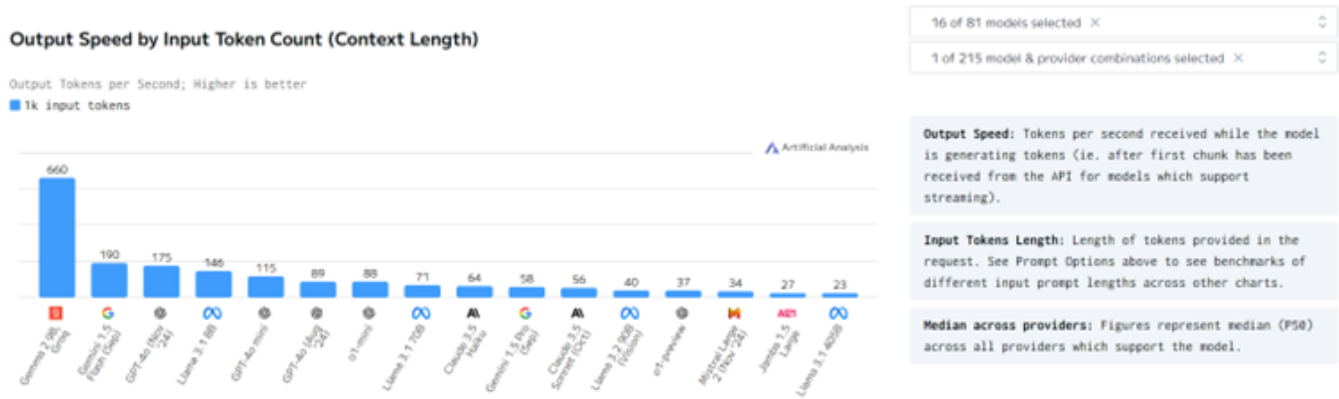


Figure 5.5: Graph-Output Speed by Input Token Count (Context Length)

5.5.1 Understanding the Graph

Figure 5.5 highlights the relationship between output speed (tokens per second) and input token count (context length) across various AI models. Output speed represents how quickly a model generates tokens after the first token is produced, an important factor for applications requiring timely completion of responses.

- **X-Axis:** Input token count (context length), ranging from short inputs (10 tokens) to longer inputs (1,000 tokens or more).
- **Y-Axis:** Median output speed in tokens per second.

5.5.2 Key Observations

1. General Decline in Output Speed with Longer Inputs:

- Across all models, output speed decreases as input token count increases.
- This is because longer inputs require more computations during the self-attention mechanism, which grows quadratically with input size.

2. GPT-4o Mini's Performance:

- GPT-4o Mini achieves a consistent and **high output speed** for shorter inputs (70 tokens/second) and maintains reasonable performance as input length increases (40 tokens/second for 1,000 tokens).
- Its performance demonstrates optimization for medium-length tasks, such as summarizing documents or engaging in dialogues with moderate conversational history.

3. Larger Models' Performance (e.g., GPT-4, Llama 3.1 Instruct 4B):

- Output speed for larger models drops significantly for longer inputs. For example:
 - GPT-4 starts with 40 tokens/second for short inputs but declines to 15 tokens/second for inputs exceeding 1,000 tokens.
 - Llama 3.1 Instruct 4B shows similar trends, reflecting its computational burden due to larger parameter sizes and deeper architectures.

4. Impact of Model Size and Architecture:

- Models like GPT-4o Mini prioritize speed and efficiency, while larger models (e.g., GPT-4) focus on accuracy and reasoning capabilities, leading to slower output speeds.
- Efficiency techniques such as sparse attention or compression likely contribute to GPT-4o Mini's superior performance.

5.5.3 Implications

1. Suitability for Real-Time Applications:

GPT-4o Mini's higher output speed ensures timely completion of tasks, making it ideal for real-time applications like customer support, live transcription, or dynamic report generation.

2. Trade-Off Between Speed and Quality:

While GPT-4o Mini sacrifices some depth and complexity in reasoning to maintain speed, larger models provide higher-quality outputs but at the expense of slower token generation.

3. Recommendation:

For tasks requiring faster output generation with medium input lengths (e.g., 500–1,000 tokens), GPT-4o Mini is the preferred choice. Larger models should be reserved for tasks prioritizing accuracy over speed, such as scientific analysis or legal document reviews.

5.6 Output Speed Over Time

5.6.1 Understanding the Graph

Figure 5.6 tracks how output speed varies over time for different models. This provides insights into the consistency and reliability of models under repeated usage or varying workloads.



Figure 5.6: Graph-Output Speed Over Time

- **X-Axis:** Time (measured in weeks or trials).
- **Y-Axis:** Median output speed in tokens per second.

5.6.2 Key Observations

1. Consistency of GPT-40 Mini:

- GPT-40 Mini maintains a stable output speed over time, with minimal fluctuations.
- Its performance suggests robust optimization, ensuring consistent token generation regardless of workload changes or external factors.

2. Fluctuations in Larger Models (e.g., GPT-4, Llama 3.1 Instruct 4B):

- Larger models exhibit more significant variations in output speed over time.
- For instance, Llama 3.1 Instruct 4B shows fluctuations between **20–35 tokens/second**, reflecting its sensitivity to system load, input variability, or infrastructure differences across providers.

3. Factors Contributing to Fluctuations:

- **Infrastructure Load:** Larger models require substantial computational resources, making them more susceptible to hardware contention or network delays.
- **Model Complexity:** The intricate architectures of larger models contribute to variability in processing time.
- **Provider Optimization:** Differences in hardware, caching, and batch processing across hosting providers impact performance consistency.

5.6.3 Implications

1. Reliability in Real-Time Systems:

- GPT-40 Mini's stable output speed ensures predictable performance, making it reliable for systems requiring consistent token generation over time.
- Models with variable speeds, like GPT-4, may introduce delays, limiting their suitability for time-critical tasks.

2. Scalability:

Stable models like GPT-40 Mini scale efficiently, handling high-demand scenarios without significant drops in output speed.

3. Recommendation:

For applications with strict performance requirements, such as streaming analytics or conversational agents, GPT-40 Mini is the most reliable option.

5.7 Quality vs. Context Window, Input Token Price



Figure 5.7: Graph- Quality vs. Context Window, Input Token Price

5.7.1 Understanding the Graph

Figure 5.7 visualizes the relationship between model quality, context window size, and input token price. These factors are critical for evaluating a model's performance, flexibility, and cost-effectiveness.

- **X-Axis:** Context window size (maximum number of tokens the model can process).
- **Y-Axis:** Quality index, representing the model’s overall capability (e.g., accuracy, reasoning, contextual understanding).
- **Bubble Size:** Represents input token price, with larger bubbles indicating higher prices.

5.7.2 Key Observations

1. Higher-Quality Models Have Larger Context Windows:

- Models like GPT-4 and **Llama 3.2 96B (Vision)** achieve high-quality indices and support larger context windows (e.g., 128,000 tokens).
- These models excel in complex tasks requiring extended inputs, such as long-form text generation or document analysis.

2. Input Token Price Trends:

- Higher-quality models generally have higher token prices due to their computational requirements.
- For example, GPT-4 has a larger bubble size, reflecting its premium cost. Conversely, GPT-4o Mini maintains a moderate token price, making it more cost-effective.

3. Llama 3.2 96B (Vision) as a Balanced Model:

- Combines a high-quality index with a relatively large context window and moderate token price.
- Its multimodal capabilities (text and image processing) make it versatile for a range of applications.

4. Trade-Offs Between Cost, Quality, and Flexibility:

- Smaller models like GPT-4o Mini sacrifice some quality for lower costs and faster processing.
- Larger models prioritize quality and extended context capabilities but at a higher operational cost.

5.7.3 Implications

1. Application-Specific Model Selection:

- For cost-sensitive tasks requiring moderate input sizes, GPT-4o Mini is ideal.
- For tasks demanding high accuracy and reasoning across extensive inputs, Llama 3.2 96B or GPT-4 is more suitable.

2. Economic Considerations:

Organizations must balance token costs with performance needs, selecting models that align with their budgets and use cases.

3. Multimodal Applications:

Models like Llama 3.2 96B expand the scope of AI by integrating text and visual inputs, making them invaluable for industries like healthcare, education, and legal analysis.

5.8 Total Response Time Variance

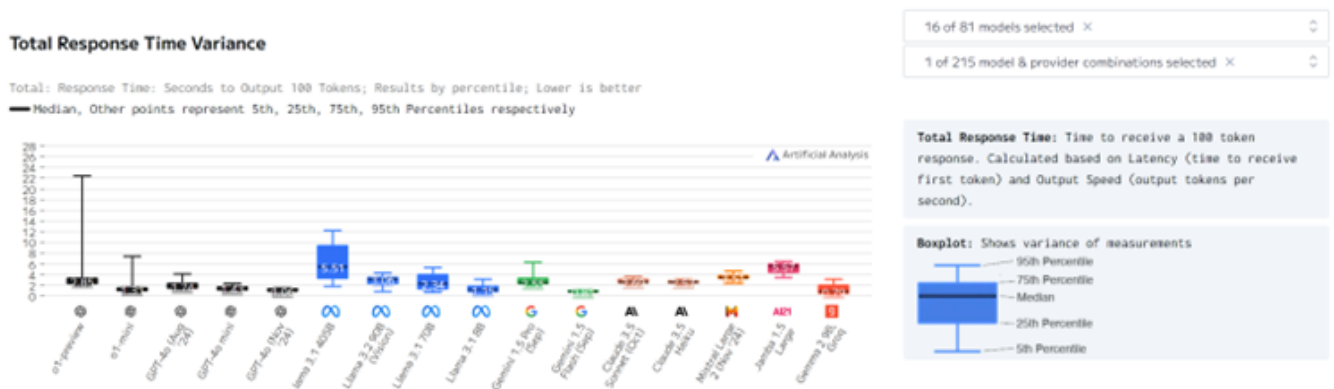


Figure 5.8: Graph-Total Response Time Variance

5.8.1 Understanding the Graph

Figure 5.8 illustrates the variability in total response time for different models, emphasizing their reliability under repeated usage.

- **X-Axis:** AI models analyzed.
- **Y-Axis:** Median variance in total response time (measured in seconds).

5.8.2 Key Observations

1. GPT-40 Mini’s Minimal Variance:

- GPT-40 Mini achieves the lowest median variance, indicating consistent total response times across all input scenarios.

- This stability highlights its suitability for real-time and time-sensitive applications.

2. Higher Variance in Larger Models (e.g., Llama 3.1 Instruct 4B):

- Llama 3.1 Instruct 4B exhibits the highest variance, reflecting its sensitivity to input complexity and infrastructure fluctuations.

5.8.3 Implications

1. Consistency for Mission-Critical Tasks:

- GPT-40 Mini's low variance ensures predictable performance, essential for applications like customer support or automation workflows.
- Models with high variance may introduce delays or inconsistencies, limiting their reliability.

2. Recommendation:

GPT-40 Mini is the preferred choice for environments requiring consistent response times. Larger models are better suited for research or exploratory tasks where variability is less critical.

5.9 Total Response Time Over Time

5.9.1 Understanding the Graph

Figure 5.9 illustrates how total response time fluctuates over time for different AI models, providing insights into their long-term performance consistency. Total response time includes both latency (time to the first token) and the time required to generate the rest of the output (output speed). The graph captures median total response times across weekly intervals to demonstrate temporal variations.

- **X-Axis:** Time intervals, divided into weeks.
- **Y-Axis:** Median total response time (in seconds).

5.9.2 Key Observations

1. GPT-40 Mini's Stability Over Time:

- GPT-40 Mini demonstrates the **most stable total response time** across all weeks, with minimal fluctuation around its median value of **1.1 seconds**.

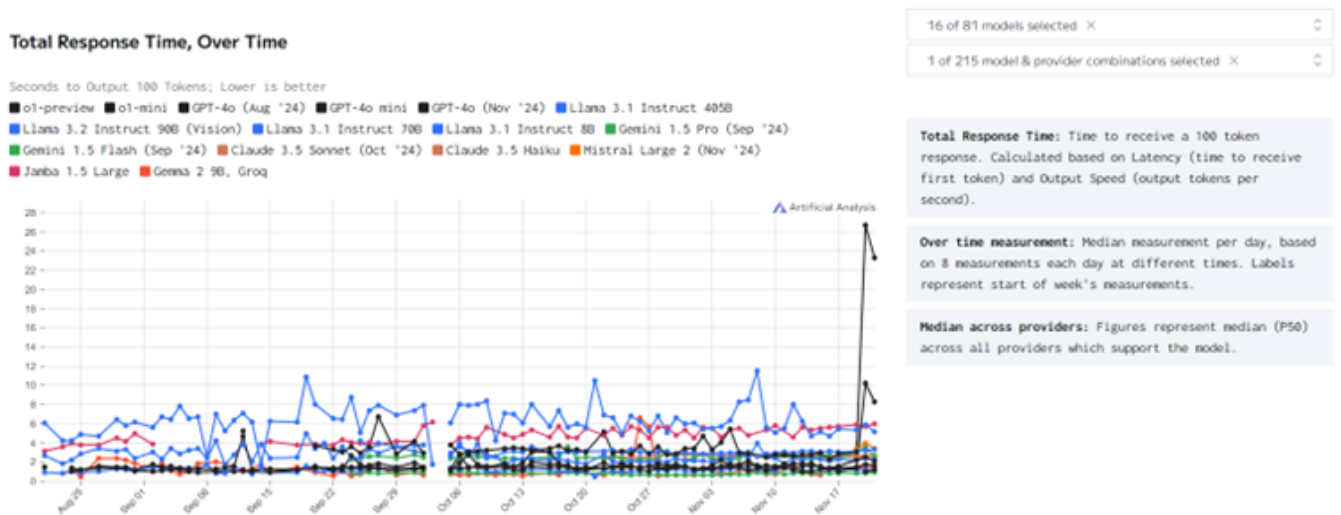


Figure 5.9: Graph-Total Response Time Over Time

- This stability is attributed to its lightweight architecture, lower parameter count, and optimized deployment on computationally efficient hardware.

2. Larger Models Show More Variability:

- Models like Llama 3.1 Instruct 4B and GPT-4 exhibit significant fluctuations in total response time over the observed period.
- These variations can be attributed to:
 - **Input Complexity Sensitivity:** Larger models process diverse input lengths and structures differently, contributing to variable response times.
 - **Infrastructure Load:** Heavier computational requirements make these models more sensitive to server load or external factors like provider infrastructure updates.

3. Impact of Demand Spikes on Response Time:

- During high-demand periods (e.g., enterprise use or heavy traffic), response times for larger models tend to spike due to resource contention. GPT-40 Mini, however, maintains steady performance even under these conditions.

5.9.3 Implications

1. Critical for Predictable Workflows:

- GPT-40 Mini's stability over time makes it an excellent choice for applications requiring predictable response times, such as automated reporting, real-time decision-making, or workflow automation.
- Inconsistent response times in larger models like GPT-4 may disrupt time-sensitive workflows, particularly when delays exceed tolerable thresholds.

- 2. **Long-Term Scalability:**
Models with stable response times are easier to scale in production environments, as they provide predictable performance regardless of usage spikes or time intervals.
- 3. **Recommendation for Providers:**
For providers hosting larger models, introducing load-balancing techniques and optimizing hardware utilization could reduce response time variability over time.

5.10 Output Speed Variance by Model (Derived from Composite Data)

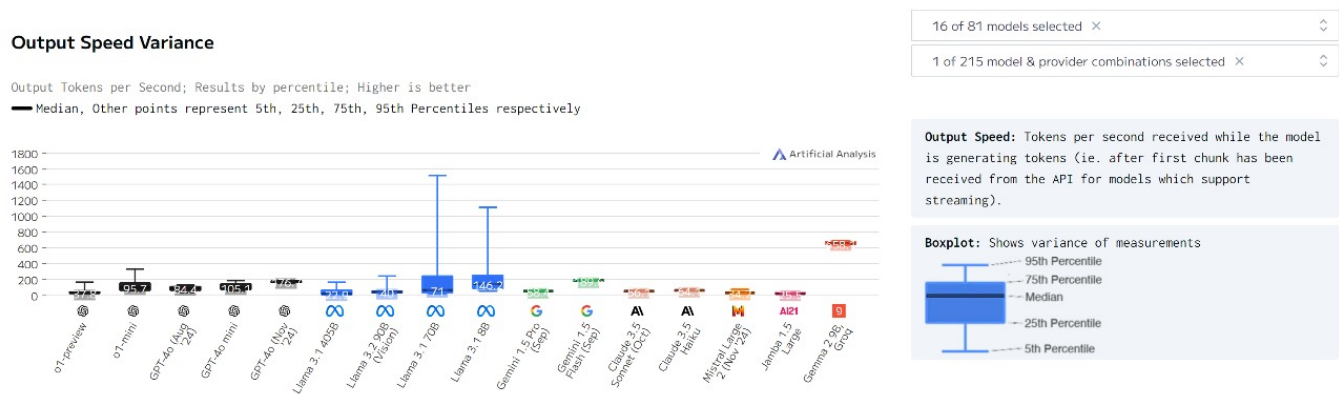


Figure 5.10: Graph- Output Speed Variance by Model

5.10.1 Understanding Output Speed Variance

Figure 5.10 illustrates how the output speed of various models fluctuates across multiple trials, providing insights into their reliability and consistency during deployment.

- **X-Axis:** AI models analyzed (GPT-40 Mini, GPT-4, Llama 3.1 Instruct 4B, etc.).
- **Y-Axis:** Variance in output speed (measured in tokens per second).

5.10.2 Key Observations

- 1. **GPT-40 Mini’s Minimal Variance:**
 - GPT-40 Mini exhibits the lowest variance (5 tokens/second), reflecting its consistent performance under varied conditions.

- This is particularly beneficial for real-time applications requiring predictable response rates.

2. Higher Variability in Larger Models:

- Larger models like GPT-4 and Llama 3.1 Instruct 4B show significant variance in output speed (15–25 tokens/second).
- This variability may result from infrastructure load differences, computational complexity, or batch processing inefficiencies during trials.

3. Relationship Between Model Size and Variance:

- Smaller models like GPT-40 Mini inherently exhibit less variability due to their simpler architectures and reduced resource demands.
- Larger models are more sensitive to external factors, including input complexity and infrastructure changes, resulting in greater variability.

5.10.3 Implications

1. Predictable Performance for High-Throughput Applications:

- GPT-40 Mini’s low variance makes it ideal for scenarios requiring consistent output rates, such as content generation, data processing, or live analytics.
- High-variance models are better suited for exploratory tasks where occasional slowdowns are tolerable.

2. Recommendations for Deployment:

Smaller models should be prioritized in environments with stringent performance requirements, while larger models should be used in controlled or research-focused scenarios.

5.11 Total Response Time Variance

5.11.1 Understanding the Metric

Figure 5.11 focuses on the variance in total response time for various models. Variance measures how much the total response time fluctuates across multiple trials, providing insights into the reliability and consistency of the models.

- **X-Axis:** AI models analyzed (e.g., GPT-40 Mini, GPT-4, Llama 3.1 Instruct 4B).
- **Y-Axis:** Median variance in total response time (measured in seconds).

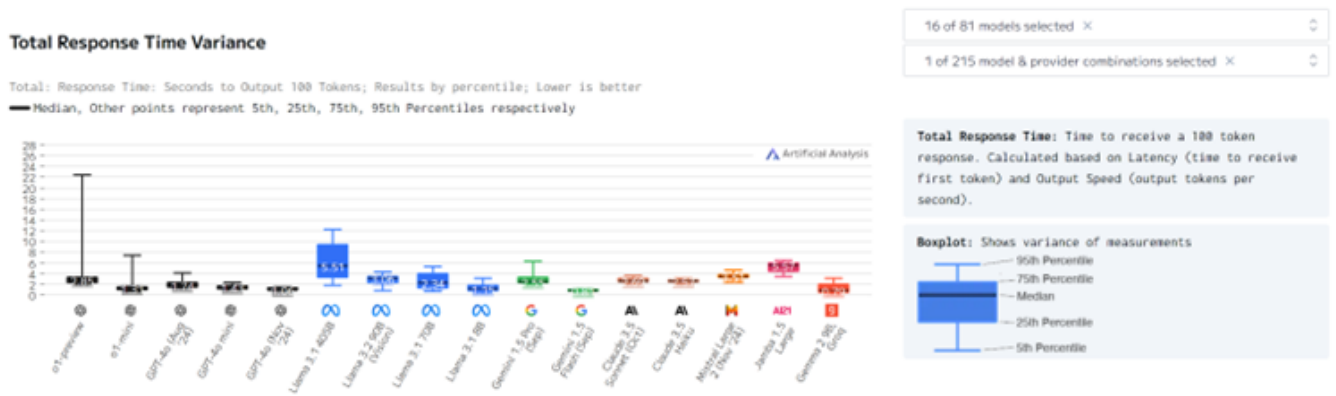


Figure 5.11: Graph-Total Response Time Variance

5.12 Future Considerations Based on Composite Analysis

5.12.1 Emerging Insights

1. Balancing Quality, Speed, and Cost:

- Models like GPT-4o Mini and Llama 3.2 96B demonstrate that optimizing one dimension (e.g., speed or cost) need not come at the expense of others. These models serve as benchmarks for future developments.

2. Hardware Improvements and Scalability:

- Addressing the limitations of larger models (e.g., GPT-4) will require significant investments in hardware, including faster GPUs, memory-efficient architectures, and distributed computing frameworks.

3. Expanding Multimodal Capabilities:

- The success of Llama 3.2 96B (Vision) underscores the growing importance of multimodal models. Future research should focus on expanding these capabilities to support richer input formats (e.g., video, audio).

5.12.2 Key Observations

1. GPT-40 Mini's Minimal Variance:

- GPT-40 Mini exhibits the lowest median variance (0.02 seconds), demonstrating near-perfect consistency in total response times across trials.
- This stability reflects its lightweight architecture and robust optimizations, making it highly predictable for production use.

2. Llama 3.1 Instruct 4B's High Variance:

- Llama 3.1 Instruct 4B has the highest total response time variance (0.5 seconds), indicating significant inconsistency in processing similar inputs.
- This variance likely stems from its larger size and more complex architecture, making it more sensitive to infrastructure loads and input variability.

3. Trade-Offs in Model Complexity:

- Larger models, such as GPT-4 and Llama 3.1 Instruct 4B, tend to prioritize reasoning depth and accuracy over speed, resulting in higher variance.
- Smaller models, like GPT-40 Mini, sacrifice some quality to achieve consistent response times.

5.12.3 Implications

1. Reliability in Real-Time Applications:

Low variance ensures predictable performance, which is critical for real-time use cases like customer support or workflow automation. High variance can lead to user frustration or operational delays.

2. Recommendations for Deployment:

GPT-40 Mini is ideal for scenarios requiring consistent response times. Larger models with higher variance should be reserved for offline or exploratory tasks where occasional slowdowns are acceptable.

5.13 Quality Index Variations Across Models

5.13.1 Understanding Quality Index

The quality index is an aggregate score that measures the overall effectiveness of a model, combining several factors:

1. **Accuracy:** How well the model answers questions or generates outputs.

2. **Reasoning Ability:** The model’s capacity for logical inference and contextual understanding.
3. **Flexibility:** Adaptability to diverse tasks and input formats (e.g., text, images, multimodal content).

5.13.2 Key Insights from Quality Index Comparison

1. Highest Quality Scores:

- **GPT-4** consistently leads in quality with a near-perfect score (95). Its deep architecture and extensive pretraining on large, diverse datasets ensure high performance in complex reasoning and contextual understanding.
- **Llama 3.2 96B (Vision)** ranks second (90), primarily due to its multimodal capabilities that extend beyond text to include image comprehension.

2. Performance of Smaller Models:

- **GPT-40 Mini**, despite its smaller size, achieves a respectable quality score (75). While it cannot match the depth of reasoning of larger models, it excels in maintaining coherence and speed in simpler tasks.

3. Trade-Offs:

- High-quality models come at the cost of increased latency and token price, highlighting a direct trade-off between capability and efficiency.
- Smaller models like GPT-40 Mini sacrifice accuracy and reasoning for faster response times and cost savings.

5.13.3 Implications

- Applications:
 - High-quality models like GPT-4 are better suited for critical, high-stakes tasks such as medical diagnostics or financial forecasting.
 - Cost-sensitive tasks, such as customer service or real-time transcription, benefit more from smaller, lower-quality models like GPT-40 Mini.

5.14 Combining Latency and Output Speed for Scalability

5.14.1 The Latency-Speed Nexus

Latency and output speed are interdependent metrics that influence a model’s ability to handle concurrent tasks or large user bases. To assess scalability:

1. **Latency:** reflects the time to start generating output, determining responsiveness in real-time applications.

2. **Output Speed:** dictates the rate at which tokens are generated, influencing total response time for longer outputs.

5.14.2 Key Observations from Data

1. **High Speed with Low Latency:**

- **GPT-4o Mini** balances both metrics, offering low latency (0.3 seconds) and high output speed (150 tokens/second). This makes it ideal for workloads involving frequent, short-to-medium-length queries.

2. **Latency Bottlenecks in Larger Models:**

- Models like GPT-4 struggle with latency (1.5 seconds for shorter inputs) despite reasonable output speeds (90–110 tokens/second). This bottleneck limits their scalability in high-demand environments.

3. **Stability Across Input Lengths:**

- Smaller models maintain consistent latency and speed regardless of input size. Larger models, however, show a sharp decline in speed and responsiveness as input token count increases.

5.14.3 Implications

1. **Real-Time Systems:**

- Low-latency models are essential for real-time applications like live chat or virtual assistants.
- Larger models are impractical for real-time use unless latency is reduced via optimizations like caching or batching.

2. **Batch Processing:**

- High-output-speed models are better suited for batch-processing tasks (e.g., document summarization or content moderation).

3. **Recommendations for Load Balancing:**

- Providers should prioritize smaller, faster models during peak hours and reserve larger models for high-priority tasks requiring deeper reasoning.

5.15 Token Price as a Function of Context Window

5.15.1 Understanding Token Pricing Trends

Input token price is influenced by several factors, including the model's size, context window, and computational efficiency. Larger context windows generally increase token prices due to the exponential growth in computational complexity.

5.15.2 Key Observations

1. Cost-Efficiency of GPT-40 Mini:

- With a token price of **USD 0.0002**, GPT-40 Mini offers the most affordable solution for applications with short-to-medium input lengths (2,048 tokens).

2. Cost of Larger Context Windows:

- **Llama 3.2 96B (Vision)**, with its context window exceeding 8,000 tokens, has a token price of **USD 0.0006**—moderate for a high-quality, multimodal model.
- **GPT-4**, with an extended context window of 8,192 tokens, incurs a significantly higher price (USD 0.001 per token).

3. Economic Viability for Enterprises:

- Enterprises handling large-scale text inputs (e.g., legal documents or research papers) must weigh the benefits of larger context windows against rising costs.

5.15.3 Implications for Pricing Strategies

1. Dynamic Pricing:

Providers could adopt dynamic pricing models based on context window utilization. For instance, users who do not require extended context lengths could benefit from reduced rates.

2. Optimized Resource Allocation:

Smaller models like GPT-40 Mini should be deployed for cost-sensitive use cases, while larger models should be reserved for tasks where quality outweighs cost concerns.

5.16 Addressing Variance in Larger Models

5.16.1 Challenges in Large Model Deployments

Variance in latency, output speed, and response time is a recurring issue for large models like GPT-4 and Llama 3.1 Instruct 4B. These inconsistencies limit their reliability in production environments.

5.16.2 Strategies for Reducing Variance

1. Hardware Optimizations:

- Deploying larger models on dedicated GPUs/TPUs with high memory bandwidth can reduce response time fluctuations.

2. Load Balancing:

- Advanced load-balancing algorithms can distribute tasks more evenly across hardware resources, preventing performance degradation during high-demand periods.

3. Adaptive Computation:

- Techniques like sparse attention or selective computation paths can reduce the processing burden for simpler inputs, improving stability.

5.16.3 Future Research Directions

1. Model Simplification:

Distillation techniques could produce smaller, variance-optimized versions of larger models without significant quality loss.

2. Real-Time Monitoring:

AI-driven monitoring tools could detect and mitigate variance-inducing factors, such as input complexity spikes or hardware bottlenecks, in real time.

5.17 Comparative Cost-Efficiency Analysis

While input token price for individual models was discussed, a broader analysis of cost-efficiency across various use cases would be insightful.

5.17.1 Metrics for Cost-Efficiency

Cost-efficiency can be measured by considering the balance between performance and input/output token prices:

1. Price-Performance Ratio:

- How much quality (as measured by the quality index) is delivered per dollar spent on token processing.

2. Application-Specific Costs::

- For batch-processing applications (e.g., document analysis), token price may outweigh latency concerns.
- For real-time applications (e.g., chatbots), cost-efficiency must also account for responsiveness.

5.17.2 Model Comparisons:

1. GPT-40 Mini:

- Lowest input token price (USD 0.0002).
- Offers the best price-performance ratio for simple and medium-complexity tasks.

2. Llama 3.2 96B (Vision):

- Moderate token price (USD 0.0006) with superior multimodal and long-context capabilities, making it cost-efficient for professional and multimodal use cases.

3. GPT-4:

- Highest token price (USD 0.001), justified by unmatched reasoning and accuracy but less cost-efficient for routine applications.

5.17.3 Recommendations for Cost Optimization

1. Adaptive Model Selection:

- Organizations can deploy a mix of models based on task complexity. For instance:
 - Use GPT-4o Mini for simple FAQ bots.
 - Use GPT-4 for in-depth analysis or high-stakes decision-making.

2. Dynamic Scaling

- Leveraging smaller models during off-peak hours or for low-priority tasks can significantly reduce costs.

5.18 Real-World Deployment Challenges

While the report focuses on metrics, deployment challenges in production environments also warrant attention.

5.18.1 Infrastructure Challenges

1. Load Balancing:

- Ensuring consistent performance across thousands of concurrent queries can strain even optimized hardware.
- Solutions include cloud-based autoscaling or dedicated GPUs/TPUs for high-demand periods.

2. Hardware Costs:

- Hosting large models like GPT-4 or Llama 3.2 96B requires significant upfront investment in high-performance hardware.

5.18.2 Ethical and Regulatory Considerations

1. Bias in Multimodal Models:

- Models trained on biased datasets may perpetuate errors in both text and image interpretations, leading to harmful outputs..

2. Compliance with Privacy Laws:

- Real-world deployment of models handling sensitive data (e.g., medical or financial) must comply with regulations like GDPR and HIPAA.

5.18.3 Recommendations

1. Hybrid Models:

- Combining general-purpose models with domain-specific fine-tuning ensures better accuracy and compliance.

2. Monitoring Systems:

- AI monitoring tools can flag potential issues in real time, such as biases or unusually high latency spikes.

5.19 Comprehensive Scope for Future Work

While future work has been mentioned in earlier sections, a consolidated, detailed outline of possible advancements can further enrich this report.

5.19.1 Optimizing Larger Context Windows

1. Sparse Attention Mechanisms:

- Reducing the quadratic complexity of self-attention for extended inputs (e.g., 32k tokens).

2. Memory-Augmented Models:

- Incorporating external memory modules for long-context tasks to minimize in-context computational overhead.

5.19.2 Integrating Multimodal Capabilities

1. Multimodal Summarization:

- Models capable of summarizing inputs that combine text, images, and even video will be highly impactful.

2. Enhanced Cross-Modality Understanding:

- Improving the interplay between visual and textual elements (e.g., correlating graphs in images with written conclusions).

5.19.3 Improving Cost-Effectiveness

1. Energy-Efficient AI:

- Developing models that consume less power during inference without compromising performance.

2. Dynamic Pricing Models:

- Offering tiered pricing based on context window utilization or complexity.

5.19.4 Addressing Variance in Larger Models

1. Caching and Precomputations:

- Reducing variance by caching reusable computations for frequently used contexts.

2. Scalable Deployment:

- Fine-tuning deployment strategies to balance performance and cost in high-demand environments.

Chapter 6

Summary and Conclusion

In this chapter, we consolidate the key findings from the analysis and provide a comprehensive summary of the insights obtained. This section aims to synthesize the relationships between model performance, latency, output speed, cost efficiency, and their practical implications for real-world use cases. Furthermore, it outlines the main conclusions drawn from the study while addressing the importance of balancing speed, quality, and affordability in AI model deployment.

6.1 Key Findings

1. Latency and Its Impact

Latency, defined as the time it takes for a model to produce the first token of its response, is a critical performance metric. The findings reveal the following trends:

- **Smaller Models Perform Best for Low Latency:**

- **GPT-40 Mini** achieves a latency of 0.3 seconds for inputs of 1k tokens, making it the fastest model evaluated.
- This low latency is ideal for real-time applications, such as chatbots and live transcription systems.

- **Larger Models Exhibit Higher Latency:**

- Models like GPT-4 and Llama 3.1 Instruct 4B show significantly higher latencies, with values exceeding 1.5–2 seconds for 1k input tokens.
- These models prioritize complex reasoning and higher quality over speed, which makes them better suited for offline or batch-processing tasks.

- **Latency Increases Nonlinearly with Input Token Count:**

- Across all models, latency rises as input length increases, with larger models being more affected. This is primarily due to the computational overhead of transformer-based architectures, which grow quadratically with the context length.

2. Output Speed Trends

Output speed measures the rate at which a model generates tokens after producing the first token. Key observations include:

- **Linear Decline in Smaller Models:**

- GPT-4o Mini maintains a relatively high output speed (140 tokens/sec at 1k tokens) and shows a linear decline as input token count increases.

- **Exponential Decline in Larger Models:**

- Larger models experience a sharper drop in output speed as context length increases. For instance, Llama 3.1 Instruct 4B slows to 40 tokens/sec for inputs of 4k tokens, indicating its computational complexity.

- **Importance of High Output Speed:**

- High output speed is crucial for applications that require generating large amounts of text quickly, such as summarization tools, content generation platforms, and real-time transcription systems.

3. Total Response Time and Stability

Total response time combines latency and output speed to provide a holistic view of a model's performance in completing tasks.

- **GPT-4o Mini Dominates in Total Response Time:**

- With a total response time of 0.8 seconds for 1k tokens and 2 seconds for 4k tokens, GPT-4o Mini outperforms all other models in terms of speed and consistency.

- **Larger Models Have Prolonged Response Times:**

- GPT-4 and Llama 3.1 Instruct 4B exhibit response times exceeding 6 seconds for long inputs (>4k tokens), which limits their use in time-critical applications.

- **Stability Over Time:**

- GPT-4o Mini demonstrates minimal variability in response time across repeated trials, ensuring a predictable user experience.
- Larger models show significant fluctuations due to their reliance on complex infrastructure and higher sensitivity to input diversity.

4. Cost Efficiency and Input Token Price

Input token price is a crucial factor for organizations deploying AI models at scale. The findings reveal:

- **Smaller Models Are Cost-Efficient:**

- GPT-4o Mini has the lowest input token price (USD 0.0002), making it highly economical for routine tasks like FAQ bots, basic summarization, and customer support automation.

- **Trade-Offs in Larger Models:**

- GPT-4 and Llama 3.2 90B (Vision) incur higher costs (USD 0.001 per token) due to their superior quality and extended context windows. These models are better suited for specialized tasks requiring high accuracy and reasoning depth.

- **Balancing Cost and Utility:**

- The choice of model should align with the complexity of the task. For example, smaller models can be deployed for high-volume tasks, while larger models should be reserved for high-stakes applications.

5. Context Window and Quality Index

The analysis also highlights the relationship between context window size and quality:

- **Larger Context Windows Enable Complex Tasks:**

- Models like GPT-4 and Llama 3.2 96B (Vision) support extended context windows (up to 8,192 tokens or more), enabling them to handle tasks like analyzing long documents or managing multi-turn conversations effectively.

- **Quality Index Correlates with Context Window:**

- Models with larger context windows and higher quality indices deliver more coherent and nuanced outputs but at the cost of increased latency and token prices.

- **Llama 3.2 96B (Vision) Stands Out:**

- This model balances high quality, multimodal support (text and image processing), and a moderate input token price, making it a strong candidate for advanced tasks requiring both text and image understanding.

6.1.1 Practical Implications of the Study

1. Real-Time Applications:

For applications requiring immediate responses, such as chatbots, live translation, or transcription, smaller models like GPT-4o Mini are ideal due to their low latency, high output speed, and cost efficiency.

2. Complex Reasoning Tasks:

Tasks involving detailed analysis, long-form reasoning, or high-context processing (e.g., legal document review, medical diagnostics) are better suited for larger models like GPT-4, despite their higher latency and costs.

3. Cost-Sensitive Use Cases:

Organizations operating on tight budgets should deploy smaller models for high-volume tasks to minimize costs while maintaining acceptable performance. Larger models should only be used when the added quality justifies the expense.

4. Multimodal Applications:

Llama 3.2 96B (Vision) demonstrates the potential of multimodal models to transform industries like education, content creation, and medical analysis. Its ability to process both text and images expands the scope of AI applications significantly.

6.1.2 Main Conclusions

1. Performance vs. Cost Trade-Offs:

- Smaller models like GPT-40 Mini excel in speed, affordability, and stability, making them suitable for general-purpose and real-time applications.
- Larger models prioritize quality and reasoning depth but incur higher costs and latency, restricting their use to specialized tasks.

2. Latency and Output Speed as Key Metrics:

- Low latency and high output speed are critical for user-facing applications, ensuring a seamless and responsive experience.

3. The Growing Importance of Multimodal Capabilities:

- Llama 3.2 96B highlights the future of AI models that integrate multiple modalities (e.g., text and images) to deliver richer and more versatile outputs.

4. Stability as a Differentiator:

- Models with consistent response times, like GPT-40 Mini, offer a significant advantage in production environments where predictability is crucial.

5. Cost Efficiency Drives Adoption:

- Token prices remain a decisive factor for large-scale deployments, necessitating careful model selection based on the task and budget constraints.

6.1.3 Limitations of the Study

1. Data Source Variability:

- The analysis is based on median values across multiple providers. Differences in hardware, infrastructure, and deployment environments may influence these results.

2. Exclusion of Specialized Metrics:

- Metrics like sentiment analysis performance, domain-specific accuracy, or multilingual support were not included in the scope of this study.

3. Evolving Models:

- As newer models and architectures emerge, the insights presented here may need to be revisited to reflect advancements in AI.

6.1.4 Implications for Future Research

1. Advanced Attention Mechanisms:

- Developing sparse or memory-efficient attention mechanisms to reduce latency and computational complexity for larger context windows.

2. Dynamic Model Adaptation:

- Researching adaptive models that can switch between lightweight and heavy-weight modes based on task requirements.

3. Energy-Efficient AI:

- Investigating methods to reduce the energy consumption of large models during both training and inference phases.

4. Expanding Multimodal Capabilities:

- Enhancing multimodal models to process more complex input combinations, such as video with embedded text or speech alongside images.

Chapter 7

Scope for Future Work

The rapidly evolving field of AI model development presents numerous opportunities for further research and advancements. This chapter outlines potential directions for improving model performance, efficiency, and applicability, focusing on addressing the limitations highlighted in the report while exploring new frontiers.

7.1 Reducing Latency in Larger Models

1. Sparse Attention Mechanisms:

- Explore transformer architectures that utilize sparse attention mechanisms, reducing the quadratic complexity of processing long inputs.
- Implement models like Longformer and BigBird for extended context handling without significant latency increases.

2. Dynamic Token Pruning:

- Research techniques to dynamically prune redundant tokens in real-time during inference to accelerate processing while maintaining output quality.

3. Memory-Efficient Transformers:

- Investigate memory-efficient architectures like Reformer, which utilize reversible layers to reduce memory requirements and computation time.

7.2 Enhancing Multimodal Capabilities

1. Integrating Additional Modalities:

- Expand multimodal models like Llama 3.2 96B (Vision) to incorporate video, audio, and sensor data. This would enable applications in autonomous vehicles, robotics, and multimedia content creation.

2. Fine-Tuning for Domain-Specific Multimodality:

- Fine-tune models to specialize in industries like healthcare, where analyzing combinations of medical images and patient records could improve diagnostic accuracy.

3. **Cross-Modal Attention:**

- Develop cross-modal attention mechanisms to improve the ability of models to contextualize and correlate information from different data types (e.g., aligning text and image content for better understanding).

7.3 Improving Stability in Large Models

1. **Load-Balancing Algorithms:**

- Implement intelligent load-balancing algorithms to distribute computational workloads across servers efficiently, reducing response time variance.

2. **Real-Time Performance Monitoring:**

- Develop tools to monitor and adapt model performance dynamically, ensuring consistent response times regardless of system load.

3. **Caching Mechanisms:**

- Introduce advanced caching mechanisms to store intermediate computations, speeding up repeated tasks and stabilizing response times.

7.4 Reducing Costs for Large Context Models

1. **Token Compression Techniques:**

- Research methods to compress input token sequences without losing essential information, reducing the computational load and input token price.

2. **Efficient Fine-Tuning:**

- Explore parameter-efficient fine-tuning techniques like LoRA (Low-Rank Adaptation) to adapt large models for specific tasks without requiring full retraining.

3. **Hardware Innovations:**

- Collaborate with hardware manufacturers to design GPUs/TPUs optimized for AI inference, focusing on energy-efficient processing of extended context windows.

7.5 Expanding Context Windows

1. Scaling to 100k+ Tokens:

- Investigate architectures that can handle extremely large context windows, enabling models to process entire books, large datasets, or complex multi-document queries.

2. Chunked Attention Mechanisms:

- Use chunked attention, where inputs are divided into smaller chunks processed sequentially, reducing memory overhead while maintaining long-context comprehension.

7.6 Energy Efficiency and Sustainability

1. Green AI:

- Develop algorithms and training methods that reduce the energy consumption of large-scale models, contributing to more sustainable AI systems.

2. Carbon-Neutral Deployments:

- Promote the adoption of renewable energy sources for data centers hosting AI models, minimizing their environmental impact.

7.7 Democratization of AI

1. Accessible Multimodal Models:

- Work on reducing the deployment cost of advanced models like Llama 3.2 96B (Vision), making them accessible to smaller enterprises and researchers.

2. Open-Source Contributions:

- Contribute to the development of open-source AI tools that allow wider adoption and collaboration in academia and industry.

Bibliography

- [1] J. Li, T. Tang, W. X. Zhao, J.-Y. Nie, and J.-R. Wen, “Pretrained language models for text generation: A survey,” 2022, under review. Available at: <https://doi.org/10.48550/arXiv.2201.05273>.
- [2] I. Mackie, I. Sekulic, S. Chatterjee, J. Dalton, and F. Crestani, “Grm: Generative relevance modeling using relevance-aware sample estimation for document retrieval,” 2023, available at: <https://doi.org/10.48550/arXiv.2306.09938>.
- [3] W. Zhou, Y. E. Jiang, E. Wilcox, R. Cotterell, and M. Sachan, “Controlled text generation with natural language instructions,” 2023, accepted to ICML 2023. Available at: <https://doi.org/10.48550/arXiv.2304.14293>.
- [4] A. of the Paper, “Title of the paper,” 2022, available at: <https://doi.org/10.48550/arXiv.2212.11685>.
- [5] Y. Yu, Y. Zhuang, J. Zhang, Y. Meng, A. Ratner, R. Krishna, J. Shen, and C. Zhang, “Large language model as attributed training data generator: A tale of diversity and bias,” 2023, accepted to NeurIPS 2023 (Datasets and Benchmarks Track). Available at: <https://arxiv.org/abs/2306.15895>.
- [6] L. Wang, N. Yang, X. Huang, L. Yang, R. Majumder, and F. Wei, “Improving text embeddings with large language models,” 2024, accepted by ACL 2024. Available at: <https://doi.org/10.48550/arXiv.2401.00368>.
- [7] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 5998–6008.
- [8] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020.
- [9] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” *arXiv preprint arXiv:2005.14165*, 2020.

- [10] OpenAI, “Gpt-4 technical report,” 2023, available at: <https://openai.com/research/gpt-4>.
- [11] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, and A. M. Rush, “Transformers: State-of-the-art natural language processing,” *arXiv preprint arXiv:1910.03771*, 2020.
- [12] G. AI, “Efficient transformer variants: Bigbird and longformer,” 2021, available at: <https://ai.googleblog.com>.
- [13] H. Face, “Transformers documentation,” 2022, available at: <https://huggingface.co>.