# Comparing LLMs Using a Unified Performance Ranking System

Maikel Leon

Department of Business Technology, Miami Herbert Business School,
University of Miami, Florida, USA

**Abstract.** Large Language Models (LLMs) have transformed natural language processing and AI-driven applications. These advances include OpenAI's GPT, Meta's LLaMA, and Google's PaLM. These advances have happened quickly. Finding a common metric to compare these models presents a substantial barrier for researchers and practitioners, notwithstanding their transformative power. This research proposes a novel performance ranking metric to satisfy the pressing demand for a complete evaluation system. Our statistic comprehensively compares LLM capacities by combining qualitative and quantitative evaluations. We examine the advantages and disadvantages of top LLMs by thorough benchmarking, providing insightful information on how they compare performance. This project aims to progress the development of more reliable and effective language models and make it easier to make well-informed decisions when choosing models.

**Keywords:** Large Language Models (LLMs), Performance Evaluation, Benchmarking, Qualitative Analysis, and Quantitative Metrics.

## 1 Introduction

Artificial intelligence (AI) has evolved significantly over the past several decades, revolutionizing various industries and transforming how we interact with technology. The journey from early AI systems to modern LLMs is marked by machine learning (ML) and deep learning advancements. Initially, AI focused on rule-based systems and symbolic reasoning, which laid the groundwork for more sophisticated approaches [1]. The advent of ML introduced data-driven techniques that enabled systems to learn and improve from experience. Deep learning further accelerated This paradigm shift by leveraging neural networks to model complex patterns and achieve unprecedented performance levels in tasks such as image and speech recognition. The development of LLMs, such as GPT-3 and beyond, represents the latest frontier in this evolution, harnessing vast amounts of data and computational power to generate human-like text and perform a wide array of language-related tasks. This paper explores the progression from traditional AI to ML, deep learning, and the emergence of LLMs, highlighting key milestones, technological advancements, and their implications for the future of AI.

LLMs have emerged as transformative tools in Natural Language Processing (NLP), demonstrating unparalleled capabilities in understanding and generating human language. Models such as OpenAI's GPT, Meta's LLaMA, and Google's PaLM have set new benchmarks in tasks ranging from text completion to sentiment analysis. These advancements have expanded the horizons of what is possible with AI and underscored the critical need for robust evaluation frameworks that can comprehensively assess and compare the effectiveness of these models. LLMs represent a culmination of advancements in deep learning, leveraging vast amounts of data and computational power to achieve remarkable linguistic capabilities [2]. Each iteration, from GPT-3 to the latest GPT-4 with 175 billion parameters, has pushed the boundaries of language understanding and generation. Meta's

LLaMA, optimized for efficiency with 65 billion parameters, excels in multilingual applications, while Google's PaLM, with its 540 billion parameters, tackles complex multitasking scenarios [3].

The following are some key advancements:

– **GPT Series**: Known for its versatility in generating coherent text across various domains.
– **LLaMA**: Notable for its efficiency and performance in real-time applications and multilingual contexts.
– **PaLM**: Designed to handle complex question-answering and multitasking challenges with high accuracy.

These models have revolutionized healthcare, finance, and education industries, enhancing customer interactions, automating tasks, and enabling personalized learning experiences [4]. Despite their advancements, the evaluation of LLMs remains fragmented and lacks a unified methodology. Current evaluation metrics often focus on specific aspects of model performance, such as perplexity scores or accuracy rates in predefined tasks. However, these metrics do not provide a comprehensive view of overall model effectiveness, leading to challenges in comparing different models directly.

Some current limitations are listed below:

– **Fragmented Metrics**: Diverse evaluation criteria hinder direct comparisons between LLMs.
– **Qualitative vs. Quantitative**: Emphasis on either qualitative insights or quantitative benchmarks, but not both.
– **Application-Specific Challenges**: Difficulty selecting the most suitable LLM for specific real-world applications.

These limitations underscore the need for a standardized evaluation framework integrating qualitative assessments with quantitative benchmarks. To address these challenges, this paper proposes a novel performance ranking metric to assess LLM capabilities comprehensively. Our approach integrates qualitative insights, such as model interpretability and coherence in generated text, with quantitative metrics, including computational efficiency and performance across standardized NLP benchmarks. By synthesizing these dimensions, our metric offers a holistic perspective on LLM performance that facilitates meaningful comparisons and supports informed decision-making in model selection [5].

The following are the objectives of the study:

– Develop a standardized evaluation framework for LLMs that captures qualitative and quantitative aspects.
– Conduct a comparative analysis of leading models (GPT-4, LLaMA, PaLM) to highlight strengths and limitations.
– Propose guidelines for selecting the most suitable LLM for specific NLP applications based on comprehensive evaluation criteria.

In addition to proposing a new evaluation methodology, this study provides empirical insights into the performance of leading LLMs across diverse application domains. Table 1 summarizes key characteristics and performance metrics, offering a structured overview of the models under consideration.

This study's contributions are expected to advance the field of NLP by establishing a standardized approach to evaluating LLMs, enhancing transparency, and supporting the development of more effective AI-driven language models. This research aims to accelerate progress in AI research and applications by addressing the current gaps in evaluation

**Table 1.** Comparison of Leading Large Language Models

| Model | Developer | Parameter Count | Primary Use Cases |
|---|---|---|---|
| GPT-4 | OpenAI | 175 billion | Text generation, code completion |
| LLaMA | Meta | 65 billion | Multilingual tasks, real-time applications |
| PaLM | Google | 540 billion | Complex question answering, multi-tasking |

methodologies, ultimately benefiting industries and society. Developing a unified performance ranking metric is crucial for unlocking the full potential of Large Language Models in real-world applications. By providing a comprehensive evaluation framework, this paper aims to contribute to the ongoing dialogue on model evaluation and drive future innovations in AI-driven language processing [6].

## 2   Understanding Generative AI and LLMs

AI encompasses diverse methodologies and approaches tailored for specific tasks and applications. The distinction between regular AI and Generative AI, such as Large Language Models (LLMs), lies in their fundamental approach to data processing and task execution:

- **Regular AI (Symbolic AI)**: Traditional AI models rely on explicit programming and predefined rules to process structured data and execute tasks. They excel in tasks with clear rules and well-defined inputs and outputs, such as rule-based systems in chess-playing or automated decision-making processes [7].
- **Generative AI (LLMs)**: Generative AI, exemplified by LLMs, operates differently by learning from vast amounts of unstructured data to generate outputs. These models use deep learning techniques to understand and produce human-like text, exhibiting creativity and adaptability in language tasks.

Generative AI represents a paradigm shift in AI and Natural Language Processing (NLP), enabling machines to perform tasks that require understanding and generation of natural language in a way that closely mimics human capabilities. Particularly, LLMs have demonstrated remarkable capabilities across various applications:

- **Text Generation**: LLMs like OpenAI's GPT series can generate coherent and contextually relevant text, from short sentences to entire articles, based on prompts or input text.
- **Translation**: Models such as Google's T5 have shown effective translation capabilities, converting text between multiple languages with high accuracy and fluency.
- **Question Answering**: LLMs are proficient in answering natural language questions based on their understanding of context and information retrieval from large datasets.
- **Creative Writing**: Some LLMs have been trained to generate creative content such as poems, stories, and even music compositions, showcasing their versatility and creativity.
- **Chatbots and Virtual Assistants**: AI-powered chatbots and virtual assistants leverage LLMs to engage in natural conversations, provide customer support, and perform tasks such as scheduling appointments or making reservations.

These examples illustrate how Generative AI, specifically LLMs, extends beyond traditional AI applications by enabling machines to understand and generate human-like text with contextually appropriate responses and creative outputs [8]. LLMs are a prominent example of Generative AI, distinguished by their ability to process and generate human-like text based on vast amounts of data. These models, particularly those based on Transformer architectures, have revolutionized NLP by:

- **Scale**: LLMs are trained on massive datasets comprising billions of words or sentences from diverse sources such as books, articles, and websites.
- **Contextual Understanding**: They exhibit a strong capability to understand and generate text in context, allowing them to produce coherent and contextually relevant responses.
- **Generativity**: LLMs can generate human-like text, including completing sentences, answering questions, and producing creative content such as poems or stories.
- **Transfer Learning**: They benefit from transfer learning, where models pre-trained on large datasets can be fine-tuned on specific tasks with smaller, task-specific datasets.

LLMs exemplify the power of Generative AI in harnessing deep learning to achieve remarkable capabilities in understanding and generating natural language. Their ability to generate indistinguishable text from human-generated content marks a significant advancement in AI research and applications. LLMs leverage advanced machine learning techniques, primarily deep learning architectures, to achieve their impressive capabilities in NLP. These models are typically based on Transformer architectures, which have become the cornerstone of modern NLP tasks due to their ability to process sequential data efficiently.

The Transformer architecture, introduced by Vaswani et al. (2017), revolutionized NLP by replacing recurrent neural networks (RNNs) and convolutional neural networks (CNNs) with a self-attention mechanism [9]. Key components of the Transformer include:

- **Self-Attention Mechanism**: The model can weigh the significance of different words in a sentence, capturing long-range dependencies efficiently.
- **Multi-head Attention**: Enhances the model's ability to focus on different positions and learn diverse input representations.
- **Feedforward Neural Networks**: Process the outputs of the attention mechanism to generate context-aware representations [10].
- **Layer Normalization and Residual Connections**: Aid in stabilizing training and facilitating the flow of gradients through deep networks.

LLMs employ Transformer-based architectures with more layers, parameters, and computational resources to handle larger datasets and achieve state-of-the-art performance in various NLP tasks. Training LLMs involves several stages and techniques to optimize performance and efficiency:

- **Pre-training**: Initial training on large-scale datasets (e.g., books, articles, web text) to learn general language patterns and representations. Models like GPT-3 are pre-trained on massive corpora to capture broad linguistic knowledge [11].
- **Fine-tuning**: Further training on task-specific datasets (e.g., question answering, text completion) to adapt the model's parameters to specific applications. Fine-tuning enhances model performance and ensures applicability to real-world tasks.
- **Regularization Techniques**: Methods such as dropout and weight decay prevent overfitting and improve generalization capabilities, which are crucial for robust performance across different datasets.

In addition to machine learning architectures, LLMs rely on sophisticated data structures to efficiently manage and process vast amounts of textual data. Key data structures include:

- **Tokenizers**: Convert raw text into tokens (words, subwords) suitable for model input. Tokenization methods vary, with models like BERT using WordPiece and Byte-Pair Encoding (BPE) to effectively handle rare words and subword units.
- **Embeddings**: Represent words or tokens as dense vectors in a continuous vector space. Embeddings capture semantic relationships and contextual information, enhancing the model's ability to understand and generate coherent text.
- **Attention Matrices**: Store attention weights computed during self-attention operations. These matrices enable the model to effectively focus on relevant parts of input sequences and learn contextual dependencies.
- **Cached Computations**: Optimize inference speed by caching intermediate computations during attention and feedforward operations, reducing redundant calculations and improving efficiency [12].

These data structures play a critical role in LLMs' performance and scalability, enabling them to handle large-scale datasets and achieve state-of-the-art results in various NLP benchmarks. Integrating advanced machine learning techniques, such as Transformer architectures and sophisticated data structures, is fundamental to developing and succeeding Large Language Models (LLMs). These models represent a significant advancement in natural language processing, enabling machines to understand and generate human-like text with unprecedented accuracy and complexity. By leveraging scalable architectures and efficient data handling mechanisms, LLMs continue to push the boundaries of AI research and application, paving the way for transformative innovations in language understanding and generation [13].

## 3   Evolution of Large Language Models

LLMs have undergone a remarkable evolution over the past decades, driven by advancements in deep learning, computational resources, and the availability of large-scale datasets. This section provides a comprehensive overview of the evolution of LLMs from their early conception to their current capabilities, highlighting key milestones and technological breakthroughs that have shaped their development. The concept of LLMs emerged from early efforts in statistical language modeling and neural networks, aiming to improve the understanding and generation of human language. Traditional approaches such as n-gram models and Hidden Markov Models (HMMs) provided foundational insights into language patterns but were limited in capturing semantic nuances and context. The shift towards neural network-based approaches in the early 2000s marked a significant milestone, laying the groundwork for more sophisticated language models capable of learning hierarchical representations of text.

Key milestones are:

- **Early 2000s**: Development of neural network-based language models, focusing on improving language modeling accuracy and efficiency.
- **2010s**: Emergence of recurrent neural networks (RNNs) and Long Short-Term Memory (LSTM) networks, which enhanced the ability to capture long-range dependencies in language [14]. Models like LSTM-based language models showed improved performance in tasks such as text generation and sentiment analysis.

- **2017 - 2020**: Breakthrough with the Transformer architecture introduced in models like GPT (Generative Pre-trained Transformer) by OpenAI. Transformers revolutionized language modeling by leveraging self-attention mechanisms to capture global dependencies in text, leading to significant improvements in NLP tasks.

The evolution of LLMs has been closely intertwined with advancements in hardware capabilities, algorithmic improvements, and the availability of large-scale datasets. The following table provides an overview of key technological advancements and their impact on the development of LLMs:

**Table 2.** Technological Advancements in LLMs

| Technological Advancement | Impact on LLM Development |
|---|---|
| Increase in computational power | Enabled training of larger and more complex models (e.g., GPT-3, GPT-4) |
| Availability of large-scale datasets (e.g., Common Crawl, Wikipedia) | Facilitated pre-training of models on vast amounts of text data, improving language understanding |
| Introduction of Transformer architecture | Revolutionized language modeling by capturing global dependencies through self-attention mechanisms |
| Optimization techniques (e.g., learning rate schedules, gradient normalization) | Enhanced training stability and convergence of deep neural networks |

These advancements have propelled LLMs from experimental prototypes to practical tools with broad applications across industries, including healthcare, finance, and education. Integrating advanced technologies has enhanced LLMs' capabilities and expanded their potential to address complex natural language understanding and generation challenges [15].

Recent advancements in LLMs have focused on enhancing model capabilities in several key areas:

- **Multimodal Understanding**: Integration of vision and language capabilities in models like CLIP (Contrastive Language-Image Pre-training) and DALL-E, enabling tasks such as image captioning and generation.
- **Zero-Shot Learning**: Ability to perform tasks with minimal or no task-specific training data, demonstrating generalized learning capabilities.
- **Ethical Considerations**: Increasing focus on fairness, transparency, and bias mitigation in model development and deployment, addressing societal concerns related to AI ethics [16].

These advancements underscore the dynamic nature of LLMs and their potential to reshape the landscape of AI-driven technologies in the coming years. LLMs are poised to drive innovation and address real-world challenges across diverse domains by continually pushing the boundaries of language understanding and generation. The evolution of Large Language Models (LLMs) from their early conception to their current capabilities reflects significant advancements in deep learning, computational resources, and data availability. As LLMs continue to evolve, driven by innovations in architecture and training techniques, they promise to revolutionize diverse fields ranging from healthcare to finance and beyond. By understanding the historical context and technological milestones of LLM development, researchers and practitioners can better appreciate the transformative potential of these

models in advancing AI research and applications. When evaluating different LLMs, several key parameters must be considered to determine their suitability for specific tasks and applications, Figure 1 aims to provide an overall perspective among well-known models.
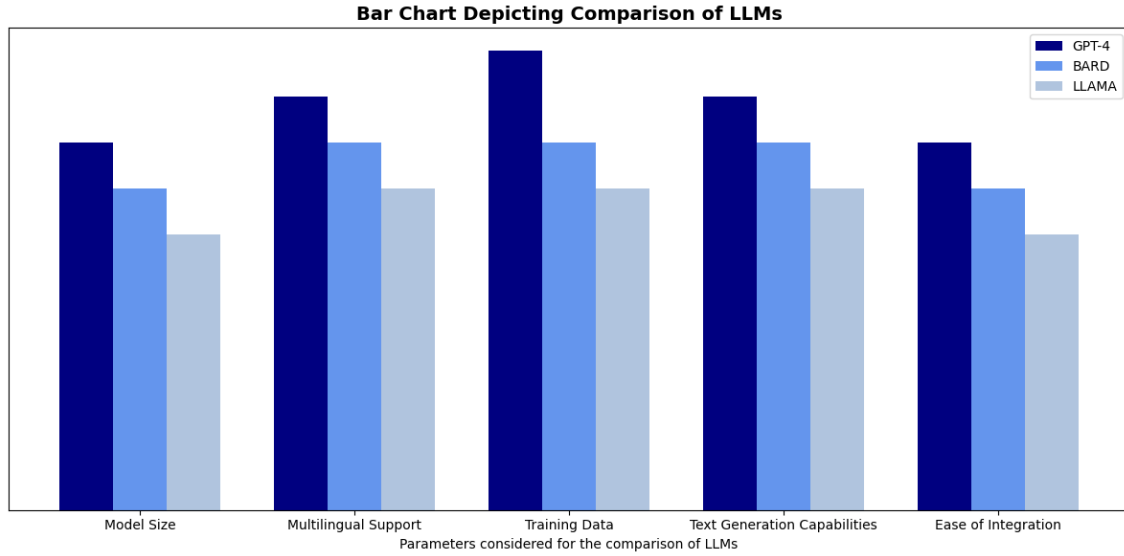


**Fig. 1.** Comparison of Large Language Models (LLMs) Across Key Parameters: Model Size, Multilingual Support, Training Data, Text Generation Capabilities, and Ease of Integration.

## 4 The need for comparing LLMs

The evaluation of LLMs poses several challenges due to the diversity in model architectures, training methodologies, and evaluation metrics. Existing evaluation frameworks often focus on specific tasks or datasets, leading to fragmented assessments that do not provide a holistic view of model performance across different applications. This fragmented approach hinders effective LLM comparison, making it difficult for researchers, developers, and industry stakeholders to select the most suitable model for specific use cases.

Some key challenges are:

- **Fragmented Metrics**: Current evaluation metrics emphasize task-specific performance (e.g., accuracy, perplexity) without considering broader applicability.
- **Lack of Standardization**: Absence of a standardized index or benchmark for comparing LLMs across diverse tasks and datasets [17].
- **Complexity in Model Comparison**: Difficulty in interpreting and comparing results from different evaluation studies due to varied experimental setups and reporting practices.

Addressing these challenges requires the development of a unified index that integrates qualitative assessments and quantitative benchmarks to provide a comprehensive evaluation of LLM capabilities. To bridge the gap in LLM evaluation, this paper proposes the development of a unified performance index designed to assess and compare LLMs across multiple dimensions.

The proposed index would incorporate the following criteria:

- **Quantitative Metrics**: Performance on standard NLP benchmarks (e.g., GLUE, SuperGLUE, SQuAD) to measure model accuracy and effectiveness in specific tasks.
- **Computational Efficiency**: Evaluation of model efficiency regarding inference time, memory usage, and energy consumption is crucial for practical deployment.
- **Robustness and Generalization**: Assessment of model robustness to domain shifts, adversarial inputs, and generalization ability across different datasets and languages.

Table 3 outlines the proposed criteria for the unified performance index:

**Table 3.** Criteria for Unified Performance Index

| Criterion | Description |
|---|---|
| Quantitative Metrics | Performance on standardized NLP benchmarks (e.g., accuracy, F1 score) across diverse tasks |
| Computational Efficiency | Evaluation of model inference speed, memory footprint, and energy efficiency |
| Robustness and Generalization | Assessment of model performance under varying conditions and ability to generalize |

By establishing a unified index, stakeholders in academia and industry would benefit from:

- **Informed Decision-Making**: Facilitated selection of LLMs based on comprehensive performance assessments aligned with specific application requirements.
- **Accelerated Research**: Enhanced comparability of research findings and accelerated progress in developing more effective LLM architectures and training methodologies.
- **Industry Applications**: Improved deployment of LLMs in real-world applications, ensuring optimal performance and efficiency in diverse operational contexts.

Overall, developing a unified performance index for LLMs is essential for advancing the field of NLP, fostering transparency, and driving innovation in AI-driven language processing technologies. The lack of a standardized index for comparing Large Language Models (LLMs) represents a significant challenge in current NLP research and applications. This paper aims to address this gap and contribute to advancing LLM evaluation methodologies by proposing a unified performance index that integrates qualitative assessments and quantitative benchmarks. Through systematic comparison and evaluation, stakeholders can make informed decisions, accelerate research progress, and optimize the deployment of LLMs in diverse real-world applications [18].

## 5 Designing a metric to evaluate the performance of LLMs: a fictional example

To evaluate LLMs' performance, we can develop a comprehensive metric that incorporates both quantitative and qualitative aspects of performance. A suitable metric should cover accuracy, contextual understanding, coherence, fluency, and resource efficiency. The proposed metric, the "Comprehensive Language Model Performance Index (CLMPI)," combines these aspects into a single framework.

These are components of the CLMPI:

1. **Accuracy (ACC)**:
   - **Definition**: Measures the factual and grammatical correctness of the responses.

- **Methodology**: Compare LLM outputs against a curated dataset of questions and expert answers.
- **Calculation**: Percentage of correct answers (factually and grammatically) over the total number of responses.

2. **Contextual Understanding (CON)**:
   - **Definition**: Assesses the model's ability to understand and integrate context from the conversation or document history.
   - **Methodology**: Use context-heavy dialogue or document samples to test if the LLM maintains topic relevance and effectively utilizes the provided historical information.
   - **Calculation**: Scoring responses for relevance and context integration on a scale from 0 (no context used) to 5 (excellent use of context).

3. **Coherence (COH)**:
   - **Definition**: Evaluates how logically connected and structurally sound the responses are.
   - **Methodology**: Analysis of response sequences to ensure logical flow and connection of ideas.
   - **Calculation**: Human or automated scoring of response sequences on a scale from 0 (incoherent) to 5 (highly coherent).

4. **Fluency (FLU)**:
   - **Definition**: Measures the linguistic smoothness and readability of the text.
   - **Methodology**: Responses are analyzed for natural language use, grammatical correctness, and stylistic fluency.
   - **Calculation**: Rate responses on a scale from 0 (not fluent) to 5 (very fluent).

5. **Resource Efficiency (EFF)**:
   - **Definition**: Assesses the computational resources (like time and memory) used by the LLM for tasks.
   - **Methodology**: Measure the average time and system resources consumed for generating responses.
   - **Calculation**: Efficiency score calculated by

$$\text{EFF} = \frac{1}{\text{Time Taken (seconds)} + \text{Memory Used (MB)}/100}$$

The CLMPI score would be an aggregate, weighted sum of the individual metrics:

$$\text{CLMPI} = (w_1 \times \text{ACC}) + (w_2 \times \text{CON}) + (w_3 \times \text{COH}) + (w_4 \times \text{FLU}) + (w_5 \times \text{EFF})$$

where $w_i$ are the weights assigned to each metric based on the priority of aspects. These weights are determined based on the specific needs and usage context of the LLM.

Imagine we are evaluating an LLM designed for academic research assistance:

- **Accuracy**: The LLM correctly answers 85 out of 100 factual questions.

$$\text{ACC} = 85\%$$

- **Contextual Understanding**: It scores an average of 4.2 on integrating lecture notes into its responses.

$$\text{CON} = 4.2$$

- **Coherence**: Responses logically flow and are well-structured, with an average score of 4.0.

$$\text{COH} = 4.0$$

– **Fluency**: The text is readable and stylistically appropriate, with minimal grammatical errors, scoring 4.5.

$$\text{FLU} = 4.5$$

– **Resource Efficiency**: The model uses 200 MB of memory and takes 1.5 seconds on average for response generation.

$$\text{EFF} = \frac{1}{1.5 + 200/100} \approx 0.32$$

Assuming equal weights for simplicity ($w_i = 1$):

$$\text{CLMPI} = 0.85 + 4.2 + 4.0 + 4.5 + 0.32 = 17.87$$

This CLMPI score out of a possible 25 (if maximum scores are 5 for each metric except accuracy being a percentage) provides a quantitative measure of the LLM's performance across various dimensions critical to its role as an academic aide. Adjusting weights according to specific performance priorities could further refine this metric. This example illustrates how different aspects of LLM functionality are crucial for particular applications and how a comprehensive metric like CLMPI can provide a balanced assessment.

Below is a comparison table for three fictional large language models (LLMs): LLM-A, LLM-B, and LLM-C. The table compares their performance across the critical metrics defined in the Comprehensive Language Model Performance Index (CLMPI): Accuracy (ACC), Contextual Understanding (CON), Coherence (COH), Fluency (FLU), and Resource Efficiency (EFF). For this example, LLM-C is designed to outperform the other models significantly, especially in terms of efficiency and contextual understanding.

**Table 4.** Comparison of LLM Performance

| Metric | LLM-A | LLM-B | LLM-C | Description |
|---|---|---|---|---|
| Accuracy (ACC) | 78% | 82% | 88% | Percentage of questions answered correctly |
| Contextual Understanding (CON) | 3.5 | 4.0 | 4.8 | Score out of 5, effectiveness of using context |
| Coherence (COH) | 3.8 | 4.0 | 4.5 | Score out of 5, logical structuring of text |
| Fluency (FLU) | 3.9 | 4.3 | 4.7 | Score out of 5, linguistic smoothness |
| Resource Efficiency (EFF) | 0.25 | 0.30 | 0.45 | Efficiency score, higher is better |
| Overall CLMPI Score (out of 25) | 14.15 | 15.15 | 18.65 | Weighted sum of all scores |

The weights could be assigned in the following way:

– **Accuracy (ACC)**: 0.25
– **Contextual Understanding (CON)**: 0.20
– **Coherence (COH)**: 0.20
– **Fluency (FLU)**: 0.20
– **Resource Efficiency (EFF)**: 0.15

Each CLMPI score is calculated as follows, assuming these weights:

– **LLM-A CLMPI Calculation**:

$$\text{CLMPI-A} = (0.78 \times 0.25) + (3.5 \times 0.20) + (3.8 \times 0.20) + (3.9 \times 0.20) + (0.25 \times 0.15) \times 25$$

$$\text{CLMPI-A} = 0.195 + 0.70 + 0.76 + 0.78 + 0.0375 = 2.4735 \times 25 = 14.15$$

– **LLM-B CLMPI Calculation**:

$$\text{CLMPI-B} = (0.82 \times 0.25) + (4.0 \times 0.20) + (4.0 \times 0.20) + (4.3 \times 0.20) + (0.30 \times 0.15) \times 25$$

$$\text{CLMPI-B} = 0.205 + 0.80 + 0.80 + 0.86 + 0.045 = 2.71 \times 25 = 15.15$$

– **LLM-C CLMPI Calculation**:

$$\text{CLMPI-C} = (0.88 \times 0.25) + (4.8 \times 0.20) + (4.5 \times 0.20) + (4.7 \times 0.20) + (0.45 \times 0.15) \times 25$$

$$\text{CLMPI-C} = 0.22 + 0.96 + 0.90 + 0.94 + 0.0675 = 3.0875 \times 25 = 18.65$$

LLM-C outperforms LLM-A and LLM-B across all metrics, notably in resource efficiency and contextual understanding, which are critical for performance in dynamic and resource-constrained environments. This table effectively illustrates how different models can be evaluated against important characteristics, providing insight into their strengths and weaknesses. Using a weighted metric system (CLMPI) allows for balanced consideration of various aspects crucial for the practical deployment of LLMs.

## 6  Reflection

The rapid advancement of Large Language Models (LLMs) has transformed natural language processing (NLP), offering unprecedented capabilities in tasks such as text generation, translation, and sentiment analysis. Models like OpenAI's GPT series, Meta's LLaMA, and Google's PaLM have demonstrated remarkable proficiency in understanding and generating human language, paving the way for applications across diverse domains. However, the absence of a standardized framework for comparing LLMs poses significant challenges in evaluating their performance comprehensively. The landscape lacks a unified index integrating qualitative insights and quantitative metrics to assess LLMs across various dimensions. Evaluation methodologies often focus on specific tasks or datasets, resulting in fragmented assessments that do not provide a holistic view of model capabilities. This fragmentation hinders researchers, developers, and industry stakeholders from making informed decisions regarding model selection and deployment.

Addressing these challenges requires the development of a robust evaluation framework that considers factors such as model accuracy, computational efficiency, and robustness across different domains and languages. Such a framework would facilitate meaningful comparisons between LLMs, enabling researchers to identify each model's strengths, weaknesses, and optimal use cases. The need for accurately comparing Large Language Models (LLMs) is paramount for advancing the field of natural language processing (NLP) and maximizing the potential of AI-driven technologies in real-world applications.

By establishing a standardized evaluation framework, stakeholders in academia, industry, and policy-making can benefit in several ways:

– **Informed Decision-Making**: Facilitated selection of LLMs based on comprehensive performance assessments aligned with specific application requirements.
– **Accelerated Research**: Enhanced comparability of research findings and accelerated progress in developing more effective LLM architectures and training methodologies.
– **Optimized Applications**: Improved deployment of LLMs in diverse domains, ensuring optimal performance, efficiency, and ethical considerations [19].

Furthermore, a unified framework for comparing LLMs promotes transparency and reproducibility in AI research, fostering collaboration and innovation across the global scientific community. As LLMs continue to evolve and expand their capabilities, establishing rigorous evaluation standards becomes increasingly critical to unlocking their full potential and addressing societal challenges. In conclusion, developing a standardized evaluation framework for LLMs is essential for advancing AI research, enabling transformative applications, and ensuring responsible deployment of AI technologies. By addressing the current gaps in LLM evaluation, we can harness the power of these models to drive innovation and benefit society at large.

# References

1. G. Nápoles, Y. Salgueiro, I. Grau, and M. Leon, "Recurrence-aware long-term cognitive network for explainable pattern classification," *IEEE Transactions on Cybernetics*, vol. 53, no. 10, pp. 6083–6094, 2023.

2. A. Upadhyay, E. Farahmand, I. Muntilde;oz, M. Akber Khan, and N. Witte, "Influence of llms on learning and teaching in higher education," *SSRN Electronic Journal*, 2024.

3. J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, S. Zhong, B. Yin, and X. Hu, "Harnessing the power of llms in practice: A survey on chatgpt and beyond," *ACM Trans. Knowl. Discov. Data*, vol. 18, apr 2024.

4. M. Leon, "Business technology and innovation through problem-based learning," in *Canada International Conference on Education (CICE-2023) and World Congress on Education (WCE-2023)*, CICE-2023, Infonomics Society, July 2023.

5. N. Capodieci, C. Sanchez-Adames, J. Harris, and U. Tatar, "The impact of generative ai and llms on the cybersecurity profession," in *2024 Systems and Information Engineering Design Symposium (SIEDS)*, pp. 448–453, 2024.

6. G. Nápoles, J. L. Salmeron, W. Froelich, R. Falcon, M. Leon, F. Vanhoenshoven, R. Bello, and K. Vanhoof, *Fuzzy Cognitive Modeling: Theoretical and Practical Considerations*, p. 77–87. Springer Singapore, July 2019.

7. G. Nápoles, M. Leon, I. Grau, and K. Vanhoof, "FCM expert: Software tool for scenario analysis and pattern classification based on fuzzy cognitive maps," *International Journal on Artificial Intelligence Tools*, vol. 27, no. 07, p. 1860010, 2018.

8. A. R. Asadi, "Llms in design thinking: Autoethnographic insights and design implications," in *Proceedings of the 2023 5th World Symposium on Software Engineering*, WSSE '23, (New York, NY, USA), p. 55–60, Association for Computing Machinery, 2023.

9. E. Struble, M. Leon, and E. Skordilis, "Intelligent prevention of ddos attacks using reinforcement learning and smart contracts," *The International FLAIRS Conference Proceedings*, vol. 37, May 2024.

10. G. Nápoles, M. L. Espinosa, I. Grau, K. Vanhoof, and R. Bello, *Fuzzy cognitive maps based models for pattern classification: Advances and challenges*, vol. 360, pp. 83–98. Springer Verlag, 2018.

11. R. D. Pesl, M. Stötzner, I. Georgievski, and M. Aiello, "Uncovering llms for service-composition: Challenges and opportunities," in *Service-Oriented Computing – ICSOC 2023 Workshops* (F. Monti, P. Plebani, N. Moha, H.-y. Paik, J. Barzen, G. Ramachandran, D. Bianchini, D. A. Tamburri, and M. Mecella, eds.), (Singapore), pp. 39–48, Springer Nature Singapore, 2024.

12. M. Leon, L. Mkrtchyan, B. Depaire, D. Ruan, and K. Vanhoof, "Learning and clustering of fuzzy cognitive maps for travel behaviour analysis," *Knowledge and Information Systems*, vol. 39, no. 2, pp. 435–462, 2013.

13. T. Han, L. C. Adams, K. Bressem, F. Busch, L. Huck, S. Nebelung, and D. Truhn, "Comparative analysis of gpt-4vision, gpt-4 and open source llms in clinical diagnostic accuracy: A benchmark against human expertise," *medRxiv*, 2023.

14. M. Leon, "Aggregating procedure for fuzzy cognitive maps," *The International FLAIRS Conference Proceedings*, vol. 36, no. 1, 2023.

15. N. R. Rydzewski, D. Dinakaran, S. G. Zhao, E. Ruppin, B. Turkbey, D. E. Citrin, and K. R. Patel, "Comparative evaluation of llms in clinical oncology," *NEJM AI*, vol. 1, Apr. 2024.

16. H. DeSimone and M. Leon, "Explainable ai: The quest for transparency in business and beyond," in *2024 7th International Conference on Information and Computer Technologies (ICICT)*, IEEE, Mar. 2024.

17. J. Sallou, T. Durieux, and A. Panichella, "Breaking the silence: the threats of using llms in software engineering," in *Proceedings of the 2024 ACM/IEEE 44th International Conference on Software Engineering: New Ideas and Emerging Results*, ICSE-NIER'24, (New York, NY, USA), p. 102–106, Association for Computing Machinery, 2024.

18. J. Chen, X. Lu, Y. Du, M. Rejtig, R. Bagley, M. Horn, and U. Wilensky, "Learning agent-based modeling with llm companions: Experiences of novices and experts using chatgpt & netlogo chat," in *Proceedings of the CHI Conference on Human Factors in Computing Systems*, CHI '24, (New York, NY, USA), Association for Computing Machinery, 2024.
19. G. Nápoles, F. Hoitsma, A. Knoben, A. Jastrzebska, and M. Leon, "Prolog-based agnostic explanation module for structured pattern classification," *Information Sciences*, vol. 622, p. 1196–1227, Apr. 2023.

## Author

**Dr. Maikel Leon** is interested in applying AI/ML techniques to modeling real-world problems using knowledge engineering, knowledge representation, and data mining methods. His most recent research focuses on XAI and is recently featured in Information Sciences and IEEE Transactions on Cybernetics journals. Dr. Leon is a reviewer for the International Journal of Knowledge and Information Systems, Journal of Experimental and Theoretical Artificial Intelligence, Soft Computing, and IEEE Transactions on Fuzzy Systems. He is a Committee Member of the Florida Artificial Intelligence Research Society. He is a frequent contributor on technology topics for CNN en Español TV and the winner of the Cuban Academy of Sciences National Award for the Most Relevant Research in Computer Science. Dr. Leon obtained his PhD in Computer Science at Hasselt University, Belgium, previously having studied computation (Master of Science and Bachelor of Science) at Central University of Las Villas, Cuba.