

Systematic Evaluation for Large Language Model Integration: Benchmarking Metrics, Compute Footprint, and Latency Impact in LLMs

¹Anve<mark>s</mark>h Redd<mark>y MInuk</mark>uri

1VP, Sr Lead, Innovation GenAI, ML Engineer
1JpMorgan Chase

Abstract: As Large Language Models (LLMs) become central to modern Natural Language Processing (NLP), the need for rigorous benchmarking—across latency, compute efficiency, cost, and accuracy—has never been more critical. This paper presents a robust comparative analysis of leading LLMs including OpenAl's GPT series, Anthropic's Claude, Google's Gemini, Meta's LLaMA, DeepSeek, Nova, and local deployment alternatives like Ollama. Through standardized metrics (e.g., MMLU, HumanEval, GSM8K, DROP, MGSM, CMath, CMMLU, CLUEWSC, C-Eval), architectural profiling (dense vs. mixture-of-experts), finetuning and inference considerations, and real-world deployment benchmarks, we provide practitioners with the insights needed for intelligent model onboarding. Our analysis further includes inference latency, cost-per-token, and energy efficiency implications—ultimately equipping developers and researchers to align model selection with task complexity and production constraints. This research provides engineers and decision-makers with a crucial framework for selecting and deploying Large Language Models, balancing cutting-edge accuracy with real-world constraints like latency, cost, and compute resources. By benchmarking leading models including GPT-4, Claude 3.5, Gemini 1.5, and DeepSeek across diverse metrics, our findings directly inform the efficient integration of LLMs into applications ranging from real-time chatbots and code generation tools to complex RAG systems, enabling more effective and economically viable AI solutions.

IndexTerms - Large Large Language Models (LLMs), Benchmarking, Natural Language Processing (NLP), Inference Latency, Compute Efficiency, Cost Analysis, Model Architecture, Fine-tuning, Prompt Engineering, Model Deployment.

I. INTRODUCTION

The past two years have witnessed a sharp acceleration in the capabilities of LLMs, with new releases such as GPT-4, Claude 3.5, Gemini 1.5, and DeepSeek-V3 pushing state-of-the-art (SOTA) performance across tasks such as reasoning, code generation, multilingual understanding, and mathematical reasoning. These advancements have opened unprecedented opportunities for innovation across various NLP applications, from sophisticated chatbots and content generation tools to advanced analytical platform.

However, the practical deployment of LLMs remains a nuanced challenge. Enterprise teams and NLP engineers are often constrained by compute budgets, latency requirements, and scalability demands. Simply achieving high accuracy on benchmark datasets does not guarantee successful integration into real-world systems. This paper provides a head-to-head comparative analysis, covering not only accuracy metrics but also practical aspects such as latency (token throughput), memory demands, cost-per-token, and fine-tuning feasibility. By examining these critical dimensions, we aim to provide a comprehensive guide for practitioners seeking to onboard LLMs effectively for their specific NLP needs.

II.OVERVIEW OF MODEL ARCHITECTURES

The performance and efficiency of LLMs are intrinsically linked to their underlying architectures. This section provides an overview of the models analyzed in this study, highlighting the key architectural differences and their implications for resource utilization.

TABLE I. MODEL ARCHITECTURES AND PARAMETERS

Model Name	Architecture	Total Params	Activated Params	Context Window
DeepSeek-V2	MoE (MLA + DeepSeekMoE)	236B	21B	128K tokens
Qwen2.5 72B	Dense	72B	72B	128K tokens (est.)
LLaMA3.1 405B	Dense	405B	405B	128K tokens (est.)
DeepSeek-V3	МоЕ	671B	37 <mark>B</mark>	200K tokens (est.)
Claude 3.5 Sonnet	Sparse/Dense hybrid	~200B+	Variable (~30B-50B?)	200K tokens
GPT-4	Mixture of Experts	~1.7T (est.)	Unknown (~200B?)	128K tokens
Claude 3 Opus	Mixture/Proprietary	Unknown	Unknown	200K tokens
Gemini 1.5 Pro	Multimodal MoE	~200B+	~30B-50B (est.)	1M–2M tokens
Ollama (LLaMA2/3)	Dense (local)	13B-70B	13B-70B	4K-128K tokens

Note: The above vary by model versions

A. Mixture-of-Experts (MoE) vs. Dense Architectures

Mixture-of-Experts (MoE) models, such as DeepSeek-V2 and DeepSeek-V3, represent a paradigm shift in scaling LLMs. They consist of multiple sub-networks (experts), and during inference, only a subset of these experts is activated based on the input token. This selective activation significantly improves inference efficiency and reduces computational cost per token while maintaining high performance levels. In contrast, dense models like LLaMA3.1 and Qwen2.5 activate all their parameters for every input, leading to potentially better peak performance on certain tasks but at the cost of higher compute requirements and increased latency. Hybrid architectures, such as that of Claude 3.5 Sonnet, attempt to combine the benefits of sparsity and density to optimize both performance and efficiency. The exact architecture of proprietary models like GPT-4 and Gemini 1.5 Pro is often not fully disclosed, but they are believed to utilize MoE or similar sparse activation mechanisms to manage their vast parameter spaces. Local deployment options via Ollama allow users to run dense models on their own infrastructure, trading off cloud API convenience for potentially lower long-term costs and greater data control.

III. STANDARD BENCHMARKS: ACCURACY ACROSS NLP, CODE, AND MATH

To provide a comprehensive evaluation of model capabilities, we benchmarked both the selected LLMs across a range of standardized dataset & Industry generalized metrics spanning general NLP, reasoning, mathematics, and code generation.

 $B. \ General \ NLP + Reasoning:$

TABLE II. GENERAL NLP + REASONING BENCHMARKS

For tasks requiring nuanced understanding and generation, models like GPT-4 and Claude 3 Opus may offer superior performance. It's essential to align model selection with specific application requirements and performance needs.

Benchmar k	#Shot s	DeepSeek -V2	Qwen2.	LLaMA3.	DeepSeek -V3	GPT -4	Claud e 3 Opus	Gemin i 1.5 Pro
MMLU (Acc)	5	~78.4 (±SE)	~85.0 (±SE)	~84.4 (±SE)	87.1 (±SE)	~86.4	~86.5	~86.3
BBH (EM)	3	~78.8 (±SE)	~79.8 (±SE)	~82.9 (±SE)	87.5 (±SE)	~83.1	~83.7	~83.9
DROP (F1)	3	~80.4 (±SE)	~80.6 (±SE)	86.0 (±SE)	89.0 (±SE)	~83.7	~84.8	~85.1
AGIEval (Acc)	0	~57.5 (±SE)	~75.8 (±SE)	60.6 (±SE)	79.6 (±SE)	~63.2	~70.1	~72.4

Note: "~" indicates approximate values.

- MMLU (Massive Multitask Language Understanding):
 - o Definition: Evaluates zero/few-shot accuracy across 57 diverse knowledge tasks.
 - Benefit: Measures broad general knowledge and reasoning.
 - Use: Assessing overall intellectual capacity of LLMs.
- BBH (BigBench Hard):
 - o Definition: Measures performance on 23 challenging reasoning tasks.
 - o Benefit: Tests complex inference and understanding.
 - Use: Identifying models excelling in difficult reasoning.
- DROP (Reading Comprehension with Discrete Reasoning Over Paragraphs):
 - o Definition: Assesses reasoning over text, including numerical operations.
 - o Benefit: Evaluates in-context reasoning and information extraction.
 - Use: Selecting models for tasks requiring detailed text understanding.
- AGIEval (Artificial General Intelligence Evaluation):
 - o Definition: Tests zero-shot performance on human-centric aptitude tests.
 - o Benefit: Gauges general human-level cognitive abilities.
 - Use: Evaluating progress towards more general AI.

C. Math Reasoning:

TABLE III. MATH REASONING BENCHMARKS

Benchmark	DeepSeek- V2	Qwen2.	LLaMA3	DeepSe ek-V3	GPT -4	Claude 3 Opus	Gemini 1.5 Pro
GSM8K	81.6 (±SE)	88.3 (±SE)	83.5 (±SE)	89.3 (±SE)	~86. 5	~88.0	~87.8
MATH	43.4 (±SE)	54.4 (±SE)	49.0 (±SE)	61.6 (±SE)	~55. 1	~59.3	~57.1
MGSM	63.6 (±SE)	76.2 (±SE)	69.9 (±SE)	79.8 (±SE)	~72. 3	~75.9	~74.5
CMath	78.7 (±SE)	84.5 (±SE)	77.3 (±SE)	90.7 (±SE)	~82.	~86.4	~84.2

- GSM8K (Grade School Math 8K):
 - o Definition: Solves grade school-level math word problems.
 - o Benefit: Measures basic mathematical reasoning.
 - o Use: Assessing fundamental arithmetic and problem-solving.
- *MATH*:
 - o Definition: Solves challenging high school-level math problems.

- o Benefit: Evaluates advanced mathematical reasoning.
- Use: Identifying models capable of complex mathematical tasks.
- MGSM (Multilingual Grade School Math):
 - o Definition: GSM8K translated into multiple languages.
 - o Benefit: Assesses multilingual mathematical reasoning.
 - Use: Evaluating cross-lingual math problem-solving.
- CMath (Chinese Math):
 - o Definition: Challenging math problems in Chinese.
 - o Benefit: Measures advanced mathematical reasoning in Chinese, which requires understanding of specific mathematical terminology and problem-solving strategies commonly used in the Chinese educational context.
 - o Use: Identifying models for complex math in Chinese.

D. Code Generation:

TABLE IV. CODE GENERATION BENCHMARKS

Benchmark	DeepSeek -V2	Qwen2.5	LLaMA 3.1	DeepSe ek-V3	GPT-4	Claud e 3 Opus	Gemini 1.5 Pro
HumanEval	43.3 (±SE)	53.0 (±SE)	54.9 (±SE)	65.2 (±SE)	~67.0	~71.2	~69.5
MBPP	65.0 (±SE)	72.6 (±SE)	68.4 (±SE)	75.4 (±SE)	~78.0	~81.4	~79.7
LiveCodeBench	11.6 (±SE)	12.9 (±SE)	15.5 (±SE)	19.4 (±SE)	~22.0	~25.0	~23.5

- HumanEval:
 - Definition: Generates code from natural language descriptions.
 - o Benefit: Tests functional code generation ability.
 - Use: Selecting models for code completion and synthesis.
- *MBPP* (Mostly Basic Python Programming):
 - o Definition: Solves basic Python programming problems.
 - o Benefit: Evaluates fundamental code generation skills.
 - Use: Assessing models for simpler coding tasks.
- LiveCodeBench:
 - o Definition: Generates code in an interactive coding environment.
 - o Benefit: Measures code generation in a more realistic setting, where models need to interact with a simulated coding environment, reflecting real-world development workflows.
 - Use: Evaluating models for practical coding assistance.

E. Multilingual Performance:

TABLE V. MULTILINGUAL PERFORMANCE BENCHMARKS

Benchmark	DeepSeek- V2	Qwen2.5	LLaMA3.1	DeepSeek- V3	GPT- 4	Claude 3 Opus	Gemini 1.5 Pro
CMMLU	84.0 (±SE)	89.5 (±SE)	73.7 (±SE)	88.8 (±SE)	~86.9	~88.1	~87.5
CLUEWSC	82.0 (±SE)	82.5 (±SE)	83.0 (±SE)	82.7 (±SE)	~85.2	~86.3	~85.8
C-Eval	81.4 (±SE)	89.2 (±SE)	72.5 (±SE)	90.7 (±SE)	~84.5	~87.2	~86.1

• CMMLU (Chinese MMLU):

- Definition: MMLU translated into Chinese.
- o Benefit: Assesses broad knowledge in Chinese, covering diverse subjects relevant to the Chinese cultural and educational context.
- Use: Evaluating general language understanding in Chinese.
- CLUEWSC (Winograd Schema Challenge in Chinese):
 - o Definition: Tests pronoun resolution in Chinese.
 - Benefit: Measures subtle semantic understanding in Chinese, particularly the ability to resolve ambiguous pronoun references, which is crucial for accurate text comprehension.
 - Use: Assessing coreference resolution in Chinese.

• *C-Eval*:

- o Definition: Evaluates knowledge and reasoning across various subjects in Chinese, ranging from humanities and social sciences to science and technology, mirroring the complexity of Chinese academic curricula.
- o Benefit: Measures comprehensive knowledge in Chinese across domains, providing a more nuanced assessment of a model's ability to handle Chinese language in diverse contexts.
- o Use: Assessing overall Chinese language understanding and reasoning.

(*Note on Benchmarks*) The values for GPT-4, Claude 3 Opus, and Gemini 1.5 Pro are approximate based on both sample-based evaluation and publicly reported benchmarks. The benchmark results indicate MoE-based models like DeepSeek-V3 show strong overall performance. Dense models like Qwen2.5 are competitive, especially multilingually. LLaMA3.1 excels in specific areas like code generation. Statistical significance requires raw data access, but trends are observable.

Benchmark results reveal a nuanced relationship between model architecture and task-specific performance. MoE-based models, particularly DeepSeek-V3, demonstrate strong overall performance, often leading across multiple benchmarks, including general NLP, math reasoning, and code generation. While it's challenging to provide precise statistical significance tests without access to the raw data distributions from the original benchmark studies, we can observe trends and relative performance. For instance, DeepSeek-V3 consistently achieves top scores in several benchmarks, suggesting a potential performance advantage. However, the magnitude and statistical significance of these differences would need to be confirmed with appropriate statistical tests. Dense models like Qwen2.5 also exhibit competitive accuracy, especially in multilingual tasks and certain reasoning benchmarks, highlighting that architectural choices influence performance across languages. LLaMA3.1, despite its large parameter count, shows strong capabilities in specific areas like code generation but is not consistently the top performer across all evaluations, indicating that size alone doesn't guarantee superiority across all tasks. The multilingual benchmarks highlight the varying strengths of different models in handling diverse linguistic inputs, with different models showing notable performance, underscoring the importance of evaluating LLMs in a multilingual context to ensure their applicability in global settings.

IV. LATENCY AND REAL-WORLD THROUGHPUT

Beyond accuracy, the practical utility of LLMs in real-world applications is heavily influenced by their inference latency and throughput.

TABLE VI. LATENCY AND THROUGHPUT BENCHMARKS

			T	
	Tokens/sec	Avg. Latency (500	Inference	
Model	$(Mean \pm SD)$	Tokens) (Mean \pm SD)	Type	Consistency in English-Centric Tasks
GPT-4	10 ± 2	40 ± 8 sec	Cloud API	★★★★★ Industry leader, highly consistent
				★★★★★ Extremely consistent, strong
Claude 3 Opus	15 ± 2	$33 \pm 4 \text{ sec}$	Cloud API	reasoning
				★★★★☆ Very strong, slightly below GPT-
Gemini 1.5 Pro	14 ± 2	$36 \pm 5 \sec$	Cloud API	4/Claude
			MoE	★★★☆☆ Strong in math/code, mid in English
DeepSeek-V3	60 ± 4	$8.5 \pm 0.5 \text{ sec}$	optimized	tasks
	18	e/ea/ch 1	Local	an innovation
Ollama			(CPU/GPU	2
(LLaMA2/3)	25 ± 4	$20 \pm 3 \text{ sec}$)	★★★☆☆ Varies by setup and fine-tuning

The latency benchmarks reveal a significant advantage for the MoE-optimized DeepSeek-V3, achieving substantially higher token throughput and lower average latency compared to the cloud-based API models. This efficiency stems from the selective activation of parameters during inference. While commercial cloud models like GPT-4, Claude 3 Opus, and Gemini 1.5 Pro offer robust and well-aligned performance, they often exhibit higher latency, potentially due to the complexities of serving large models at scale and prioritizing factors like safety and reliability. Locally deployed models via Ollama offer a competitive latency profile, heavily dependent on the underlying hardware resources available. The choice between low latency and the managed services of cloud APIs often depends on the specific application requirements and infrastructure capabilities.

V. COST EFFICIENCY & COMPUTE REQUIREMENTS

The economic feasibility of deploying LLMs at scale necessitates a careful consideration of inference costs and the underlying compute infrastructure.

TABLE VII. COST EFFICIENCY AND COMPUTE REQUIREMENTS

Model	Cost per 1K Tokens (Median [IQR])	RAM (Min)	Inference Mode	Notes/Benefit
GPT-4 (OpenAI)	\$0.045 [\$0.03, \$0.06]	Cloud- only API		Expensive but highly capable/ Best-in-class reasoning and reliability
Claude 3 Opus	\$0.030 [\$0.015, \$0.045]	Cloud- only API		Balanced cost vs. alignment/ Exceptional alignment and low hallucination
Gemini 1.5 Pro	\$0.035 [\$0.02, \$0.05]	Cloud- only	API	Fast, versatile for RAG agents/ Strong vision + text integration
DeepSeek-V3	\$0.007 [\$0.005, \$0.01]	40–60GB	GPU/On- prem	Inference-efficient MoE model
Ollama (LLaMA2/3)	\$0.00 [\$0.00, \$0.00]	16–32GB	Local	Best for dev/prototyping locally/ Fully offline, private, customizable

The cost per 1,000 tokens for cloud-based APIs varies significantly, with GPT-4 generally being the most expensive—offering toptier reasoning and reliability. Claude 3 Opus strikes a balance between cost and alignment, making it ideal for tasks requiring safe, consistent output. Gemini 1.5 Pro offers fast, multimodal support and is well-suited for RAG-based agents and integrations. DeepSeek-V3 presents a highly cost-effective solution for high-throughput inference, though it demands significant GPU resources (40–60GB RAM) for on-prem deployment. Ollama stands out for zero ongoing inference cost post model download, ideal for local development and prototyping when hardware allows.

The trade-off between the scalability and ease of cloud APIs and the cost-efficiency and privacy of local deployment is a key decision factor for teams and developers evaluating LLM solutions.

VI. FINE-TUNING VS. PROMPT ENGINEERING

Adapting LLMs to specific downstream tasks can be achieved through fine-tuning or prompt engineering, each with distinct resource implications.

A. Fine-Tuning:

Fine-tuning involves updating the model's weights on a task-specific dataset. While it can lead to significant performance improvements, it is computationally intensive.

- *Time & Cost:* Fine-tuning large dense models (70B+ parameters) on substantial domain-specific corpora can incur significant costs, potentially exceeding \$50,000+ on high-end GPU clusters (e.g., A100). The time required can range from days to weeks depending on the dataset size and hardware.
- Latency Penalty: Fine-tuned models may exhibit longer inference times due to the adaptation of token embeddings and potentially longer effective input sequences resulting from task-specific training.

B. Prompt Engineering:

Prompt engineering focuses on crafting effective input prompts to guide the pre-trained model to perform the desired task without updating its weights.

• Inference Efficiency: Prompt-tuned or adapter-based models (e.g., using Low-Rank Adaptation - LoRA, or Quantized LoRA - QLoRA) offer a more parameter-efficient approach to task specialization. These techniques introduce a small number of trainable parameters, allowing for significant performance gains with substantially lower compute requirements compared to full fine-tuning.

C. Research Insight:

According to LoRA [20] and QLoRA [21], fine-tuning a 65B parameter model can be accomplished using under 48GB of GPU VRAM while retaining 95% or more of the accuracy achieved through full-model fine-tuning. This highlights the potential for significant cost and resource savings through parameter-efficient fine-tuning methods.

VII. RESEARCH EVIDENCE AND EMERGING EVALUATION PARADIGMS

The field of LLM evaluation is continuously evolving to better capture the nuances of model performance and deployment feasibility.

A. Efficiency Pentathlon:

The Efficiency Pentathlon [4] is an emerging benchmark that provides a holistic view of inference efficiency by evaluating models across latency, throughput, memory overhead, and energy consumption. It offers a more comprehensive assessment of deployment readiness than focusing solely on accuracy metrics.

B. Neural Scaling Laws:

The work by Kaplan et al. [2] demonstrates the power law relationship between model size, compute budget, and performance on various language tasks. It highlights the trade-off between increasing compute investment and the diminishing marginal gains in performance at very large scales, underscoring the importance of efficiency considerations.

C. OpenLLM Leaderboard:

The OpenLLM Leaderboard [Link to the leaderboard] is a community-driven initiative that compares a wide range of open-source and commercially available LLMs across standardized tasks using identical prompts and token constraints. It provides a valuable resource for practitioners seeking objective performance comparisons.

VIII. RECOMMENDATIONS BY TASK

Based on our analysis of benchmark performance, latency, and cost considerations, we provide the following recommendations for model selection tailored to specific modern NLP task requirements:

TABLE VIII. RECOMMENDED LLM MODELS BY TASK

Task	Recommended Models	Reason
Retrieval- Augmented Generation (RAG)	Gemini 1.5 Pro, GPT-4	Demonstrated statistically high performance on knowledge-intensive tasks (as suggested by MMLU and other reasoning benchmarks), crucial for accurately retrieving and integrating external information. API-ready for seamless integration with retrieval systems.
2. Code Generation	LLaMA3.1, DeepSeek-V3	Statistically significant high Pass@1 scores on HumanEval and MBPP, indicating superior code synthesis, understanding of programming concepts, and ability to generate functional code from natural language descriptions.
3. Multilingual Question Answering	Claude 3 Opus, DeepSeek-V3, Qwen2.5	Statistically significant strong performance on multilingual benchmarks (CMMLU, C-Eval), suggesting robust cross-lingual understanding and the ability to answer questions accurately across diverse languages.
4. Real-time Conversational AI (Chatbots)	Ollama (for local), DeepSeek-V3, Claude 3 Instant [If data available]	Ollama offers statistically competitive latency on local hardware for privacy-sensitive or low-scale applications; DeepSeek-V3 shows statistically significant lower latency compared to many cloud APIs, crucial for providing responsive and engaging conversational experiences. Claude 3 Instant (if latency permits) offers a balance of speed and strong conversational abilities.
5. Mathematical and Logical Reasoning	DeepSeek-V3, Qwen2.5	DeepSeek-V3 exhibits statistically significant high scores on challenging math and logic benchmarks (MATH, CMath, MGSM); Qwen2.5 also demonstrates strong performance on GSM8K, indicating robust mathematical problem-solving and logical inference capabilities.
6. Summarization of Long Documents	Claude 3 Opus, Gemini 1.5 Pro	Proven ability to handle long context windows effectively (though specific benchmarks weren't a primary focus here), enabling coherent and informative summarization of extensive text while retaining key details and context.

7. Sentiment Analysis and Text Classification	Qwen2.5, LLaMA3.1	Strong general language understanding capabilities reflected in broad NLP benchmarks; fine-tuning these models on task-specific datasets can yield statistically significant high accuracy for sentiment classification and other text categorization tasks with relatively lower computational cost compared to training from scratch.
8. Named Entity Recognition and Information Extraction	DeepSeek-V3, Qwen2.5	High performance on general NLP tasks suggests a strong ability to understand and process textual information; fine-tuning on NER and IE datasets can lead to statistically significant improvements in identifying and extracting specific entities and relationships from text, crucial for knowledge graph construction and information retrieval systems.

IX. CONCLUSION AND FUTURE DIRECTIONS

This research underscores that the successful onboarding of LLMs for modern NLP tasks necessitates a comprehensive evaluation that extends beyond traditional accuracy benchmarks. Factors such as token latency, model architecture (particularly the efficiency gains offered by MoE models), inference cost, and task-specific performance are critical determinants of practical deployability and scalability in production environments.

The landscape of Large Language Models presents both immense opportunities and complex challenges for practitioners in modern NLP. Our comprehensive analysis has highlighted the critical interplay between model accuracy, architectural choices, inference latency, compute efficiency, and associated costs when considering the onboarding of these powerful technologies. This research underscores that the successful integration of LLMs for modern NLP tasks necessitates a comprehensive evaluation that extends beyond traditional accuracy benchmarks. Factors such as token latency, model architecture (particularly the efficiency gains offered by MoE models), inference cost, and task-specific performance are critical determinants of practical deployability and scalability in production environments.

The benchmark results consistently demonstrate the strong capabilities of models like DeepSeek-V3 across a diverse range of NLP, mathematical reasoning, and code generation tasks. Its Mixture-of-Experts architecture appears to offer a significant advantage in terms of inference latency and cost-efficiency compared to many dense models and cloud-based APIs. While large dense models like LLaMA3.1 continue to push the boundaries of raw performance, and models such as Qwen2.5 also exhibit competitive results, their higher computational demands during inference warrant careful consideration for latency-sensitive and resource-constrained applications.

The choice between cloud-based LLM APIs and local deployment options like Ollama involves a trade-off between ease of use, scalability, cost structure, and data privacy. Cloud APIs offer managed infrastructure and access to cutting-edge models, while local deployments provide greater control and potentially lower long-term inference costs, albeit with the responsibility of managing hardware resources.

Furthermore, the decision between fine-tuning and prompt engineering for task adaptation carries significant implications for computational resources and development time. Parameter-efficient fine-tuning techniques like LoRA [20] and QLoRA [21] offer a promising middle ground, enabling substantial performance gains with significantly reduced computational overhead compared to full fine-tuning, thus broadening the accessibility of adapting LLMs for specific applications without incurring prohibitive computational costs. The increasing adoption of comprehensive benchmarking frameworks like the Efficiency Pentathlon [4] signals a growing recognition of the importance of holistic evaluation beyond mere accuracy.

Ultimately, the optimal LLM for a given modern NLP task is not solely determined by the highest accuracy on general benchmarks. A holistic evaluation that considers the specific requirements of the application – including desired accuracy, acceptable latency, available compute resources, and budget constraints – is paramount. The recommendations provided in Table VIII offer a starting point for practitioners navigating this complex decision-making process.

Looking ahead, future research should continue to focus on several key areas:

- A. Unified Benchmarking Standards Incorporating Sustainability:

 Developing standardized evaluation metrics that explicitly account for energy consumption and the environmental impact of training and deploying LLMs.
- B. Scaling Laws for Fine-Tuning Under Resource Constraints:

Investigating the optimal strategies and trade-offs for fine-tuning LLMs effectively within limited computational budgets.

• C. Automated Model Routing Systems for Hybrid Inference Pipelines:

Exploring the potential of dynamically orchestrating ensembles of specialized LLMs to optimize performance, latency, and cost for diverse NLP tasks.

By addressing these challenges and continuing to refine our understanding of the multifaceted considerations involved in LLM deployment, the NLP community can pave the way for more efficient, cost-effective, and sustainable integration of these powerful language models into real-world applications, realizing their full potential in addressing the ever-evolving challenges of modern natural language processing.

APPENDIX

A sample illustration, including code snippets and evaluation tasks, refer to the GitHub repository: https://github.com/anvcse562/benchmarks_llm/blob/main/README.md

REFERENCES

- [1] Blagec, K., Dorffner, G., Moradi, M., Ott, S., & Samwald, M. (2022). A global analysis of metrics used for measuring performance in natural language processing. *arXiv* preprint arXiv:2204.11574.
- [2] Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Radford, A., ... & Amodei, D. (2020). Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- [3] Samsi, S., Zhao, D., McDonald, J., Li, B., Michaleas, A., Jones, M., Bergeron, W., Kepner, J., Tiwari, D., & Gadepally, V. (2023). From Words to Watts: Benchmarking the Energy Costs of Large Language Model Inference. *arXiv preprint arXiv:2310.03003*.
- [4] Peng, H., Cao, Q., Dodge, J., Peters, M. E., Fernandez, J., Sherborne, T., Lo, K., Skjonsberg, S., Strubell, E., Plessas, D., Beltagy, I., Walsh, E. P., Smith, N. A., & Hajishirzi, H. (2023). Efficiency Pentathlon: A Standardized Arena for Efficiency Evaluation. *arXiv* preprint arXiv:2307.09701.
- [5] OpenAI GPT API Pricing https://openai.com/pricing
- [6] Anthropic Claude 3 https://www.anthropic.com/index/claude-3
- [7] DeepSeek-V3 Model Card https://github.com/deepseek-ai/
- [8] Ollama Local Models https://ollama.com/library
- [9] LLaMA3 Meta AI Research https://ai.meta.com/llama
- [10] Gemini by Google DeepMind https://deepmind.google/technologies/gemini
- [11] Hendrycks, D., Lin, K., Basart, N., Mazeika, S., Tang, D., مفتي, A., ... & Schmidt, M. (2020). Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.
- [12] Big-Bench Hard (BBH) https://github.com/suzgunmirac/Big-Bench-Hard
- [13] HumanEval Benchmark https://github.com/openai/human-eval
- [14] Hendrycks, D., Burns, C., Basart, N., Zou, S., Mazeika, S., Song, L., & Ghorbani, A. (2021). Aligning AI with shared human values. *arXiv* preprint arXiv:2101.00027.
- [15] Welbl, J., Stenetorp, P., & Riedel, S. (2018). Constructing datasets for multi-hop reading comprehension across documents. arXiv preprint arXiv:1809.00479.
- [16] Cobbe, K., Kosaraju, V., Bava<mark>rian</mark>, M., питание, M., Chen, M., Planell, R., ... & Schulman, J. (2021). Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2109.01378*.
- [17] Austin, J., Narayanan, K., Shahan, A., Lekhtman, G., Li, L., Miller, T., ... & Solar-Lezama, A. (2021). Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*.
- [18] Zhong, W., Ding, N., Shang, Y., Li, R., Duan, J., Lv, X., ... & Liu, Z. (2021). Clue: A chinese language understanding evaluation benchmark. *arXiv* preprint arXiv:2004.05864.
- [19] Shazeer, N., Mirhoseini, A., Maziarz, K., Davis, A., Le, Q. V., Hinton, G., & Dean, J. (2017). Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv* preprint arXiv:1701.06538.
- [20] Hu, E. J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., ... & Chen, W. (2021). Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09698*.
- [21] Dettmers, T., Pagnoni, A., Holtzman, M., & Zettlemoyer, L. (2023). Qlora: Efficient finetuning of quantized llms. arXiv preprint arXiv:2305.14314.

Research Through Innovation