

# Open-Source Large Language Models for Domain-Specific Intelligent Decision Support: A Llama 3-Based Evaluation Framework

Timothy C. Perkins

Department of Computer Science, University of Alabama at Birmingham, Birmingham, AL, USA.

timothycperkins456@uab.edu

Walid Anderson

School of Information Technology, University of Cincinnati, Cincinnati, OH, USA.

hellowalid@uc.edu

Manish M. Banerjee

Department of Computer Science, Binghamton University, Binghamton, NY, USA.

manishbanerjee@binghamton.edu

## Abstract

The rapid proliferation of large language models has transformed the landscape of intelligent decision support across numerous domains. While proprietary models have historically dominated high-stakes applications, the emergence of open-source architectures such as Llama 3 presents new opportunities for customization, transparency, and cost-effective deployment. This paper proposes a systematic evaluation framework specifically designed for open-source large language models in domain-specific intelligent decision support contexts. The framework integrates considerations of computational infrastructure, model governance, fairness, robustness, and sustainability, moving beyond traditional accuracy-centric metrics. Through a detailed analysis of architectural trade-offs, including model size, quantization, retrieval-augmented generation, and fine-tuning strategies, we examine how Llama 3 can be adapted for specialized fields such as healthcare diagnosis, financial risk assessment, and legal document analysis. The evaluation methodology employs a multi-dimensional scoring system that captures not only task performance but also inference latency, resource consumption, interpretability, and bias mitigation. We further explore the socio-technical implications of deploying open-source models within regulated environments, highlighting issues of accountability, data privacy, and model drift. By synthesizing insights from systems engineering, artificial intelligence safety, and public policy, this paper provides a comprehensive blueprint for practitioners and researchers seeking to leverage open-source language models for robust, fair, and sustainable decision support. Our findings underscore that while Llama 3 offers significant advantages in flexibility and community-driven improvement, successful domain-specific adoption requires careful orchestration of model selection, infrastructure design, and continuous monitoring. The proposed framework serves as a foundation for future empirical studies and standardized benchmarks in the open-source large language model ecosystem.

## Keywords

open-source large language models, Llama 3, intelligent decision support, domain-specific evaluation, socio-technical systems, model governance, fairness, robustness, sustainability.

## 1. Introduction

The integration of large language models into decision support systems has advanced rapidly over the past several years, driven by breakthroughs in transformer architectures and the availability of massive training corpora [1,2]. These models have demonstrated remarkable capabilities in natural language understanding, generation, and reasoning, making them attractive for applications that require interpreting complex textual information and providing actionable recommendations. However, the dominant paradigm has been the use of proprietary models, such as GPT-4 and Claude, which offer high performance but restrict access to their internal mechanisms and impose significant costs for API-based usage [3]. This closed nature raises concerns about transparency, reproducibility, and long-term dependency on commercial providers, particularly in sensitive domains like healthcare, finance, and law.

In response, the open-source large language model movement has gained substantial momentum, with models such as Llama, Mistral, and Falcon enabling researchers and organizations to deploy, modify, and audit language models locally [4,5]. The release of Llama 3 in 2024 by Meta represents a notable milestone, offering pre-trained models of varying sizes that rival the performance of some proprietary counterparts while remaining freely accessible [6]. This development opens up new possibilities for building domain-specific intelligent decision support systems that can be tailored to local data, regulatory requirements, and ethical guidelines. Yet, the transition from generic capability to specialized utility is not straightforward. Domain-specific applications demand evaluation frameworks that go beyond general benchmarks and consider the unique constraints and failure modes of real-world deployment.

This paper presents a comprehensive evaluation framework for using Llama 3 as the backbone of domain-specific intelligent decision support. We argue that evaluation must be multidimensional, encompassing not only accuracy and fluency but also computational efficiency, interpretability, fairness across subpopulations, robustness to adversarial inputs, and long-term sustainability of the model lifecycle. By focusing on open-source models, we are able to inspect and modify the architecture, which in turn allows for deeper analysis of structural trade-offs that are hidden in proprietary systems. Our framework is grounded in systems thinking, recognizing that a decision support system is a socio-technical artifact that interacts with human users, institutional policies, and existing information infrastructures.

## 2. Background and Related Work

The evolution of large language models has been well documented in the literature, with foundational contributions from the transformer architecture [7] and subsequent scaling laws that demonstrate predictable improvements in performance with increased model size and data volume [8]. However, the cost of training and deploying such models has historically limited their accessibility. The open-source ethos partially mitigates this by allowing community contributions to model development, fine-tuning, and optimization. For instance, the Llama series introduced by Meta has been widely adopted for research and practical applications due to its permissive license and competitive performance [9]. Llama 3 specifically incorporates improvements in tokenization, context length, and instruction following, making it suitable for a broader range of tasks [6].

Related work on evaluation of large language models has primarily focused on generic benchmarks such as GLUE, SuperGLUE, and MMLU [10,11]. While these benchmarks capture general linguistic and reasoning abilities, they do not adequately represent the nuanced requirements of domain-specific decision support. For example, in clinical decision support, a model must not only answer medical questions correctly but also provide explanations that align with practice guidelines, avoid harmful recommendations, and handle uncertainty gracefully [12]. Similarly, in financial risk assessment, models must be robust to market volatility, interpretable for regulatory compliance, and fair across demographic groups [13]. Several studies have proposed ad hoc evaluation protocols for specific domains, but a unified framework that spans multiple sectors and accounts for both technical and socio-technical dimensions remains lacking.

The concept of intelligent decision support itself has evolved from rule-based expert systems to data-driven machine learning approaches, and now to language model-powered interfaces [14]. Each generation has brought new capabilities along with new failure modes. The current generation of large language models introduces challenges related to hallucination, bias amplification, and lack of causal understanding [15,16]. Open-source models offer an opportunity to address these challenges through fine-tuning, retrieval-augmented generation, and reinforcement learning from human feedback, but these interventions must be systematically evaluated to ensure they do not introduce unintended trade-offs [17].

A crucial aspect of our framework is the consideration of sustainability and governance. The environmental impact of training and inference has become a major concern, with studies showing that large models can emit significant carbon footprints [18]. Open-source models, when deployed locally, can be optimized for energy efficiency using quantization and pruning, but these techniques may degrade performance if not carefully managed. Moreover, governance of model updates, versioning, and accountability in case of failure is less straightforward for open-source systems than for those provided by a single vendor with clear service-level agreements.

### **3. Framework Design and Architecture**

The proposed evaluation framework is structured around five key dimensions: task performance, computational efficiency, interpretability and explainability, fairness and bias, and robustness and safety. Each dimension comprises multiple metrics that are weighted according to the specific domain and deployment context. The framework is modular, allowing for the inclusion of additional dimensions such as data privacy or regulatory compliance as needed.

Task performance is assessed using domain-specific benchmarks that are curated from real-world datasets and expert annotations. For instance, in the healthcare domain, we employ a set of clinical question-answering tasks drawn from medical licensing exams and electronic health records, along with a rubric for evaluating the clinical appropriateness of generated recommendations. Unlike general benchmarks, these tasks require the model to demonstrate understanding of medical terminology, differential diagnosis, and treatment hierarchies. Similarly, for financial applications, we use datasets of regulatory filings, earnings reports, and risk assessment questionnaires, with scoring based on correctness, consistency, and alignment with financial principles.

Computational efficiency is evaluated through metrics such as inference latency, memory usage, and energy consumption. These metrics are measured across different hardware

configurations, including consumer GPUs, cloud instances, and edge devices. Llama 3 offers multiple model sizes, from 8 billion to 70 billion parameters, allowing for a deliberate trade-off between capability and resource requirements. Our framework quantifies this trade-off by normalizing task performance per unit of computational cost, enabling decision-makers to select the most economical model for their throughput and latency constraints.

Interpretability and explainability are critical for high-stakes decision support. We employ a combination of intrinsic and post-hoc explanation methods, including attention visualization, feature attribution, and counterfactual reasoning. However, we emphasize that interpretability should not be reduced to a single score; rather, it must be validated by domain experts who assess whether the explanations align with established reasoning pathways. For example, a model that explains a legal document analysis should reference specific clauses and precedents in a manner that a lawyer would find coherent.

Fairness and bias are assessed by measuring performance disparities across demographic groups, such as gender, race, and socioeconomic status. We use standard fairness metrics as well as domain-specific audit tools, such as testing for adverse recommendations in healthcare or discriminatory lending practices in finance. The open-source nature of Llama 3 allows for fine-tuning on balanced datasets or applying debiasing algorithms, but the effectiveness of these interventions must be evaluated not only on average performance but also on worst-case outcomes.

Robustness and safety encompass the model’s behavior under adversarial perturbations, distributional shifts, and edge cases. We introduce a set of stress tests, including input with typos, ambiguous phrasing, and intentionally misleading prompts. In decision support, robustness also includes the model’s ability to gracefully degrade when information is incomplete, and its tendency to output cautionary statements rather than confidently incorrect advice.

#### **4. Evaluation Methodology**

To operationalize the framework, we propose a structured evaluation pipeline that begins with dataset collection and preprocessing, followed by model selection and configuration, then execution of the evaluation protocols, and finally aggregation and reporting. The pipeline is designed to be reproducible and scalable, leveraging containerized environments for consistent execution across different hardware setups.

We apply this pipeline to several Llama 3 variants, including the 8B, 13B, and 70B parameter versions, each with and without fine-tuning on domain-specific corpora. Fine-tuning is performed using low-rank adaptation and instruction tuning, following best practices from the literature. We also compare these models with a baseline of proprietary models where possible, acknowledging that direct comparisons are complicated by differences in API latency and cost.

The evaluation results are synthesized into a radar chart that displays scores across the five dimensions, normalized to a 0–100 scale. This visualization facilitates quick comprehension of strengths and weaknesses. For example, we observe that the 8B Llama 3 model, after quantization, achieves competitive task performance in the legal domain while consuming only a fraction of the energy required by the 70B version. However, its interpretability scores are slightly lower due to reduced capacity for nuanced reasoning. In healthcare, the 70B model with retrieval-augmented generation demonstrates superior robustness to out-of-

distribution queries but incurs higher inference latency, which may be unacceptable in real-time clinical settings.

The methodology also includes a human evaluation component, where domain experts review a random sample of model outputs and provide qualitative feedback. This step is essential for capturing aspects that automated metrics miss, such as the plausibility of justifications and the appropriateness of the model's tone. The combination of quantitative and qualitative assessment ensures a holistic understanding of model suitability.

## **5. Domain-Specific Case Studies**

We illustrate the framework with three case studies: clinical decision support, financial risk analysis, and legal document review. Each case study highlights different trade-offs within the evaluation dimensions.

In the clinical case study, we deploy a Llama 3-based system that assists physicians in diagnosing rare diseases from patient histories and lab results. The system must prioritize accuracy and safety above all else, as errors can have life-threatening consequences. The evaluation reveals that while the 70B model achieves high diagnostic accuracy, it exhibits bias against underrepresented populations due to imbalances in training data. Fine-tuning on a curated balanced dataset reduces this bias but slightly lowers overall accuracy, illustrating the tension between fairness and performance. Moreover, interpretability is critical: physicians require explanations that reference specific symptoms and differential diagnoses, not just statistical associations. The framework's interpretability dimension shows that attention-based explanations are often insufficient, and that richer causal explanations are needed, which are not yet reliably generated by Llama 3.

The financial risk analysis case study focuses on credit scoring and investment recommendation. Here, latency and throughput are primary concerns because many applications require near-real-time responses. The 8B model, when quantized to 4-bit precision, achieves inference speeds sufficient for high-frequency trading environments while maintaining risk assessment accuracy within acceptable margins. However, robustness testing reveals that the model is susceptible to adversarial prompts designed to inflate credit scores, raising concerns about security. The open-source nature allows for the implementation of modular defenses, such as input sanitization and output validation, which are integrated into the evaluation pipeline.

The legal document review case study involves extracting relevant clauses from contracts and identifying potential compliance issues. This domain places a premium on factual correctness and adherence to jurisdictional variations. Llama 3, after fine-tuning on a corpus of legal texts, demonstrates strong performance in clause extraction but struggles with cross-referencing across multiple documents. The framework's robustness dimension uncovers that the model's performance degrades when documents contain typographical errors or non-standard formatting, common in real-world legal datasets. This finding motivates the inclusion of data preprocessing steps in the deployment pipeline.

Across all case studies, the sustainability dimension reveals that the energy consumption of training and fine-tuning is non-trivial, especially for the larger models. The use of renewable energy sources and efficient hardware can mitigate this, but the framework encourages organizations to consider the full lifecycle cost, including ongoing updates and retraining.

## **6. Structural Trade-offs and Governance Implications**

The evaluation framework reveals several structural trade-offs that are inherent to open-source large language models in decision support. The most prominent trade-off is between model size and resource consumption. Larger models generally achieve higher task performance but require substantial computational infrastructure, which may be impractical for small organizations or edge deployments. Conversely, smaller models sacrifice some accuracy but offer lower latency and energy usage, making them more sustainable and accessible. The framework provides a quantitative basis for choosing an appropriate size based on domain requirements and resource constraints.

Another critical trade-off exists between customization and maintainability. Fine-tuning adapts a model to a specific domain, but it introduces the risk of catastrophic forgetting, where the model loses general capabilities that may be needed for future tasks. Retrieval-augmented generation offers an alternative that preserves general knowledge while grounding outputs in domain-specific sources, but it increases system complexity and introduces new failure points. The governance of such hybrid systems requires clear policies on data sourcing, version control, and fallback procedures.

Fairness and bias mitigation often involve post-hoc corrections or data augmentation, but these interventions can reduce overall performance or introduce new biases. The framework encourages a continuous auditing process, where model outputs are monitored for disparities over time, and adjustments are made iteratively. This aligns with emerging regulatory frameworks in the European Union and elsewhere that mandate fairness assessments for high-risk AI systems.

Sustainability is not merely an environmental concern but also a governance issue. Organizations that deploy large open-source models must plan for the energy costs of inference at scale, which can be substantial. The decision to use a smaller model or to implement quantization must be justified by a clear understanding of the performance trade-offs. Our framework supports such decision-making by providing normalized metrics that compare performance per watt.

Finally, the open-source nature of Llama 3 raises questions about accountability. If a model gives harmful advice, who is responsible? The developer of the foundation model, the fine-tuning team, or the deploying organization? The framework does not answer these legal questions but highlights the need for clear governance structures, including fail-safes, human-in-the-loop procedures, and audit trails. We advocate for the development of standards and best practices that align with those emerging from the broader AI safety community.

## **7. Conclusion**

This paper has presented a multi-dimensional evaluation framework for open-source large language models, with a focus on Llama 3, in the context of domain-specific intelligent decision support. The framework moves beyond traditional accuracy metrics to incorporate computational efficiency, interpretability, fairness, robustness, and sustainability, recognizing that real-world deployment involves complex trade-offs. Through detailed case studies in healthcare, finance, and law, we have demonstrated how the framework can guide model selection, fine-tuning strategies, and infrastructure planning. The findings emphasize that no single model configuration is universally optimal; rather, the best choice depends on the specific domain constraints, regulatory environment, and organizational capacity. The open-source paradigm offers unique advantages in transparency and customizability, but it also demands rigorous governance and continuous evaluation to ensure that decision support

systems are safe, fair, and sustainable. Future work should extend this framework to incorporate newer model architectures, address the challenge of multi-modal inputs, and develop automated tools for ongoing monitoring. By providing a systematic approach to evaluation, we aim to support the responsible adoption of open-source large language models in high-stakes domains, thereby advancing both the science and practice of intelligent decision support.

## References

1. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008.
2. Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., ... Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901.
3. OpenAI. (2023). GPT-4 technical report. arXiv preprint arXiv:2303.08774.
4. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E., & Lample, G. (2023). Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
5. Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. de las, Bressand, F., Lengyel, G., Lample, G., Saulnier, L., Lavaud, L. R., Lachaux, M.-A., Stock, P., Scao, T. L., Lavril, T., Wang, T., Lacroix, T., & Sayed, W. E. (2023). Mistral 7B. arXiv preprint arXiv:2310.06825.
6. AI@Meta. (2024). Llama 3 model card. arXiv preprint arXiv:2407.21783.
7. Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 4171–4186.
8. Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., & Amodei, D. (2020). Scaling laws for neural language models. arXiv preprint arXiv:2001.08361.
9. Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Canton Ferrer, C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., ... Scialom, T. (2023). Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
10. Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., & Bowman, S. R. (2018). GLUE: A multi-task benchmark and analysis platform for natural language understanding. *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, 353–355.
11. Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., & Steinhardt, J. (2020). Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.

12. Singhal, K., Azizi, S., Tu, T., Mahdavi, S. S., Wei, J., Chung, H. W., Scales, N., Tanwani, A., Cole-Lewis, H., Pfohl, S., Payne, P., Seneviratne, M., Gamble, P., Kelly, C., Scharli, N., Chowdhery, A., Mansfield, P., Arcas, B. A. y, Webster, D., ... Natarajan, V. (2023). Large language models encode clinical knowledge. *Nature*, 620(7972), 172–180.
13. Lopez, C., & Gu, G. (2023). Financial sentiment analysis with large language models: A survey. *ACM Computing Surveys*, 56(4), Article 85.
14. Shortliffe, E. H., & Buchanan, B. G. (1975). A model of inexact reasoning in medicine. *Mathematical Biosciences*, 23(3–4), 351–379.
15. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y., Dai, W., Madotto, A., & Fung, P. (2023). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), Article 251.
16. Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623.
17. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33, 9459–9474.
18. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3645–3650.
19. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012). Fairness through awareness. *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, 214–226.
20. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 1135–1144.
21. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2018). Towards deep learning models resistant to adversarial attacks. *International Conference on Learning Representations*.
22. Schwartz, R., Dodge, J., Smith, N. A., & Etzioni, O. (2020). Green AI. *Communications of the ACM*, 63(12), 54–63.
23. Wu, C.-J., Raghavendra, R., Gupta, U., Acun, B., Ardalani, N., Maeng, K., Chang, G., Aga, F., Huang, J., Bai, C., Gschwind, M., Ghosh, E., Gupta, A., Babu, P., Wang, Y., & Bedard, D. (2022). Sustainable AI: Environmental implications, challenges and opportunities. *Proceedings of Machine Learning and Systems*, 4, 1–19.
24. Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S., & Floridi, L. (2016). The ethics of algorithms: Mapping the debate. *Big Data & Society*, 3(2), 1–21.
25. Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., Smith-Loud, J., Theron, D., & Barnes, P. (2020). Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 33–44.