

Leveraging the Complexity and Criticality Matrix as a Governance Framework for Compliance in Generative AI Systems

Saiprapul R Thotapally

Global Payments, Georgia, USA

Abstract - Generative AI systems, particularly those powered by Large Language Models (LLMs), are increasingly integrated into compliance-heavy enterprise workflows, including financial audits, healthcare compliance checks, and data privacy reviews. While these models excel at producing contextually rich and fluent responses, their reliance on probabilistic patterns rather than semantic understanding poses significant challenges. Errors such as hallucinations, biases, and factual inaccuracies risk severe legal, financial, and reputational consequences.

This thesis introduces the Complexity and Criticality Matrix, a governance framework designed to guide responsible AI deployment in regulated environments. By classifying tasks according to complexity (domain specificity) and criticality (compliance impact), the framework prescribes oversight strategies: direct LLM responses for low-risk tasks, Retrieval-Augmented Generation (RAG) for improved factual grounding where needed, and Human-in-the-Loop (HITL) interventions for high-stakes, compliance-critical scenarios. This structured approach bridges the gap between high-level AI governance principles and practical operationalization, enabling organizations to integrate oversight seamlessly within DevOps/MLOps pipelines.

While this work is primarily conceptual, it lays the groundwork for future empirical validation and scalability. By providing a method to balance automation with accuracy and regulatory adherence, this thesis aims to empower enterprises to confidently leverage generative AI in even the most stringent compliance contexts.

Key Words: Generative AI, Compliance, Retrieval-Augmented Generation (RAG), Human-in-the-Loop (HITL), DevOps, MLOps, AI Governance

1. INTRODUCTION

Large Language Models (LLMs) such as GPT-4 and BERT have demonstrated remarkable capabilities in producing contextually relevant and fluent responses [1], [2]. As these systems find use in regulated sectors—financial auditing, healthcare compliance, and data privacy requests—organizations must ensure not only efficiency but also strict adherence to legal and ethical standards [3], [4]. Misinterpretations, hallucinations, and biases can lead

to serious repercussions, including legal penalties, financial losses, and reputational damage [5], [6].

Motivation

High-level governance frameworks (e.g., NIST AI RMF [11], OECD AI Principles [15], ISO/IEC standards [12]) emphasize transparency, accountability, and fairness. However, they offer limited guidance on turning these abstract principles into daily operational practices. Enterprises struggle with questions like: When can LLM outputs be trusted as-is? When is it necessary to bolster responses with RAG for factual correctness? When should a human reviewer intervene for compliance-critical decisions?

Contribution

This thesis addresses these gaps by introducing the Complexity and Criticality Matrix, a governance framework that classifies tasks according to their complexity (domain-specific knowledge required) and criticality (severity of errors). Each quadrant of the matrix is associated with a different oversight strategy:

- Low Complexity/Low Criticality: Direct LLM outputs for efficiency.
- High Complexity/Low Criticality: Selective RAG to ensure factual grounding.
- Low Complexity/High Criticality: Prompt templates with minimal HITL checks for compliance accuracy.
- High Complexity/High Criticality: Combined RAG and HITL interventions plus auditing for top-tier regulatory adherence.

This structured approach operationalizes governance principles, enabling organizations to integrate compliance checks into DevOps/MLOps pipelines [7], [8]. While no empirical results are presented here, the framework sets a foundation for future validations, including heuristic or zero-shot classification approaches for automating task categorization.

2. BACKGROUND AND RELATED WORK

2.1 LLM Limitations and RAG

LLMs are known to hallucinate and produce unsupported statements, primarily because they generate text from probabilistic patterns rather than verified facts [3], [4]. This limitation is risky in compliance-heavy contexts, where accuracy is paramount [5]. Retrieval-Augmented Generation (RAG) has emerged as a solution to enhance factual grounding by integrating authoritative external knowledge sources [9]. While RAG improves accuracy, it introduces additional latency and infrastructure complexity, necessitating selective application based on task requirements.

2.2 Human Oversight, Prompt Engineering, and Beyond

Prompt engineering, model cards, and bias audits have improved transparency and reduced unwanted biases [17]–[19]. However, these approaches alone may not suffice for high-risk scenarios. Human-in-the-Loop (HITL) interventions introduce expert judgment and domain expertise, serving as a last line of defense against non-compliant outputs [10], [13]. Yet, no widely accepted method exists to determine when to activate HITL or RAG optimally, highlighting the need for a structured framework like the Complexity and Criticality Matrix.

2.3 Governance Frameworks and Implementation Gaps

Established governance frameworks and standards—NIST AI RMF [11], ISO/IEC guidelines [12], and proposed AI regulations [15], [16]—emphasize accountability, transparency, and risk management. However, these documents lack granular, actionable tools for daily operations in regulated environments. Enterprises need a mechanism to translate abstract principles into operational workflows.

The Complexity and Criticality Matrix proposed in this thesis addresses these challenges by mapping task attributes—complexity and criticality—to tailored oversight strategies. This structured approach bridges the gap between abstract governance principles and practical implementation, enabling organizations to deploy generative AI responsibly and effectively.

3. METHODOLOGICAL APPROACH

3.1 Framework Overview

The Complexity and Criticality Matrix classifies tasks by their domain complexity and compliance-criticality. Complexity scales with domain specificity and required expertise, while criticality measures the severity

of errors. By defining these attributes for each query, organizations can determine the appropriate oversight level.

3.2 Quadrant-Based Oversight Strategies

- **Low Complexity, Low Criticality:** (e.g., basic FAQs)
 - Oversight: Direct LLM outputs, prioritizing speed and cost-efficiency.
- **High Complexity, Low Criticality:** (e.g., complex informational queries)
 - Oversight: Moderate RAG or refined prompt engineering for improved factual grounding.
- **Low Complexity, High Criticality:** (e.g., eligibility checks under GDPR)
 - Oversight: Strict prompt templates, partial HITL for compliance verification.
- **High Complexity, High Criticality:** (e.g., interpreting regulatory filings under HIPAA or SEC rules)
 - Oversight: RAG for authoritative information plus HITL validation and audit logs to ensure compliance and audit logs.

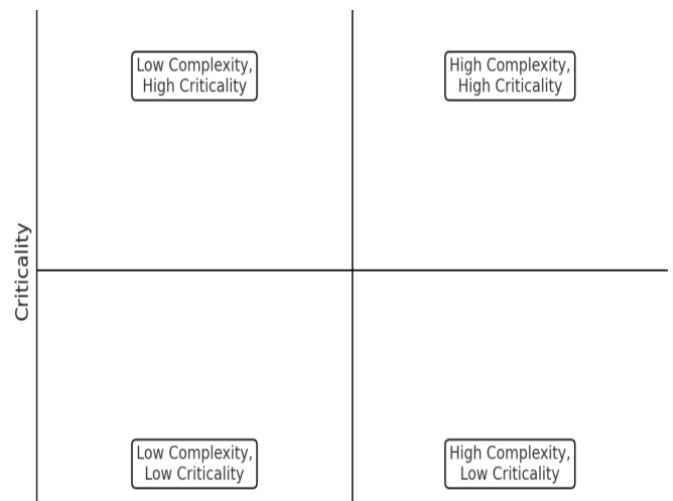


Fig -1: Complexity and Criticality Matrix visualization

3.3 Theoretical Foundation

This approach aligns with risk-based governance, where oversight intensity scales with potential harm [11], [15]. By explicitly categorizing tasks, the framework

remains flexible, scalable, and adaptable to evolving regulations. It operationalizes governance principles, providing a path from abstract guidelines to concrete workflows.

4. Practical Applications and Use Cases

4.1 Matrix as a Query Redirection and Governance Framework

The Complexity and Criticality Matrix serves as both a query redirection tool and an internal governance framework in AI systems:

- **Query Redirection:** By classifying queries based on their complexity and criticality, the matrix routes them to the most suitable model, whether that's LLM, RAG, or HITL, ensuring accurate and efficient responses.
- **Internal Governance:** The matrix also acts as a governance mechanism, ensuring that AI-generated responses meet industry-specific and region-specific compliance standards, applying the appropriate levels of oversight (e.g., RAG for factual grounding or HITL for compliance validation).

This dual functionality makes the matrix a powerful tool for organizations, allowing them to handle queries with compliance and efficiency while ensuring human oversight where necessary.

Retrieval-Augmented

Generation:

Implement RAG by maintaining a vector database of authoritative documents. On receiving a query, perform a vector similarity search and feed the retrieved passages into the LLM prompt. Although this boosts accuracy, it adds infrastructure and latency overhead.

Human-in-the-Loop:

For critical outputs, route preliminary RAG-based responses to compliance officers or domain experts. This ensures human judgment and accountability in high-stakes scenarios.

4.2 Use Cases Across Various Industries

The Complexity and Criticality Matrix adapts to different industries, each with unique compliance needs, by categorizing tasks and applying oversight mechanisms to ensure regulatory adherence.

- **Healthcare** (e.g., HIPAA): High complexity/high criticality tasks, such as explaining privacy law penalties, use RAG plus HITL. Routine summaries (low complexity/low criticality) rely on direct LLM outputs.

- **Finance** (e.g., SEC Filings): Complex regulatory interpretations trigger RAG and HITL, while simple financial FAQs are managed by LLM outputs.

- **Enterprise Monitoring:** Routine log summaries (LL quadrant) are handled by LLMs, while suspicious security anomalies (HC quadrant) undergo HITL review.

4.3 Technical Integration and Workflow in DevOps/MLOps Pipelines

In DevOps/MLOps pipelines [7], [8], the matrix classification can be automated. Queries are tagged according to complexity and criticality, triggering the appropriate oversight strategy. Gated deployments and continuous integration checks ensure compliance standards are maintained.

4.4 Summary

The Complexity and Criticality Matrix serves as a dual-purpose framework that not only directs queries to the appropriate model but also ensures compliance with industry and regional regulations. It acts as a governance tool by applying oversight where necessary, whether for high-complexity/high-criticality tasks requiring RAG or HITL or for low-complexity/low-criticality tasks handled by LLM outputs. This ensures that organizations can scale AI systems while maintaining accuracy, accountability, and regulatory compliance across industries like healthcare, finance, and enterprise monitoring.

MLOps Pipeline for Complexity and Criticality Matrix

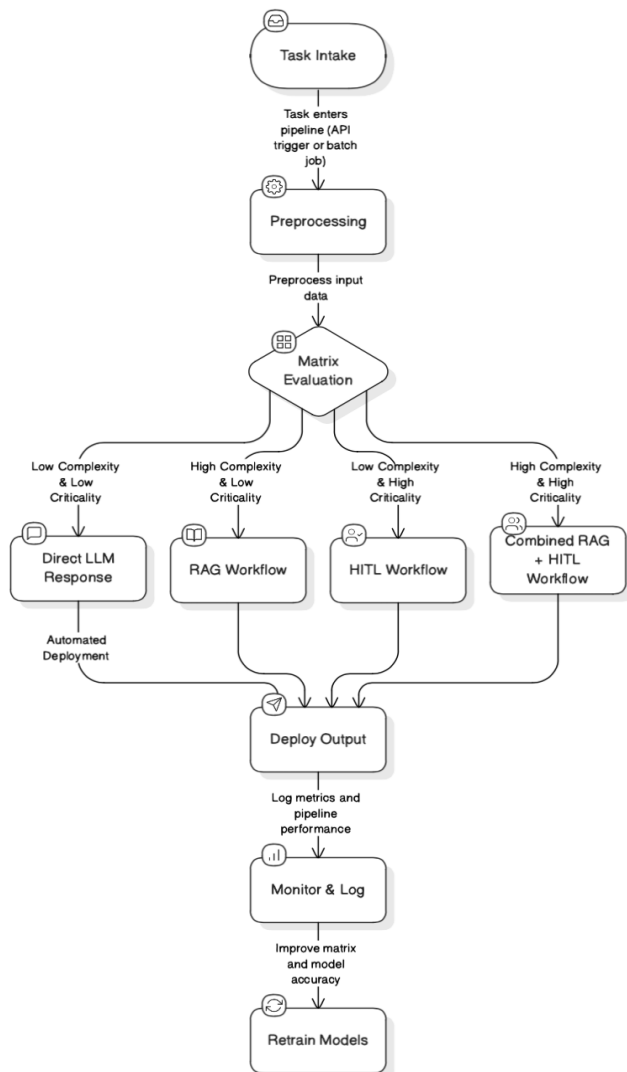


Fig -2: MLOps Pipeline for Complexity and Criticality Matrix

5. COST AND TIME EFFICIENCY CONSIDERATIONS

Baseline LLM inference may be cost-effective but less reliable for complex or critical tasks [2], [4]. RAG introduces retrieval operations, adding infrastructure costs and milliseconds of latency [9]. HITL adds human review time and labor costs. By applying these resources selectively—based on the matrix—organizations can minimize overhead and ensure that costs scale only with risk level.

For instance, vector retrieval may cost ~\$0.002–\$0.01 per query with slight latency additions [9], while LLM inference might cost ~\$0.03–\$0.10 per 1k tokens [2]. The matrix ensures RAG and HITL are reserved for queries where improved accuracy justifies extra cost and latency.

6. ALIGNMENT WITH GOVERNANCE STANDARDS

The Complexity and Criticality Matrix operationalizes governance principles from frameworks like NIST AI RMF [11], ISO/IEC standards [12], and proposed EU AI legislation [16]. While governance frameworks emphasize transparency, accountability, and fairness, they often lack actionable guidelines for practical implementation in daily AI operations. Many organizations struggle to operationalize these principles, often finding themselves caught in long cycles of planning and documentation.

This thesis introduces a practical governance solution that enables organizations to align with high-level governance principles while providing concrete steps for deployment.

Addressing the Gap Between Theory and Practice:

- Governance Principles:** Existing frameworks outline governance principles, but organizations face challenges in applying these principles in real-world scenarios, particularly in regulated industries.
- The Practical Solution:** The Complexity and Criticality Matrix operationalizes governance, providing clear guidelines on when to trust LLM outputs, when to use RAG, and when to involve HITL review. This approach makes it easier for organizations to implement and scale AI governance without getting bogged down in theoretical frameworks or over-engineered processes.

6.1 Compliance Alignment:

By categorizing tasks according to complexity and criticality, the matrix provides a straightforward way for organizations to ensure compliance at scale. Continuous oversight, audit trails, and HITL interventions can be integrated in a manner that aligns with governance frameworks and regulatory requirements. This reduces risks of non-compliance and allows organizations to adjust oversight dynamically as they gain trust in the AI system's performance.

6.2 Governance Efficiency:

The Matrix provides an operational framework that allows organizations to meet governance standards efficiently, balancing automated processes with human oversight. Unlike theoretical frameworks that may leave organizations in an implementation gap, the Complexity and Criticality Matrix serves as a practical tool, ensuring organizations are both compliant and efficient in deploying AI.

7. Challenges and Limitations

- **Task Classification Accuracy:** Automatically classifying queries into complexity/criticality tiers remains non-trivial. Future research may explore zero-shot classification or heuristic rule-based approaches to improve accuracy [20], [21]. Feedback loops from HITL interventions can refine classification models over time.
- **Balancing Automation and Oversight:** While the matrix provides a structured approach, finding the optimal balance between automation and human review may be challenging. Ongoing evaluation and adjustments are necessary as trust in AI systems evolves.
- **Scalability and Dynamic Regulations:** As regulations evolve, thresholds and classification criteria may need recalibration. Future iterations could adapt dynamically, updating the RAG index and adjusting complexity/criticality tiers in response to regulatory changes [16].
- **Experimental Validation and Tooling:** While this thesis lays the conceptual foundation, future work will provide empirical validation, user studies, and reference implementations. Subsequent research can test heuristic or zero-shot classification models for query categorization and integrate explainability tools to improve user trust [20], [21].

8. FUTURE DIRECTIONS

Rigorous User Studies: Validate the Complexity and Criticality Matrix by collaborating with domain experts, testing it on real-world compliance documents, and assessing its performance across various industries and regulatory contexts.

Advanced Reasoning Integration: Explore symbolic reasoning or knowledge-graph-based reasoning to reduce reliance on probabilistic outputs and further enhance interpretability and accuracy.

Dynamic Thresholds: Implement adaptive mechanisms that automatically adjust complexity and criticality thresholds as models and organizational requirements evolve, minimizing manual recalibration.

Explainability Tools: Develop tools to enhance user trust by indicating which knowledge sources were retrieved during RAG and providing transparency about HITL interventions. Such

explainability features can improve compliance confidence and user acceptance.

9. CONCLUSION

The Complexity and Criticality Matrix offers a practical governance framework for deploying generative AI systems in compliance-heavy sectors. By mapping tasks to an appropriate oversight level, the framework transforms high-level governance principles into actionable guidance. It provides a structured approach to balancing automation with regulatory adherence, ensuring each query receives proportionate scrutiny.

Although conceptual at this stage, the matrix sets the groundwork for future research and implementation. Subsequent studies can empirically validate its effectiveness, integrate heuristic or zero-shot classification techniques for task categorization, and develop explainability tools to foster trust and transparency. As enterprises increasingly rely on generative AI, this framework can serve as a cornerstone for responsible, efficient, and compliant AI-driven operations.

REFERENCES

- [1] T. Miller, "Explanation in Artificial Intelligence: Insights from the Social Sciences," *Artif. Intell.*, vol. 267, 2019.
- [2] OpenAI, "GPT-4 Technical Report," OpenAI, 2023.
- [3] Z. Ji, et al., "Survey of Hallucination in Neural Machine Translation," *EMNLP*, 2022.
- [4] E. F. Bender, T. Gebru, et al., "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" *FAccT*, 2021.
- [5] M. Veale, "Governing machine learning that matters: Adaptive regulation, responsive coordination and the risks of regulator entrepreneurship," *Internet Policy Review*, vol. 8, no. 2, 2019.
- [6] R. Calo, "Artificial Intelligence Policy: A Primer and Roadmap," *U.C. Davis Law Review*, vol. 51, 2018.
- [7] T. Sato, "DevOps Patterns and Anti-Patterns," *IEEE Software*, vol. 33, no. 6, 2016.
- [8] T. Klausen, "MLOps: ModelOps, DataOps, and Everything In Between," *InfoQ*, 2022.
- [9] P. Lewis, et al., "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks," *NeurIPS*, 2020.
- [10] K. Starek, et al., "Guiding Humans and Machines toward Better Decisions," *CHI*, 2021.

[11] NIST, "Artificial Intelligence Risk Management Framework (AI RMF)," NIST, 2023.

[12] ISO/IEC 27001:2013, "Information technology — Security techniques — Information security management systems — Requirements," ISO/IEC, 2013.

[13] T. W. Simpson and B. E. Henneman, "Human-in-the-Loop Systems Engineering Framework for Enhancing Human-Machine Teaming in Complex Systems," *Systems*, vol. 10, no. 3, 2022.

[14] US Dept. of Health & Human Services, "HIPAA Privacy Rule," HHS.gov.

[15] OECD, "Recommendation of the Council on Artificial Intelligence," OECD Legal Instruments, 2019.

[16] European Commission, "Proposal for a Regulation on Artificial Intelligence," 2021.

[17] S. Liu, et al., "Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in NLP," *ACM Comput. Surv.*, 2023.

[18] J. Buolamwini and T. Gebru, "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification," *PMLR*, 2018.

[19] M. Mitchell, et al., "Model Cards for Model Reporting," *FAT.*, 2019.

[20] Y. Xia, et al., "Zero-Shot Learning — A Comprehensive Review," *Neural Networks*, vol. 124, 2020.

[21] B. Zhang, et al., "Automated Heuristic-Based Classification in NLP Pipelines," *IEEE Access*, vol. 8, 2020