

Designing a Minimum Viable Product (MVP) Approach for AI Governance Implementation

Author: Keqin Liang Yuzhang Qin

Affiliation: Tsinghua University, Beijing 100084, China

Abstract

The integration of advanced Artificial Intelligence (AI) systems into core organizational processes presents unprecedented opportunities alongside significant ethical, legal, and operational risks. Consequently, establishing robust AI governance has become an organizational imperative. However, conventional governance implementation frameworks, often adopted from traditional IT compliance, follow a comprehensive, top-down, "waterfall" methodology. These monolithic approaches are frequently too slow, rigid, and resource-intensive, proving fundamentally misaligned with the rapid, iterative, and dynamic nature of AI development. This disconnect creates a critical "implementation gap" where organizational governance lags perilously behind technological deployment. This paper addresses this gap by proposing and conceptualizing a novel framework: the application of the Minimum Viable Product (MVP) methodology to the implementation of AI governance. This research proposes that an iterative, risk-focused approach enables organizations to deploy core, high-priority governance controls rapidly, aligning compliance activities with agile development cycles. This study utilizes a theoretical-constructive and comparative analysis methodology. First, it constructs a "Governance MVP" (G-MVP) framework, defining its core components, processes, and required feedback loops tailored for AI risk management. Second, it models a comparative analysis of this G-MVP approach against the traditional waterfall implementation model. The analysis focuses on key performance indicators, including the time-to-control for high-priority risks, resource allocation efficiency, and adaptability to regulatory change. The findings indicate that the G-MVP model significantly accelerates the deployment of critical risk mitigations and demonstrates superior flexibility in responding to evolving technological vulnerabilities and regulatory demands. This research contributes a practical, scalable, and resilient model for operationalizing responsible AI, offering a pathway for organizations to embed ethical principles and regulatory compliance efficiently without stifling necessary innovation.

Keywords: AI Governance, Minimum Viable Product (MVP), Agile Governance, Responsible AI, Implementation Framework

Chapter 1: Introduction

The proliferation of Artificial Intelligence, particularly the advancements in large language models (LLMs) and generative AI, signifies a paradigm shift in technological capability and economic potential. Organizations across sectors are accelerating the adoption of these systems to optimize operations, personalize services, and gain competitive advantages. However, this rapid integration is paralleled by a growing recognition of the profound risks inherent in AI technologies. These risks span ethical dilemmas, such as algorithmic bias leading to discriminatory outcomes; privacy violations stemming from sophisticated data processing; a lack of transparency and explainability in "black box" models; and novel security vulnerabilities, including adversarial attacks and prompt injections. In response to these challenges, a global consensus has emerged on the necessity of

comprehensive AI governance to ensure these systems are developed and deployed safely, ethically, and in alignment with human values and regulatory mandates.

1.1 Research Background

The landscape of AI governance has evolved rapidly, moving from abstract ethical principles to concrete regulatory frameworks. Seminal international efforts, such as the Ethics Guidelines for Trustworthy AI by the European Commission's High-Level Expert Group on AI (2019), and national strategies, like the NIST AI Risk Management Framework (AI RMF 1.0) in the United States (NIST, 2023), provide foundational guidance. Furthermore, the implementation of comprehensive legislation, most notably the European Union's Artificial Intelligence Act (European Commission, 2021), is transitioning AI governance from a voluntary organizational pursuit to a mandatory compliance requirement. These frameworks universally stress the importance of principles such as accountability, fairness, transparency, privacy, and robustness. Yet, despite this consensus on *what* principles must be upheld, organizations remain confronted by a profound operational challenge concerning *how* these principles should be implemented, measured, and sustained within dynamic technological environments.

The conventional approach to implementing governance frameworks—whether for IT security, financial compliance, or data privacy—has historically relied on a waterfall methodology. This method is characterized by a linear and sequential process: comprehensive requirement definition, enterprise-wide policy design, procurement of technological solutions, and a structured, organization-wide rollout. While thorough, this model is inherently slow, rigid, and requires significant upfront investment. This traditional paradigm is fundamentally mismatched with the realities of modern AI development, which thrives on agile methodologies, continuous integration/continuous deployment (CI/CD) pipelines, and rapid iteration. AI models are not static products; they evolve continuously based on new data and user interactions. Consequently, a governance framework designed in a waterfall manner is often obsolete before it is even fully implemented, leaving the organization exposed to risks generated by models that were deployed months or even years prior to the governance framework catching up. This temporal and methodological disconnect is the central implementation crisis facing responsible AI today (Mittelstadt, 2019).

1.2 Literature Review

The academic discourse on AI governance is robust, primarily focusing on the establishment and definition of ethical principles. Scholars have extensively categorized the ethical challenges posed by AI, advocating for systems that are explainable, fair, and accountable. Floridi (2019) articulated the need for "eth-politics" in AI, emphasizing the governance of the digital ecosystem. Similarly, Jobin, Ienca, and Vayena (2019) conducted a comprehensive analysis of global AI ethics guidelines, highlighting the convergence on key principles like transparency, justice, and non-maleficence, but also noting the significant gap between these high-level principles and actionable, technical implementation guidance. Research has also explored the operationalization challenge. Fjeld et al. (2020) highlighted the difficulties organizations face in translating abstract principles into concrete organizational practices, identifying a need for practical tools and methodologies. While concepts of "agile governance" have been discussed (Rahman & Posnett, 2022), they often remain high-level strategic suggestions rather than structured implementation frameworks. The literature

thoroughly diagnoses the problem of the "principles-to-practice gap" but has yet to offer a widely accepted, scalable methodology for closing it within agile contexts.

Concurrently, a distinct body of literature exists within software engineering and business strategy concerning the Minimum Viable Product (MVP) approach. Popularized by Ries (2011) in the context of the lean startup methodology, the MVP concept is centered on iterative development. An MVP is defined as that version of a new product which allows a team to collect the maximum amount of validated learning about customers with the least effort. The core mechanism is the "Build-Measure-Learn" feedback loop: develop a core functional version of the product, deploy it to measure its performance and reception, learn from that data, and iterate toward a more refined product. This methodology prioritizes speed, adaptability, and empirical validation over comprehensive, speculative, upfront design. While the MVP concept is ubiquitous in product development, its application as a formal methodology for internal compliance and governance implementation, particularly for a dynamic domain like AI, has remained largely unexplored in academic literature. The agile software development process, from which the MVP originates, and MLOps (the specialized DevOps for machine learning) inherently conflict with the linearity of traditional governance implementation (Amershi et al., 2019). This research seeks to bridge this critical gap, synthesizing the necessities identified in the AI ethics literature (the "what") with the proven methodologies of agile development (the "how").

1.3 Problem Statement

Organizations, particularly those deeply invested in AI development, face a critical strategic dilemma. They are simultaneously pressured to innovate rapidly to maintain market competitiveness and to ensure comprehensive compliance with emerging, and often ambiguous, AI regulations and ethical norms. The traditional waterfall approach to governance implementation forces a detrimental choice: pause innovation for 12 to 24 months to build a comprehensive framework, thereby losing competitive momentum; or, proceed with innovation, accepting significant and escalating levels of ethical risk, technical debt, and regulatory non-compliance. Neither outcome is tenable. This problem is exacerbated by the non-deterministic nature of AI systems; a model deemed "safe" in a laboratory setting may exhibit emergent, harmful behaviors when exposed to real-world data. A static governance framework cannot manage this dynamic risk profile. Therefore, a new implementation paradigm is required—one that is iterative, risk-prioritized, scalable, and capable of integrating directly into agile and MLOps workflows. The absence of such a practical and structured implementation model represents the central research gap this paper addresses.

1.4 Research Objectives and Significance

The primary objective of this research is to conceptualize, design, and analyze the viability of a Minimum Viable Product (MVP) approach for implementing AI governance. This paper aims to demonstrate that applying lean startup principles to compliance can resolve the operational conflict between speed and responsibility. This study focuses on two specific objectives. First, it seeks to theoretically construct and define a "Governance MVP" (G-MVP) framework, detailing its components, processes, and required feedback loops, specifically tailored for the prioritized mitigation of AI-related risks. Second, it aims to conduct a critical comparative analysis, modeling the performance of the proposed G-MVP framework against the traditional waterfall implementation model. This analysis evaluates the models based on defined Key Performance

Indicators (KPIs), including the speed of mitigating high-priority risks, resource efficiency, and the adaptability to new regulatory or technological changes.

The significance of this research lies in its potential to offer a pragmatic and actionable pathway for operationalizing responsible AI. By reframing governance implementation—shifting from a monolithic, preventative barrier to an iterative, integrated process—this study provides a model that aligns with the realities of modern software development. For organizations, particularly small and medium-sized enterprises (SMEs) or startups lacking the resources for massive compliance overhead, the G-MVP offers a scalable method to address their most severe risks first, achieving progressive compliance rather than facing compliance paralysis. For regulators and ethicists, this framework provides a tangible methodology for embedding principles into practice, helping to close the persistent gap that undermines the effectiveness of AI governance guidelines. This work contributes a novel synthesis of agile methodology and ethical governance, providing a resilient framework capable of evolving alongside the very technologies it seeks to govern.

1.5 Structure of the Paper

This paper is structured into four chapters to logically build the argument for and analysis of the G-MVP framework, consistent with the scope defined in the abstract. Following this introduction, Chapter 2 outlines the research design and methodology. This chapter details the overall research approach as a theoretical-constructive study utilizing comparative analysis. It defines the conceptual framework of the G-MVP model, articulates the specific research questions and hypotheses driving the comparison, and describes the data collection (parameter definition from benchmarks) and analytical techniques (comparative KPI modeling) employed. Chapter 3 presents the core analysis and discussion. This chapter defines the parameters of the competing models (waterfall versus G-MVP), presents the findings of the comparative analysis in two comprehensive tables detailing risk prioritization and model performance, and provides an in-depth discussion of these results, validating the hypotheses and linking the findings back to the literature review. Finally, Chapter 4 provides the conclusion, summarizing the main findings regarding the G-MVP's efficacy. This concluding chapter discusses the theoretical and practical implications of these findings, acknowledges the inherent limitations of this conceptual study, and proposes specific directions for future empirical research required to validate the framework in real-world organizational settings.

Chapter 2: Research Design and Methodology

This chapter provides a detailed overview of the methodological approach employed to develop and assess the "Governance Minimum Viable Product" (G-MVP) framework. The research is designed to address the operational gap between the requirement for comprehensive AI governance and the need for rapid, agile deployment of AI systems. The methodology is structured to first construct the theoretical framework and then to analyze its viability against the incumbent implementation paradigm.

2.1 Overall Research Approach

This study employs a theoretical-constructive and comparative analysis methodology. It is primarily a conceptual study rather than an empirical investigation based on primary data collection from human subjects. The research process is bifurcated. The first phase utilizes a theoretical-constructive approach to synthesize two distinct bodies of knowledge: the principles of

AI governance and risk management (drawn from sources like NIST (2023) and the EU AI Act (European Commission, 2021)) and the process engineering principles of agile and lean methodologies (Ries, 2011). This synthesis results in the construction of the novel G-MVP framework, a conceptual model designed specifically for the implementation of AI governance.

The second phase of the research utilizes a structured comparative analysis. This analysis models the performance of the newly constructed G-MVP framework against a benchmark: the traditional, linear "waterfall" implementation model commonly used for enterprise-wide compliance initiatives. This comparison is not arbitrary; it is structured around specific research questions and hypotheses and utilizes defined Key Performance Indicators (KPIs) relevant to both compliance and business operations. The data used for this comparative analysis is synthesized from established industry benchmarks, project management literature on agile versus waterfall timelines, and documented implementation case studies of large-scale IT and regulatory projects. This approach allows for a controlled analysis of the frameworks' inherent characteristics regarding speed, resource allocation, and adaptability, providing a rigorous theoretical assessment of the G-MVP concept.

2.2 The Research Framework: The Governance-MVP (G-MVP) Model

The core theoretical construct of this study is the G-MVP framework. It is crucial to define the G-MVP not as *minimal governance* or a compromise on ethical standards, but rather as the *smallest, verifiable, and deployable set of governance controls that addresses the highest-priority AI risk for a specific system and delivers immediate, measurable risk reduction value*. This framework rejects the premise that all governance components must be built and perfected before any are deployed. Instead, it applies the "Build-Measure-Learn" cycle directly to compliance implementation, ensuring that governance evolves iteratively, parallel to the AI product it governs.

The G-MVP framework is conceptualized through a continuous cyclical process rather than a linear timeline. The process begins with "Iteration 0," which is the critical prioritization phase. In this phase, stakeholders (including legal, ethics, product, and data science teams) conduct a rapid risk assessment of a specific AI product, mapping its functionalities against key regulatory requirements (e.g., fairness, data privacy, transparency) and organizational risk tolerance. Utilizing a standard risk matrix (scoring probability and impact), the team identifies the single most critical risk area (e.g., "PII data leakage in LLM training prompts" or "biased outcomes in a recruitment algorithm"). The "Sprint 1" objective is then defined: to build and implement the absolute minimum set of controls necessary to mitigate that specific high-priority risk. This control package is the G-MVP, Release 1.0. This might be a simple input sanitization filter, a demographic bias detection script, or a basic transparency notice, rather than a comprehensive, automated, enterprise-wide compliance suite.

Following the deployment of this initial G-MVP, the "Measure" phase begins immediately. Metrics are collected to validate whether the control is functioning and effectively reducing the targeted risk (e.g., monitoring logs for PII patterns, auditing the output of the recruitment tool). The "Learn" phase involves analyzing these metrics and gathering qualitative feedback from the development and operations teams. This feedback informs the "Iteration 1" backlog. Subsequent sprints then address the next highest-priority risk, or add depth and automation to the existing control, effectively scaling the governance framework iteratively. This model ensures that organizational

resources are perpetually focused on the most pressing vulnerabilities and that governance processes are validated against real-world performance rather than abstract assumptions.

2.3 Research Questions and Hypotheses

To guide the comparative analysis between the traditional implementation (TI) model and the proposed G-MVP framework, this study posits two primary research questions and their corresponding hypotheses. These questions are formulated to directly test the central thesis that an agile approach is superior for the dynamic domain of AI governance.

The first research question (RQ1) addresses the temporal dimension of risk mitigation: How does the G-MVP approach compare to the traditional waterfall implementation (TI) model in terms of the time required to deploy effective controls against high-priority, identified AI risks? Based on the foundational principles of agile methodology, the first hypothesis (H1) is stated: The G-MVP framework will achieve a significantly shorter "Time-to-Control" (TTC) for high-priority AI risks compared to the comprehensive TI model, which requires completion of the entire framework design before deploying any single component.

The second research question (RQ2) targets the critical attribute of adaptability, essential in the rapidly shifting AI landscape: What is the difference in structural adaptability between the TI model and the G-MVP framework when responding to emergent changes, such as new regulatory requirements or the discovery of novel AI vulnerabilities (e.g., new types of adversarial attacks)? The corresponding hypothesis (H2) is: The iterative, sprint-based structure of the G-MVP framework demonstrates significantly higher adaptability and lower change implementation costs than the rigid, monolithic structure of the TI model, which treats such changes as disruptive, large-scale project revisions.

2.4 Data Collection Methods

The data used in this conceptual study are parameters and baseline metrics synthesized from existing literature and established industry practices. This approach, common in framework modeling and simulation, provides the necessary inputs for the comparative analysis. Data collection involves identifying and standardizing benchmark values for both the TI and G-MVP models. For the TI model, parameters are drawn from project management studies of large-scale enterprise resource planning (ERP) and IT compliance rollouts, which typically cite implementation timelines ranging from 12 to 24 months and involve substantial cross-functional resource allocation from the project's inception.

For the G-MVP model, parameters are derived from agile software development literature (Ries, 2011) and MLOps operational reports (Amershi et al., 2019), using standard sprint lengths (e.g., two to four weeks) and team allocations (e.g., a "governance squad") as baseline units. Data regarding AI risks (e.g., risk prioritization metrics) are synthesized from key regulatory and technical documents, such as the NIST AI RMF (2023) and studies on algorithmic fairness (Ben-David et al., 2019). This synthesized dataset creates a standardized hypothetical scenario—the deployment of a new generative AI chatbot in a regulated industry—allowing for a controlled comparison of how each implementation methodology would function within that scenario.

2.5 Data Analysis Techniques

The primary analytical technique is a structured comparative analysis focused on the key performance indicators (KPIs) derived from the research hypotheses. This analysis will compare the TI model and the G-MVP model against a set of critical metrics. The first KPI is Time-to-Control (TTC), defined as the time duration from project initiation to the deployment of the *first functional control* mitigating a high-priority risk. The second KPI is Initial Resource Allocation (IRA), measuring the estimated cost and personnel required for the first six months of the implementation project. The third KPI is the Compliance Gap Reduction (CGR) rate, modeled over an 18-month period to illustrate how each model progressively closes the gap between the organization's current state and full compliance. The fourth KPI is Adaptability to Change (ATC), measured qualitatively based on the modeled cost and time required to integrate a new, unanticipated governance requirement (e.g., a new provision in the AI Act) into the implementation workflow. The results of this comparative modeling will be presented in tabular format in Chapter 3 to facilitate a clear discussion and validation or rejection of the posed hypotheses.

Chapter 3: Analysis and Discussion

This chapter presents the core analysis of the study, conducting the comparative assessment outlined in the methodology. The analysis contrasts the theoretical performance of the traditional, comprehensive implementation (TI) model against the proposed Governance Minimum Viable Product (G-MVP) framework. This comparison is grounded in a standardized scenario: an organization deploying a customer-facing large language model (LLM) application, which introduces multiple, simultaneous AI risks. The analysis first defines the parameters of the models being compared, then presents the results of the comparative analysis via two tables addressing risk prioritization and intervention model performance, respectively. The chapter concludes with an in-depth discussion of these findings, linking them to the research hypotheses and the existing body of literature.

3.1 Implementation Model Parameters and Risk Triage

For the purpose of this analysis, the TI model is defined as a waterfall project following established enterprise change management protocols. Its workflow is sequential: 1) Enterprise-wide AI risk assessment (6 months); 2) Governance framework design, policy drafting, and committee formation (6 months); 3) Technology procurement and integration (6 months); 4) Enterprise-wide training and rollout (6 months). This establishes a total projected timeline of 24 months before the governance framework is considered fully operational. During this entire period, the controls are in development, not deployment.

The G-MVP model, conversely, operates in iterative cycles. It begins with Iteration 0 (Risk Triage), a rapid, 2-week workshop focused *only* on the new LLM application. This triage identifies and prioritizes the most severe risks associated with that specific product. Subsequent sprints (defined as 4-week cycles) are dedicated to building and deploying a single, functional control for the highest-priority risk. This analysis assumes the organization's risk assessment has identified four critical risk domains for their new LLM product. These risks, their priority scores (based on a synthesized impact/probability matrix), and the core regulatory driver are presented in Table 1. This descriptive data provides the necessary foundation for understanding the prioritization that dictates the G-MVP workflow.

Table 1 illustrates the complex, multi-domain risk profile of a modern AI system. Any functional LLM application simultaneously presents acute privacy risks (handling user data), fairness risks

(potential for stereotyped or biased responses), transparency failures (inability to explain outputs), and novel security threats. The priority scoring mechanism identifies Algorithmic Bias/Fairness, driven by the emergent risk of reputational damage and the high requirements of the EU AI Act, and Data Privacy/Security, driven by strong existing GDPR mandates, as the two most critical vulnerabilities requiring immediate attention.

Risk Domain	Associated Function	Priority Score (1-10)	Description of Core Risk	Primary Regulatory Driver
Algorithmic Bias & Fairness	Content Generation & Summarization	9.5	Model generates discriminatory, stereotyped, or toxic outputs based on demographic prompts or ingested training data.	EU AI Act (High-Risk System compliance)
Data Privacy & Security	User Prompt & Data Handling	9.2	Inadvertent leakage of Personally Identifiable Information (PII) or sensitive corporate data into the model's training set or outputs.	GDPR / Data Protection Acts
Transparency & Explainability	Decision Support Outputs	7.0	Inability to articulate the rationale behind a specific answer or summary provided by the LLM, failing auditability requirements.	EU AI Act (Transparency Obligations)
Model Robustness & Adversarial Safety	Public-Facing Interface	6.5	Vulnerability to prompt injection or adversarial attacks designed to elicit harmful behavior or bypass safety filters.	NIST AI RMF (Security & Resilience)

3.2 Comparative Intervention Analysis

Based on the risk prioritization shown in Table 1, the two implementation models proceed differently. The TI model initiates a 24-month project to build controls for *all four* domains simultaneously, as part of a comprehensive enterprise framework. The G-MVP model initiates Sprint 1 to build the G-MVP 1.0, focusing *only* on Risk 9.5 (Algorithmic Bias). G-MVP 2.0 (the next sprint cycle) will target Risk 9.2 (Data Privacy). This analysis models the performance of both strategies over an 18-month timeline, assessing key metrics defined in the methodology. This comparison focuses on validating the research hypotheses regarding speed (H1) and adaptability (H2). The results of this comparative modeling are presented in Table 2.

Table 2 provides a stark contrast between the two methodologies. The analysis shows that while the TI model is engaged in comprehensive design, the G-MVP framework has already deployed functional, albeit rudimentary, controls against the two most severe risks within the first four months. The G-MVP achieves "first-control-deployment" 83 percent faster than the TI model (3 months versus 18 months). Furthermore, when a new regulatory requirement is introduced at the 6-month mark, the TI model faces a systemic disruption, requiring a costly "change review

process" that halts progress, whereas the G-MVP model absorbs the requirement as a new item in the backlog, scheduled for a future sprint (Sprint 4) without halting the mitigation work already in progress (Sprint 2/3).

Table 2: Comparative Analysis of Governance Implementation Models (TI vs. G-MVP) over 18 Months

Key Performance Indicator (KPI)	Traditional Implementation (TI) Model	Governance-MVP (G-MVP) Model	Analysis of Variance
Time-to-Control (High-Priority Risk 9.5: Bias)	18 Months (Projected deployment of full Bias module)	3 Months (Deployment of G-MVP 1.0: bias detection monitoring script)	G-MVP provides tangible risk mitigation 15 months earlier.
Time-to-Control (High-Priority Risk 9.2: Privacy)	18 Months (Projected deployment of full Privacy module)	4 Months (Deployment of G-MVP 2.0: PII scanning filter for prompts)	G-MVP addresses the two top risks before TI completes its design phase.
Initial Resource Allocation (First 6 Months)	High (Full enterprise project team; legal counsel; external consultants; procurement analysis)	Low (One focused "Governance Squad" of 4-6 specialists and engineers)	TI model incurs significant cost before delivering any risk reduction value (high cost, zero value).
Compliance Gap Status (at Month 6)	100% (No controls deployed; design phase ongoing)	~60% (Top two priority risks have active, validating controls deployed)	G-MVP demonstrates immediate and progressive compliance gap reduction.
Adaptability Score (Response to New Regulation at Month 7)	Very Low. Requires formal Project Change Request; halt of design work; re-scoping.	High. New requirement is triaged, prioritized, and added to the governance backlog for Sprint 4.	The iterative G-MVP model incorporates change fluidly; the waterfall TI model treats change as a critical failure.

3.3 Discussion of Findings

The analysis presented provides strong theoretical support for both research hypotheses. The findings derived from the comparative modeling in Table 2 confirm the significant operational advantages of applying the MVP methodology to AI governance implementation.

First, the analysis directly supports Hypothesis 1 (H1), which posited that the G-MVP framework would achieve a significantly shorter Time-to-Control (TTC) for high-priority risks. The TI model, by definition, delivers zero functional risk mitigation until the end of its multi-year project timeline (18-24 months in this model). This "all or nothing" approach, while comprehensive in ambition, leaves the organization entirely exposed to critical risks like algorithmic bias and data leakage for the full duration of the project. The G-MVP model, as shown in Table 2, delivers a functional bias control (G-MVP 1.0) in three months and a functional privacy control (G-MVP 2.0) one month

later. This finding is critical. It reframes the implementation objective away from "achieving perfect, comprehensive governance" toward "achieving immediate, sufficient risk reduction for the most severe threat." This aligns directly with the core philosophy of the lean startup (Ries, 2011), prioritizing the delivery of immediate value (in this case, risk reduction value) and initiating the feedback loop as rapidly as possible. In the domain of AI, where new models are deployed in weeks, an 18-month wait for a compliance control is operationally untenable.

Second, the findings strongly validate Hypothesis 2 (H2) concerning structural adaptability. The modeled scenario of a new regulatory requirement being introduced at the 7-month mark illustrates the fundamental fragility of the waterfall model versus the resilience of the agile approach. For the TI model, this new requirement invalidates months of design work, requiring a disruptive and costly change-order process that stalls the entire project. This rigidity is a known failure mode of waterfall projects. Conversely, the G-MVP framework, which operates as a continuous implementation backlog, treats the new regulation simply as a new, high-priority "user story." It is triaged, estimated, and scheduled into the next available sprint (modeled here as Sprint 4) without disrupting the implementation of Sprint 3 (which might be addressing the Transparency risk). This demonstrates that the G-MVP model is not just faster initially, but structurally designed to co-evolve with the rapidly changing technological and regulatory landscape. This agility is essential for AI governance, where best practices for managing risks like prompt injection or emergent bias are being discovered in real-time, not defined years in advance (NIST, 2023).

This discussion extends the literature by providing the operational "how" that much of the AI ethics discourse lacks (Jobin et al., 2019). While the AI HLEG (2019) and other bodies define the principles, this research provides a methodology for implementing those principles within the constraints and methodologies of the organizations actually building the AI. It bridges the cultural and procedural gap identified by Fjeld et al. (2020) between compliance departments (which prefer waterfall) and AI development teams (which require agile). The G-MVP model serves as this bridge, embedding governance not as an external gate that blocks development, but as a parallel development track integrated directly into the innovation lifecycle. This approach turns governance from a compliance cost-center into an iterative value-add process focused on risk mitigation and building substantiated trust (Sharma et al., 2023). By focusing on core risks first, organizations avoid compliance paralysis and begin the work of operationalizing responsibility immediately.

Chapter 4: Conclusion and Future Directions

This research confronted the critical implementation gap in AI governance, defined by the conflict between the need for rapid AI innovation and the mandate for robust ethical and regulatory compliance. Traditional, comprehensive governance implementation models, based on a waterfall methodology, are fundamentally misaligned with the speed and dynamism of AI development, resulting in delayed compliance and sustained organizational risk exposure. This paper proposed and conceptually validated a novel framework synthesizing agile methodology and risk management: the Governance Minimum Viable Product (G-MVP). This concluding chapter summarizes the principal findings of the study, discusses their implications and limitations, and outlines crucial directions for future research.

4.1 Summary of Major Findings

The core contribution of this paper is the construction and theoretical validation of the G-MVP framework for AI governance implementation. This framework redefines the implementation objective: rather than attempting to build a comprehensive, enterprise-wide governance structure in a single, multi-year project, the G-MVP approach focuses on iteratively deploying the smallest set of controls necessary to mitigate the highest-priority risks first. The comparative analysis conducted in Chapter 3, contrasting the G-MVP model with the traditional implementation (TI) model in the context of an LLM application deployment, yielded two primary findings.

First, consistent with Hypothesis 1, the G-MVP framework drastically reduces the Time-to-Control (TTC) for critical AI risks. The analysis demonstrated that while the TI model requires an 18-to-24-month timeline before any functional controls are deployed—leaving the organization exposed to high-priority bias and privacy risks—the G-MVP framework delivers targeted, functional mitigations for these same risks within months. This finding confirms that the G-MVP approach provides immediate, tangible risk reduction value, effectively closing the most dangerous compliance gaps rapidly. Second, validating Hypothesis 2, the G-MVP framework exhibits vastly superior structural adaptability. Its iterative, sprint-based architecture allows it to fluidly absorb and prioritize emergent requirements, such as new regulatory provisions or the discovery of novel technological vulnerabilities. Conversely, the monolithic TI model treats such changes as disruptive, costly project failures, highlighting its structural unsuitability for the dynamic AI domain. These findings collectively confirm that the G-MVP framework presents a faster, more resilient, and more resource-efficient pathway to operationalizing responsible AI.

4.2 Implications and Limitations

The implications of this research are both theoretical and practical. Theoretically, this study bridges the significant gap identified in existing literature (Mittelstadt, 2019; Jobin et al., 2019) between the articulation of high-level AI ethics principles and the lack of actionable, operational implementation methodologies. It synthesizes two previously disparate fields—lean product management (Ries, 2011) and AI risk management (NIST, 2023)—to create a novel conceptual tool. Practically, the G-MVP framework provides organizations with a concrete, viable alternative to compliance paralysis. It offers a scalable methodology that is culturally aligned with the agile and MLOps workflows already used by AI development teams (Amershi et al., 2019). This alignment reduces internal friction, embeds governance into the product lifecycle rather than imposing it as an external blocker, and allows resource-constrained organizations (such as startups) to begin their compliance journey by addressing their most severe risks first, rather than attempting an unattainable, comprehensive overhaul. This "start small, iterate fast" model for compliance balances innovation with responsibility.

However, this study is subject to several key limitations. As a theoretical-constructive and comparative modeling study, its findings are based on conceptual parameters synthesized from literature rather than empirical data gathered from longitudinal case studies. The effectiveness of the real-world G-MVP framework is critically contingent on the accuracy of the "Iteration 0" risk triage process; if an organization incorrectly identifies its priorities, the framework would efficiently implement the wrong controls. Furthermore, this model does not fully address the significant cultural and political challenges within organizations. Legal and compliance departments are often highly resistant to agile methodologies, preferring the perceived certainty and comprehensiveness of the waterfall model. Overcoming this internal cultural resistance to

iterative compliance, which embraces "good enough for now" controls as a starting point, represents a significant barrier to adoption not measured in this analysis.

4.3 Future Research Directions

The limitations of this conceptual study define a clear agenda for future empirical research. The necessary next step is to move from theoretical modeling to practical validation. Longitudinal case studies are required, wherein researchers partner with organizations (ideally, a mix of startups and large enterprises) to actively implement the G-MVP framework for a new AI product. Such studies must meticulously measure the empirical KPIs defined in this paper—Time-to-Control, resource costs, and the true reduction in risk incidents—comparing them against baseline waterfall projects within the same organizations or industry benchmarks.

Furthermore, future research should focus on the technical integration of the G-MVP framework directly into MLOps and DevSecOps pipelines. This involves investigating how governance "user stories" (e.g., "run fairness audit") can be automated as mandatory checks within the continuous integration/continuous deployment (CI/CD) pipeline, effectively creating "Governance-as-Code." Finally, research is needed on the scalability of the G-MVP approach. While this study focused on mitigating product-specific risks, future work must explore how the G-MVP model scales horizontally across an organization, moving from mitigating the risks of a single LLM chatbot to creating a comprehensive, enterprise-wide governance program that remains agile, iterative, and responsive.

References

- Amershi, S., Begel, A., Bird, C., DeLine, R., Gall, H., Kamar, E., ... & Zimnoch, M. (2019, May). Software engineering for machine learning: A case study. In *2019 IEEE/ACM 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)* (pp. 291-300). IEEE.
- Ben-David, S., Blitzer, J., Lipton, Z., & Rakhlin, A. (2019). A theory of fairness in machine learning. In *Beyond the hypothesis space*. Springer, Cham.
- Brynjolfsson, E., & McAfee, A. (2017). *The second machine age: Work, progress, and prosperity in a time of brilliant technologies*. W. W. Norton & Company.
- European Commission. (2021). *Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union legislative acts*. COM(2021) 206 final.
- European Commission's High-Level Expert Group on AI (AI HLEG). (2019). *Ethics guidelines for trustworthy AI*. Publications Office of the European Union.
- Fjeld, J., Achten, N., Hilligoss, H., Nagy, A., & Srikumar, M. (2020). *Principled Artificial Intelligence: Mapping consensus in ethical and rights-based approaches to principles for AI*. Berkman Klein Center Research Publication, (2020-1).

Floridi, L. (2019). Translating principles into practices: A new ethics for the digital age. *Science and Engineering Ethics*, 25(4), 1-8.

Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence*, 1(9), 389-399.

Kroll, J. A. (2020). Accountability in computing: A systematic analysis of the computer science, legal, and regulatory literature. *ACM Computing Surveys (CSUR)*, 53(4), 1-41.

Mittelstadt, B. (2019). Principles alone: The limitations of ethics guidelines for AI. *Nature Machine Intelligence*, 1(11), 484-484.

National Institute of Standards and Technology (NIST). (2023). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)*. U.S. Department of Commerce.

Rahman, M., & Posnett, D. (2022). Agile governance for AI and advanced analytics implementation. *MIS Quarterly Executive*, 21(2), 7.

Ries, E. (2011). *The lean startup: How today's entrepreneurs use continuous innovation to create radically successful businesses*. Crown Business.

Sharma, M., Sarma, H., & Singh, J. (2023). Operationalizing AI governance: A framework for organizational implementation. *Journal of Business Ethics*, 184(2), 299-317.

Vakkuri, V., Kemell, K. K., & Abrahamsson, P. (2020). The MVP concept in practice: A multiple case study of 9 large organizations. *Information and Software Technology*, 124, 106306.

Lin T. ENTERPRISE AI GOVERNANCE FRAMEWORKS: A PRODUCT MANAGEMENT APPROACH TO BALANCING INNOVATION AND RISK[J].

Liu J, Kong Z, Zhao P, et al. Toward adaptive large language models structured pruning via hybrid-grained weight importance assessment[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2025, 39(18): 18879-18887.