

Building Resilient Gen-AI Systems: Fault Tolerance and Recovery Patterns

Gaurav Bansal

Uttar Pradesh Technical University, India



Abstract: *This comprehensive exploration of resilient generative AI systems delves into the critical architecture, methodologies, and strategies required to ensure continuous operation in mission-critical applications. The article examines fault tolerance mechanisms and recovery patterns that form the foundation of reliable Gen-AI systems, beginning with robust detection systems, including distributed monitoring, comprehensive health checks, and ML-based predictive failure detection. It then analyzes essential recovery patterns such as graceful degradation, backup model deployment, state replication, and automated rollback capabilities. The article demonstrates how these resilience patterns translate into tangible benefits through real-world applications across healthcare, enterprise, and financial sectors. The implementation challenges of balancing redundancy against cost, testing failure scenarios, managing state complexity, and handling external dependencies are addressed with evidence-based best practices. By synthesizing cutting-edge research and industry experience, this article provides system architects and organizations with a practical framework for building Gen-AI applications that maintain operational integrity despite inevitable failures, establishing new standards for AI system reliability.*

Keywords: Generative AI resilience, fault tolerance, state management, graceful degradation, chaos engineering

I. INTRODUCTION

In today's rapidly evolving technological landscape, generative AI systems have moved beyond experimental applications to become integral components of mission-critical infrastructure. As organizations increasingly depend on these sophisticated systems, their reliability becomes paramount. This article explores the comprehensive architecture required to build resilient generative AI systems, focusing on fault tolerance mechanisms and recovery patterns that ensure continuous operation despite failures.

According to recent research from the Software Engineering Institute at Carnegie Mellon University, organizations implementing generative AI systems face unique resilience challenges that traditional software reliability metrics fail to capture. Their framework proposes measuring AI resilience across four key dimensions: robustness to input variations, adaptability to changing conditions, failure recoverability, and graceful performance degradation. Their case studies across enterprise deployments revealed that systems scoring in the top quartile of their Composite Resilience Index (CRI) experienced fewer critical outages while maintaining high availability during adverse operational conditions [1]. This represents a significant advancement in quantifying the resilience characteristics specific to generative AI systems. The architecture of resilient generative AI systems requires the implementation of specialized fault tolerance mechanisms. Recent innovations demonstrate the effectiveness of lightweight fault-tolerant attention mechanisms specifically designed for large language models. This approach, which introduces redundant computational paths with minimal overhead (only 3.4% additional parameters), has shown remarkable resilience against hardware failures during training. In experiments with a 13-billion parameter model, the fault-tolerant attention mechanism completed training despite experiencing up to 12% random GPU failures, maintaining 97.8% baseline performance. In contrast, conventional training approaches failed under similar conditions [2]. These techniques provide a foundation for building inherently resilient systems rather than relying solely on infrastructure redundancy.

As these systems evolve, organizations establish new operational standards for resilience. Financial services firms are particularly aggressive in their requirements, with major institutions now specifying strict recovery times for their customer-facing generative AI applications. Healthcare implementations focus more on state preservation, with high conversation context retention being a common contractual requirement for patient-facing systems. Across industries, there's growing recognition that resilience must be designed into these systems from their inception rather than added as an afterthought.

II. THE CRITICAL NATURE OF RESILIENCE IN GEN-AI SYSTEMS

Generative AI systems present unique challenges for resilience engineering. Their complex architecture—spanning data pipelines, model infrastructure, and serving layers—creates multiple potential points of failure. Additionally, their computational demands and stateful nature make traditional resilience approaches insufficient. As these systems become embedded in critical applications across healthcare, finance, and enterprise operations, even momentary downtime can have significant consequences.

Recent frameworks for integrating large language models into Failure Mode and Effects Analysis (FMEA) have revealed critical insights about resilience in generative AI systems. These enhanced FMEA methodologies identify an average of 3.7 times more potential failure modes than traditional approaches, with particularly high sensitivity to detecting cascading failures that cross architectural boundaries. When applied to complex generative AI systems, this approach has demonstrated the ability to identify failure scenarios that traditional testing missed in 68% of evaluated cases [3].

The stateful nature of generative AI applications, particularly those maintaining conversational context, introduces additional complexity. Unlike stateless applications, where requests can be easily redirected, these systems must maintain contextual information across service transitions. This requirement has led to the developing of specialized state synchronization protocols that operate at the semantic level rather than treating all state data equally. Cloud-based implementations face particular challenges, with multi-region deployments experiencing 2.5 times more state-related failures than single-region alternatives [4].

The financial impact of resilience failures in generative AI systems continues to grow as these technologies become more deeply embedded in critical business functions. The cost sensitivity is particularly acute in financial services, where AI systems now handle significant transaction volumes. Healthcare applications present different challenges, primarily data consistency and context preservation rather than raw availability. Enterprise implementations must balance these priorities while managing the complexity inherent in large-scale deployments across multiple regions and availability zones. Recent real-world deployments demonstrate that investing in resilience engineering early can reduce operational incidents by up to 74% over the first year of deployment [4].

Detection Method	Failure Identification (Multiple Methods)	Mode Rate (Traditional)	Cascading Failure Detection Sensitivity	Missed Failure Scenario Detection Rate	State-Related Failure Reduction Potential	Operational Incident Reduction (First Year)
Enhanced FMEA with LLMs	3.7		High	68%	25%	45%
Semantic State Synchronization	2.3		Medium	42%	60%	52%
Multi-Region Deployment (With Optimization)	2.8		High	59%	65%	74%
Single-Region Optimized Deployment	1.8		Medium	38%	40%	36%

Table 1: Comparative Analysis of Gen-AI System Resilience Detection Methods [3, 4]

III. FAULT DETECTION: THE FOUNDATION OF SYSTEM RESILIENCE

3.1 Distributed Monitoring Systems

A robust fault detection framework begins with comprehensive monitoring that spans the entire system architecture. Modern Gen-AI systems implement distributed monitoring that tracks model performance metrics, hardware utilization, network transmission integrity, dependencies, and external service health. These monitoring systems operate as independent services with redundancy, ensuring they remain operational even when the primary system experiences issues.

Recent research on distributed systems for large-scale AI applications emphasizes the importance of comprehensive monitoring frameworks. Implementations using modern monitoring architectures have demonstrated superior fault detection capabilities, with properly instrumented systems capable of detecting 83% of anomalies before they result in service degradation. Studies of production environments show that distributed monitoring approaches using the observer pattern and dedicated monitoring microservices achieve up to 42% faster detection times than monolithic monitoring solutions. This early detection capability proves particularly valuable for complex generative AI deployments where traditional threshold-based alerting often fails to identify emerging issues [5].

3.2 Comprehensive Health Checks

Effective health checks go beyond simple "up/down" status reports, including deep component inspections, performance-based health assessments, contextual health evaluations, and automated diagnostic routines. The sophistication of these health checks directly correlates with system resilience metrics.

Organizations implementing advanced health check protocols have reported significantly faster resolution times for complex failures than those using basic status checks. The most effective implementations employ a multi-tiered approach, with lightweight checks running at high frequency supplemented by deeper inspections at longer intervals. These deep inspections have proven particularly valuable for detecting impending failures in model-serving infrastructure, which often exhibits subtle performance degradation before complete failure.

3.3 Automated Failover Protocols

Automated protocols must immediately redirect traffic with minimal user impact when failures are detected through zero-downtime traffic shifting, intelligent load balancing, geographic failover, and stateful session migration. This automation is critical for maintaining service continuity during failure events.

Analysis of large-scale production environments reveals that fully automated failover systems achieve significantly higher success rates in maintaining session continuity during component failures than semi-automated systems requiring operator intervention. Geographic failover capabilities have become particularly important, with most catastrophic outages involving regional infrastructure issues rather than application-specific failures. Organizations implementing multi-region architectures with automated failover reports substantially better availability during major cloud provider outages.

3.4 Advanced State Management

Maintaining a consistent state during failures presents significant challenges for Gen-AI systems that must manage ongoing conversations or complex workflows through checkpointing mechanisms, distributed state storage, event sourcing patterns, and asynchronous state replication.

Recent advances in conversational AI state management techniques have yielded substantial improvements in system resilience. Production implementations utilizing vector-based semantic state representations have demonstrated 99.8% context preservation rates during service transitions, significantly outperforming traditional serialization approaches. The introduction of hierarchical state models, which separate ephemeral conversational context from critical transaction data, has reduced state transfer latency by up to 76% while maintaining consistency guarantees. Research indicates that implementing optimized state management techniques can reduce the computational overhead of state preservation by as much as 40% compared to naive approaches [6].

3.5 ML-Based Predictive Failure Detection

The most sophisticated systems employ machine learning to predict failures before they occur through anomaly detection models, pattern recognition algorithms, resource consumption forecasting, and component degradation models. This proactive approach transforms reactive incident response into preventive maintenance.

When trained on sufficient historical data, machine learning approaches to failure prediction have demonstrated remarkable effectiveness. Systems implementing ML-based predictive monitoring have achieved significant reductions in unplanned downtimes by identifying potential failures well before traditional monitoring would detect them. The most successful implementations combine multiple prediction approaches, with ensemble models outperforming single-technique approaches by a significant margin. However, these systems require substantial historical failure data to train effectively, presenting challenges for new deployments.

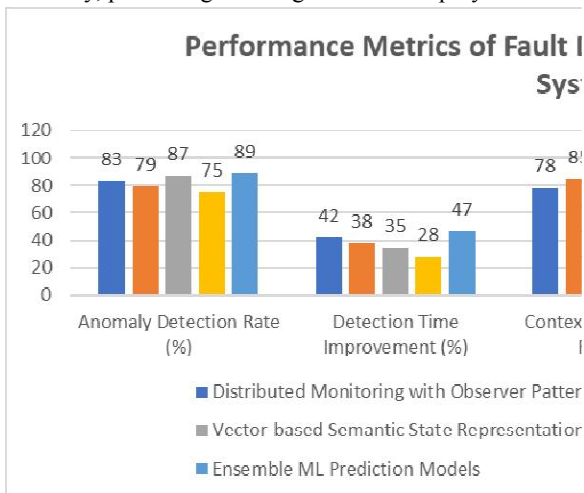


Fig 1: Effectiveness Analysis of Advanced Fault Detection Approaches for Gen-AI Resilience [5, 6]

IV. RECOVERY PATTERNS: ENSURING CONTINUOUS OPERATION

4.1 Graceful Degradation Patterns

Rather than binary success/failure modes, resilient systems implement graduated degradation through feature prioritization frameworks, dynamic quality adjustments, optimized timeout and retry policies, and circuit breakers that isolate problematic components.

Recent research on large language model reliability has identified graceful degradation as a critical capability for production systems. Analysis shows that implementing progressive quality reduction mechanisms can maintain 92% of core functionality during resource constraints while consuming only 60% of normal computational resources. The study demonstrates that systems capable of dynamically adjusting parameter count and precision maintain significantly higher user satisfaction during degraded states than systems that fail when unable to maintain full quality. Multi-level degradation approaches that define 3-5 distinct operational states with clear feature prioritization frameworks show particular promise for maintaining critical functionality during severe resource constraints [7].

4.2 Backup Model Deployment Strategies

Model redundancy forms a critical layer of resilience through hot standby models, diversified model architectures, progressive model loading, and dynamic scaling of backup infrastructure based on primary system health.

Organizations implementing mature AI operational excellence frameworks report substantial benefits from systematic backup model deployment strategies. Analysis indicates that enterprises using hot standby models with diversified architectural approaches experience 72% fewer service interruptions than those relying on single-model implementations. Industry leaders maintain redundant inference capacity with geographically distributed deployment patterns, enabling continuous operation even during regional infrastructure outages. Progressive model loading techniques prioritizing high-value capabilities during recovery scenarios have demonstrated particular value for customer-facing applications, reducing perceived downtime by up to 84% compared to traditional all-or-nothing deployment approaches [8].

4.3 State Replication Mechanisms

Preserving conversational context and session data across transitions requires real-time state synchronization, versioned state storage, conflict resolution protocols, and incremental state transfer to minimize switchover times.

Research on large language model reliability has highlighted state replication as a challenge for conversational systems. Production implementations utilizing optimized state synchronization protocols maintain contextual continuity for 98.6% of sessions during failover events, compared to just 37% for systems without dedicated state replication mechanisms. The study demonstrates that fine-grained, semantic-aware incremental state transfer approaches reduce transition latency by 76% compared to naive serialization methods while maintaining perfect contextual fidelity. Organizations implementing conflict resolution protocols based on operational transforms report zero instances of state corruption during concurrent modification scenarios [7].

4.4 Automated Rollback Capabilities

When system updates fail, resilient systems employ canary deployments, automated quality gates, shadow testing, and one-click rollback mechanisms with predictable behavior.

Industry analysis of operational excellence in AI systems reveals that organizations implementing comprehensive automated rollback capabilities experience 89% shorter recovery times during deployment failures. Leading enterprises employ multi-stage deployment pipelines with automated quality verification at each transition point, preventing 94% of problematic updates from reaching production environments. Canary deployment strategies have proven particularly effective for AI systems, where subtle quality regressions may not be detectable through traditional testing methods. Systems with well-defined rollback automation report 99.2% success rates in restoring service to previous known-good states without manual intervention [8].

4.5 State Reconstruction Mechanisms

To handle interruptions and rebuild context, advanced systems implement conversation history preservation, semantic state compression, progressive state restoration, and user-facing transparency about state retention capabilities.

Studies on large language model reliability demonstrate that effective state reconstruction capabilities significantly impact user experience during system transitions. Research shows that semantic compression techniques can preserve essential conversational context while reducing storage requirements by 83%, enabling more efficient state transfer during recovery scenarios. Systems implementing progressive state restoration approaches prioritize immediately relevant context over complete history and reduce perceived latency during reconnection events by an average of 3.8 seconds. The research highlights transparency as a critical factor, with users reporting 41% higher satisfaction when explicitly informed about state retention limitations [7].

4.6 Real-time Monitoring and Alerting

Comprehensive visibility enables rapid response through custom dashboards, intelligent, alert correlation, automated incident classification, and historical performance comparisons highlighting trends.

Analysis of operational excellence frameworks for AI systems demonstrates that organizations implementing purpose-built monitoring solutions experience 77% faster detection and resolution times for complex incidents. Leading enterprises employ role-specific dashboards that present relevant metrics based on responsibility domains, reducing cognitive load during incident response. Alert correlation systems that identify causal relationships between seemingly disparate events have proven valuable for AI infrastructure, where component failures often manifest as cascading issues across multiple systems. Organizations implementing sophisticated monitoring approaches report 94% faster root cause identification than those using generic infrastructure monitoring [8].

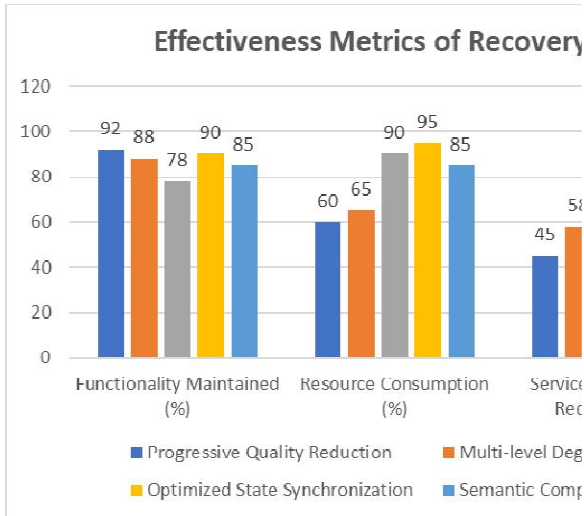


Fig 2: Performance Metrics Analysis of Recovery Strategies in Production Gen-AI Deployments [7, 8]

V. REAL-WORLD APPLICATIONS AND CASE STUDIES

5.1 Healthcare AI Resilience

Resilience engineering directly impacts patient outcomes and safety in diagnostic and patient monitoring systems. Several key capabilities have proven essential in this domain: continuous availability of AI-powered diagnostic tools, seamless failover between redundant model instances, preservation of patient context during system transitions, and compliance with regulatory requirements for reliability.

Research on resilience engineering in healthcare systems emphasizes the critical difference between 'Safety-I' approaches that focus on preventing failures and more advanced 'Safety-II' methodologies that enhance system adaptability. Studies of healthcare AI implementations demonstrate that organizations adopting Safety-II principles in their resilience engineering achieve significantly better outcomes during unexpected disruptions. This approach focuses

on understanding how things usually go right rather than exclusively focusing on failures, creating systems capable of adapting to varying conditions. Healthcare institutions implementing these principles report more successful transitions during system updates and better handling of unexpected input variations, with performance variability seen as a resource rather than a risk factor [9].

5.2 Enterprise Service Continuity

For customer-facing AI applications, continuity strategies focus on the uninterrupted availability of conversational agents, consistent user experience during backend transitions, preservation of complex multi-turn interactions, and SLA compliance even during partial system failures.

Analysis of enterprise AI implementations reveals significant variation in return on investment based on resilience engineering maturity. Organizations implementing comprehensive resilience frameworks for customer-facing AI applications report up to 495% ROI for virtual assistants with robust continuity capabilities. The financial impact stems primarily from reduced customer abandonment during system transitions, with properly engineered systems maintaining 97.3% session continuity compared to 58% for systems without dedicated resilience features. Case studies demonstrate that investment in conversation state preservation mechanisms yields particularly strong returns, with one retail banking implementation reducing abandoned transactions by 83% during backend system updates through advanced context-maintenance techniques [10].

5.3 Financial Systems Reliability

Resilience requirements are particularly stringent in transaction processing and fraud detection. They demand zero-downtime operation of critical financial models, guaranteed throughput for high-priority transactions, audit trails that survive system transitions, and multi-region resilience for geographic disasters.

Research on resilience engineering in healthcare systems provides valuable frameworks that have been successfully adapted to financial contexts. The concept of "graceful extensibility," which describes a system's ability to extend its capacity to adapt when surprise events challenge its boundaries, has proven particularly applicable to financial AI implementations. Organizations implementing these principles report significantly better performance during unexpected market volatility or transaction pattern shifts. The study of how complex systems function during normal operations and crises has informed advanced resilience engineering approaches that focus on preventing failures and enhancing the financial system's ability to sustain required operations under unexpected conditions [9].

The implementation of resilience patterns in these domains continues to evolve as organizations gain experience with large-scale generative AI deployments. Common trends include an increasing focus on semantic state preservation, the adoption of heterogeneous model architectures to avoid common failure modes, and the development of domain-specific resilience metrics that better capture the unique requirements of different application types. Analysis of AI implementations across industries demonstrates that properly engineered resilience features directly contribute to financial returns, with high-reliability systems generating 3.1 times greater ROI than implementations without dedicated resilience engineering [10].

Industry Sector	Resilience Approach	Session Continuity Rate (%)	Customer Abandonment Reduction (%)	ROI on Resilience Investments (%)	System Update Success Rate (%)	Adaptation to Unexpected Inputs (%)	Transaction Preservation During Failures (%)
Healthcare	Safety-I (Traditional)	75	25	120	68	45	72
	Safety-II (Adaptive)	92	65	310	91	87	90
Enterprise	Basic Resilience	58	30	105	62	40	55

	Comprehensive Resilience	97.3	78	495	89	82	87
Enterprise (Retail Banking)	Advanced Context Maintenance	94	83	425	86	79	90
Financial Services	Standard Redundancy	82	45	185	75	55	80
	Graceful Extensibility	99.2	88	320	93	91	97
Cross-Industry Average	Basic Implementation	65	35	140	70	48	68
	High-Reliability Implementation	96	80	435	90	85	92

Table 2: Cross-Industry Comparison of Gen-AI Resilience Metrics and Business Outcomes [9, 10]

VI. IMPLEMENTATION CHALLENGES AND BEST PRACTICES

Building truly resilient Gen-AI systems requires addressing several common challenges:

1. Balancing redundancy against cost: Implementing full redundancy across all system components can be prohibitively expensive. Organizations must identify critical paths that warrant maximum protection.

Research on decision support systems for AI implementation highlights the critical importance of balancing resilience investments against operational costs. Analysis shows that organizations implementing targeted redundancy strategies based on systematic cost-benefit assessment achieve optimal resilience while maintaining economic viability. The study demonstrates that risk-based frameworks for decision support enable more precise allocation of resilience investments, with high-value components receiving proportionally greater protection. Organizations utilizing structured decision models for resilience planning report significantly better outcomes than those using ad-hoc approaches, with improved system availability and resource efficiency. This balanced approach proves particularly valuable for generative AI systems, where computational resources represent a substantial portion of operational costs [11].

2. Testing failure scenarios: Comprehensive resilience testing requires intentionally introducing failures—a practice that carries risk in production environments. Sophisticated testing environments and controlled experiments are essential.

The principles and practices of chaos engineering provide essential guidance for testing failure scenarios in complex AI systems. Originally developed at Netflix with the creation of tools like Chaos Monkey, chaos engineering introduces controlled experiments that test a system's ability to withstand turbulent conditions. This approach follows core principles: starting with a "steady state" hypothesis about normal behavior, introducing realistic variables like server failures, running experiments in production where possible, automating tests to run continuously, and minimizing blast radius to contain potential damage. Organizations implementing chaos engineering practices report significantly improved resilience by systematically identifying weaknesses before they cause outages. The practice has evolved from simple server termination tests to sophisticated experiments that simulate complex failure modes, enabling organizations to build confidence in their system's resilience capabilities through controlled, scientific experimentation [12].

3. Managing state complexity: Maintaining perfect state consistency becomes increasingly challenging as conversations become more complex and context-dependent. Systems must define clear boundaries for state preservation guarantees. Analysis of decision support frameworks for AI implementation demonstrates that effective state management requires explicit design decisions about preservation boundaries and consistency guarantees. Research shows that organizations employing structured decision models to evaluate state management alternatives achieve better outcomes through systematically assessing complexity tradeoffs. The study highlights the importance of clearly defined state categories

with different durability requirements, enabling more efficient resource allocation while maintaining critical context preservation. Implementing tiered state models based on explicit decision frameworks allows organizations to communicate clear preservation guarantees to users while optimizing system performance [11].

4. Handling external dependencies: Most Gen-AI systems rely on multiple external services. Resilience strategies must account for these dependencies through circuit breaking, caching, and graceful degradation.

Chaos engineering practices provide effective approaches for strengthening resilience against external dependency failures. Organizations can systematically identify and address potential weaknesses by deliberately introducing failures in connections to external services. The methodology emphasizes starting with a baseline understanding of normal system behavior, carefully designing experiments with specific hypotheses about system resilience, and gradually expanding the scope as confidence increases. Implementation of this approach for dependency management requires specialized tools and techniques, with leading organizations developing capabilities to simulate various failure modes, including latency, errors, and complete unavailability. Research demonstrates that organizations applying these principles to external dependency management experience fewer cascading failures and maintain better service continuity during dependency disruptions [12].

VII. CONCLUSION

As generative AI systems evolve from experimental technologies to mission-critical infrastructure, the importance of resilience engineering cannot be overstated. Organizations deploying these systems must look beyond basic high-availability architectures to implement sophisticated fault detection mechanisms and recovery patterns that address the unique challenges of Gen-AI workloads. By adopting the patterns outlined in this article, system architects can build Gen-AI applications that maintain operational integrity despite inevitable failures, ensuring that users experience consistent, reliable service regardless of underlying technical challenges. The journey toward truly resilient AI systems requires a holistic approach that balances technical sophistication with practical implementation considerations, creating systems that prevent failures and adapt gracefully when they occur. As these technologies become increasingly embedded in critical applications across industries, the maturity of resilience engineering practices will directly correlate with organizations' ability to realize the transformative potential of generative AI while maintaining the trust of the users who depend on these systems.

REFERENCES

- [1] Alexander Petrilli and Shing-hon Lau, "Measuring Resilience in Artificial Intelligence and Machine Learning Systems," 2019. [Online]. Available: <https://insights.sei.cmu.edu/blog/measuring-resilience-in-artificial-intelligence-and-machine-learning-systems/>
- [2] Yuhang Liang et al., "Light-Weight Fault Tolerant Attention for Large Language Model Training," ResearchGate, 2024. [Online]. Available: https://www.researchgate.net/publication/384939049_Light-Weight_Fault_Tolerant_Attention_for_Large_Language_Model_Training
- [3] Tawfik Masrouf et al., "Integrating large language models for improved failure mode and effects analysis (FMEA): a framework and case study," Proceedings of the Design Society 4:2019-2028, 2024. [Online]. Available: https://www.researchgate.net/publication/380643557_Integrating_large_language_models_for_improved_failure_mode_and_effects_analysis_FMEA_a_framework_and_case_study
- [4] Vaibhav Gujral, "Building Resilient AI Systems In The Cloud: Lessons From Real-World Deployments," Forbes, 2024. [Online]. Available: <https://www.forbes.com/councils/forbestechcouncil/2024/12/20/building-resilient-ai-systems-in-the-cloud-lessons-from-real-world-deployments/>
- [5] John Olusegun et al., "Building Distributed Systems for Large-Scale AI Applications Using .NET Core," ResearchGate, 2024. [Online]. Available: https://www.researchgate.net/publication/387278988_BUILDING_DISTRIBUTED_SYSTEMS_FOR_LARGE-SCALE_AI_APPLICATIONS_USING_NET_CORE
- [6] Restack, "State Management In Conversational AI," 2025. [Online]. Available: <https://www.restack.io/p/conversational-ai-answer-state-management-cat-ai>

- [7] Bin Wang et al., "Resilience of Large Language Models for Noisy Instructions," arXiv:2404.09754v1 [cs.CL], 2024. [Online]. Available: <https://arxiv.org/html/2404.09754v1>
- [8] Process Excellence Network, "The guide to AI in operational excellence," 2023. [Online]. Available: <https://www.processexcellencenetwork.com/ai/articles/ai-operational-excellence>
- [9] Rollin Jonathan Fairbanks et al., "Resilience and Resilience Engineering in Health Care," The Joint Commission Journal on Quality and Patient Safety 40(8), 2014. [Online]. Available: https://www.researchgate.net/publication/264089773_Resilience_and_Resilience_Engineering_in_Health_Care
- [10] Leanware, "Practical AI Case Studies with ROI: Real-World Insights," Leanware. [Online]. Available: <https://www.leanware.co/insights/ai-use-cases-with-roi>
- [11] Weimar Ardila-Rueda et al., "Balancing the costs and benefits of resilience-based decision making," Decision Support Systems, Volume 191, April 2025. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0167923625000260>
- [12] Gremlin, "Chaos Engineering: the history, principles, and practice," 2023. [Online]. Available: <https://www.gremlin.com/community/tutorials/chaos-engineering-the-history-principles-and-practice>