

# **Architectural Synthesis of Agentic Swarm Intelligence and Recursive Meta-Systems for Next-Generation Autonomous DevOps Platforms**

The landscape of software engineering and operational management is undergoing a fundamental transformation, moving from deterministic automation frameworks toward autonomous, goal-oriented cognitive systems. This shift is characterized by the emergence of "Agent Swarm" architectures, exemplified by Moonshot AI's Kimi K2.5, and the recursive self-improvement methodologies pioneered by systems like the POETIQ meta-layer. For a modern DevOps platform to claim the title of the finest in the industry, it must move beyond simple command execution to embrace a unified architecture that orchestrates massive parallel intelligence while maintaining rigorous self-auditing and iterative refinement. The following analysis explores the intricate mechanisms of these technologies and outlines a blueprint for their integration into a state-of-the-art DevOps ecosystem.

## **Foundations of the Kimi K2.5 Architecture: The Mixture-of-Experts Frontier**

The release of Kimi K2.5 in January 2026 marked a pivotal moment in the accessibility of frontier-level agentic capabilities. Unlike previous generations of models that relied on dense parameter sets, Kimi K2.5 utilizes a massive Mixture-of-Experts (MoE) architecture with one trillion total parameters.<sup>1</sup> The core innovation lies in its sparse activation strategy, which ensures that only 32 billion parameters are activated for any single request. This high efficiency allows the model to be deployed in environments where computational resources must be balanced against the need for deep reasoning and vast knowledge capacity.

The architectural backbone consists of 61 layers, comprising one dense layer and 60 MoE layers. Within these layers, the model employs a configuration of 384 distinct experts, with a router selecting only eight experts per token.<sup>1</sup> This selection process is not merely a reduction in compute but a sophisticated delegation mechanism where specific experts specialize in distinct domains, such as code syntax, mathematical logic, or natural language nuances.<sup>1</sup> Complementing this is the Multi-head Latent Attention (MLA) mechanism, featuring a 7,168-dimension hidden attention space, which enables the model to handle a 256,000-token context window with high precision.<sup>2</sup> Such a context window is critical for DevOps applications, where agents must ingest entire codebases, multi-volume documentation, and extensive system logs simultaneously.

Kimi K2.5 is natively multimodal, trained on approximately 15 trillion mixed visual and text tokens.<sup>3</sup> This design choice eliminates the latency and information loss associated with "stitched" multimodal systems, where a separate vision encoder is grafted onto a language model. In a DevOps context, this allows the model to perform "Coding with Vision," where it can interpret UI designs, architectural diagrams, or video walkthroughs of software workflows and translate them directly into executable code.<sup>4</sup>

## Strategic Operation through Multi-Mode Decoding

The Kimi K2.5 system is designed to operate through four distinct modes, each adjusting its decoding strategies and tool permissions to match the specific demands of a task.<sup>1</sup> For a comprehensive DevOps platform, understanding the technical nuances of these modes is essential for optimizing cost and performance.

Operation Mode	Technical Characteristics	Optimization Objective	Primary DevOps Use Case
<b>Instant</b>	Temp: 0.6, Top_p: 0.95, 3-8s Latency	Speed and Token Efficiency	Routine CLI lookups and unit test generation
<b>Thinking</b>	Temp: 1.0, Top_p: 0.95, up to 96K Reasoning Tokens	Depth of Logic and Accuracy	Complex debugging and architectural refactoring
<b>Agent</b>	200-300 Stable Sequential Tool Calls	Autonomous Multi-step Workflow	End-to-end research or documentation synthesis
<b>Agent Swarm</b>	Up to 100 Parallel Agents, 1500 Tool Calls	Massive Parallelism and Throughput	Large-scale security audits or multi-cloud provisioning

In Thinking mode, the model utilizes internal reasoning traces, often referred to as a "thinking budget," to solve problems that require high-fidelity logic.<sup>1</sup> This mode has achieved state-of-the-art results on benchmarks such as AIME 2025 and GPQA-Diamond, demonstrating its capacity for PhD-level reasoning.<sup>1</sup> However, the most transformative capability for a next-generation platform is the Agent Swarm mode, which transitions from

single-agent scaling to a self-directed, coordinated swarm-like execution.<sup>2</sup>

## The Mechanism of Agent Swarm Methodology

The Agent Swarm methodology represents a departure from traditional orchestration, where roles and workflows are predefined by human engineers. In the Kimi K2.5 paradigm, the model acts as a self-directed orchestrator that can autonomously decompose a high-level goal into parallelizable subtasks.<sup>7</sup> This process involves the dynamic instantiation of up to 100 sub-agents, each specialized to handle a specific niche of the broader problem.<sup>3</sup>

### Parallel-Agent Reinforcement Learning (PARL)

The technical foundation of this capability is Parallel-Agent Reinforcement Learning (PARL), a novel training framework that makes parallelism a learnable skill rather than a hard-coded heuristic.<sup>5</sup> PARL addresses the two primary failure modes of naive multi-agent systems: "serial collapse" and "fake parallelism." Serial collapse occurs when an orchestrator, despite having parallel resources, defaults to a sequential processing pattern. Fake parallelism involves the instantiation of multiple agents that do not actually contribute to a reduction in end-to-end runtime due to high coordination overhead.<sup>7</sup>

To overcome these obstacles, PARL utilizes a reward function that is latency-aware, evaluating performance based on "Critical Steps" rather than total tokens or steps.<sup>5</sup> The Critical Steps metric identifies the longest path in the execution graph, and the model is incentivized to shorten this path. The reward structure can be mathematically modeled as follows:

$$R_t = \lambda_{aux}(e) \cdot r_{parallel} + (1 - \lambda_{aux}(e)) \cdot (I[\text{success}] \cdot Q(\tau))$$

In this equation, the reward  $R_t$  is a balanced function of the parallelization efficiency ( $r_{parallel}$ ) and the ultimate quality of the task success ( $Q(\tau)$ ), weighted by a scheduling factor  $\lambda_{aux}(e)$  that shifts focus from exploration to quality over the course of the training episode.<sup>5</sup> This training allows the orchestrator to manage up to 1,500 coordinated tool calls simultaneously, resulting in a 4.5x reduction in execution time compared to traditional sequential agent models.<sup>3</sup>

### The Lifecycle of a Swarm Execution

In a practical application, such as a large-scale DevOps research task or a multi-repository refactoring, the swarm execution follows a recursive, four-stage lifecycle:

1. **Task Decomposition:** The orchestrator agent analyzes the initial prompt and identifies sub-problems that can be executed independently. For a DevOps platform, this might

involve separating infrastructure provisioning, security policy application, and application deployment into distinct tracks.<sup>7</sup>

2. **Specialized Instantiation:** The model creates specialized sub-agents—such as a "Terraform Specialist," a "Policy-as-Code Auditor," and a "Cloud Cost Analyzer"—without human intervention or predefined role templates.<sup>7</sup>
3. **Concurrent Execution:** These sub-agents run in parallel, utilizing specialized tools like code interpreters, web browsers, and file managers to complete their specific subtasks.<sup>7</sup>
4. **Aggregation and Reconciliation:** The orchestrator receives the outputs from the sub-agents and performs a synthesis, reconciling any contradictions and merging the results into a final, structured deliverable.<sup>5</sup>

## The POETIQ META Layer: Recursive Self-Improvement and Auditing

While the Kimi Agent Swarm provides the "muscle" for massive parallel execution, the POETIQ meta-layer provides the "intellect" for rigorous quality control and efficiency. The POETIQ philosophy is built on the principle of "LLMs all the way down," using the recursive power of large language models to build and improve the very systems that utilize them.<sup>6</sup> This layer is model-agnostic, meaning it can leverage the best available frontier models (Gemini 3, GPT-5.1, or Kimi K2.5) to achieve state-of-the-art results across various benchmarks.<sup>6</sup>

### Principle 1: The Prompt as an Interface

POETIQ rejects the idea that a single prompt can extract the full intelligence of a model. Instead, it treats the prompt as an interface to an iterative problem-solving loop.<sup>6</sup> In this loop, the system does not just ask a question; it generates a potential solution, receives concrete feedback (often from a compiler, a test suite, or a visual verification step), analyzes that feedback, and then uses the LLM again to refine the solution.<sup>6</sup>

This "Learned Test Time Reasoning" is particularly effective for high-complexity tasks like the ARC-AGI benchmark, where the system must discover the underlying logic of a grid transformation from a few examples.<sup>6</sup> By generating executable Python code to solve these problems and iteratively refining that code based on its performance on the training examples, POETIQ achieved 54% accuracy on the ARC-AGI-2 Semi-Private set, shattering previous records at less than half the cost of competing methods.<sup>6</sup>

### Principle 2: Autonomous Self-Auditing

A central challenge in autonomous systems is knowing when to stop. The POETIQ meta-layer incorporates an autonomous auditing mechanism that allows the system to monitor its own progress.<sup>6</sup> The meta-system evaluates whether it has enough information to solve the task and determines if the current solution is satisfactory. If the criteria are met, it terminates the

process to avoid wasteful computation and minimize token costs.<sup>6</sup>

This self-monitoring is critical for creating a "Pareto-optimal" solution frontier, where accuracy is maximized while compute cost is minimized.<sup>6</sup> For a DevOps platform, this means the system can autonomously decide if a proposed infrastructure plan is robust enough to be presented for human approval, or if it requires another round of security validation.

Performance Metric	Previous Best (Gemini 3 Deep Think)	POETIQ + Gemini 3	Improvement
ARC-AGI-2 Accuracy	45%	54%	+20% (Relative)
Cost Per Problem	\$77.16	\$30.57	-60% (Cost Saving)
Efficiency Ratio	Baseline	Pareto Frontier	Significant

The success of POETIQ illustrates that the intelligence of an agentic system lies in its *scaffolding* and *recursive loops* as much as it does in the base model weights.<sup>6</sup>

## Implementation Strategy for a Unified DevOps Platform

Building the "best and finest" DevOps platform requires the integration of Kimi's swarm-based execution with POETIQ's recursive auditing. This synergy enables a platform that can handle everything from high-level architectural planning to low-level incident remediation with minimal human intervention.

### Multi-Agent Code-Orchestrated Generation (MACOG) for IaC

Infrastructure as Code (IaC) is inherently a graph-shaped problem, where resources like VPCs, subnets, and databases have complex, non-linear dependencies.<sup>14</sup> A single agent often struggles with the cognitive load of a large-scale deployment. By implementing a Multi-Agent Code-Orchestrated Generation (MACOG) architecture, the task can be decomposed into specialized roles.<sup>14</sup>

The platform should employ a "Blackboard" architecture where agents interact through a shared state managed by a finite-state orchestrator.<sup>14</sup> The roles in this DevOps swarm include:

- **Architect Agent:** Analyzes natural language intent and defines the high-level

infrastructure topology.

- **Provider Harmonizer Agent:** Ensures that the specific cloud provider requirements (AWS, Azure, GCP) are met and that resource ARNs and identifiers are correctly mapped.
- **Engineer Agent:** Generates the actual Terraform or Pulumi configurations using constrained decoding to ensure the code follows HCL grammar and provider field sets.<sup>14</sup>
- **Security Prover Agent:** Subjects the candidate configuration to a battery of policy-as-code validators (e.g., OPA, Sentinel) to detect insecure settings or missing tags.<sup>14</sup>
- **Cost and Capacity Planner Agent:** Evaluates the proposed plan against real-time cloud pricing and quota limits, proposing alternatives if budgets are exceeded.<sup>14</sup>
- **DevOps Auditor Agent:** Acts as the POETIQ-style supervisor, reviewing the entire "Proof-Carrying Bundle" before it is presented for execution.<sup>14</sup>

Evaluation of the MACOG approach shows significant gains in accuracy. For instance, on the IaC-Eval benchmark, GPT-5 improved its score from 54.90 using simple RAG to 74.02 when enhanced with the MACOG multi-agent framework.<sup>14</sup>

## Autonomous SRE and Incident Management

In the operational phase of the DevOps lifecycle, the platform can deploy "AI SRE" agents that operate 24/7 to maintain system reliability.<sup>16</sup> Unlike traditional monitoring tools that merely alert human engineers, these agents can conduct autonomous investigations, test hypotheses, and suggest remediation steps.<sup>16</sup>

The AI SRE agent utilizes the Agent Swarm mode to examine multiple failure paths simultaneously.<sup>16</sup> For example, when an API latency spike occurs, the swarm might concurrently:

1. Query metrics across the request path to identify the point of degradation.<sup>18</sup>
2. Analyze recent deployment logs and configuration changes for correlations.<sup>16</sup>
3. Check infrastructure topology for upstream dependency issues.<sup>18</sup>
4. Examine pod-level resource pressure and container logs.<sup>18</sup>

By documenting its entire chain of reasoning—including the queries executed and the logic behind its conclusions—the AI SRE provides human engineers with an "evidence chain" that transforms hours of manual diagnosis into minutes of focused validation.<sup>16</sup> Early implementations of this technology have reported incident triage speeds up to 10x faster and reductions in observability costs by up to 60% through more efficient data querying.<sup>20</sup>

## Governance, Safety, and the "Observability Conscience"

A primary risk in creating a massive agent swarm for DevOps is the potential for "uncontrolled

autonomy," where agents might make catastrophic changes to production environments. To address this, the platform must integrate rigorous governance frameworks.

## Autonomy Gates and Observability Conscience

The platform should implement "Autonomy Gates," which are conditional checkpoints that determine the level of freedom granted to an agent based on its observed reliability and the risk level of the task.<sup>21</sup> High-risk behaviors—such as deleting a database or modifying a security group—should trigger an immediate human-in-the-loop requirement, while low-risk tasks like gathering diagnostics or restarting a staging service can be executed autonomously.<sup>16</sup>

Complementing this is the "Observability Conscience," an embedded telemetry layer that records every decision an AI agent makes, along with its causal lineage and confidence score.<sup>21</sup> This ensures that all agentic activity is traceable and auditable. By using a Multimodal Artifact File Format (MAIF), every operation is inherently auditable, as the data itself carries its provenance and cryptographic signature.<sup>22</sup>

## Recursive Self-Auditing for Production Safety

The POETIQ self-auditing principle should be applied to the deployment pipeline itself. Before a code change is merged, an "Auditor Swarm" can simulate what-if scenarios, such as a 5% traffic hike or a regional cloud outage, to evaluate the resilience of the proposed change.<sup>23</sup> The system can then generate a "Confidence Score" for the deployment. If the score falls below a certain threshold, the "Autonomy Gate" blocks the merge and provides the engineer with a detailed report on the identified risks.<sup>21</sup>

## Technical Stack and Architectural Selection

For a startup aiming to build the finest DevOps platform, selecting the right orchestration framework is a critical decision. In 2026, the ecosystem has converged on several leading frameworks, each with unique trade-offs.

Framework	Core Model	Architecture	Best Use Case
LangGraph	Graph-based DAGs	Explicit State Transitions	Complex, multi-step workflows with branching and retries
AutoGen	Conversation-base	Actor Model /	Brainstorming and

	d	Message Passing	human-AI collaborative workflows
<b>CrewAI</b>	Role-based	Organizational Structure	Team-like workflows (e.g., Researcher -> Writer)
<b>PydanticAI</b>	Model-agnostic	Highly Abstracted SDK	Developers needing maximum flexibility across LLM providers

For the proposed platform, **LangGraph** is the recommended choice for the core orchestration engine. Its graph-based design allows for precise control over the iterative loops required by the POETIQ layer, and its built-in support for "check-pointing" and "time travel" enables reliable human-in-the-loop interventions.<sup>25</sup> Furthermore, LangGraph's ability to scale horizontally by distributing nodes across Kubernetes workers makes it suitable for the massive parallel execution of Kimi-style swarms.<sup>27</sup>

## Future Outlook: The Evolution toward Autonomous Resilience

The integration of Kimi's Agent Swarm and the POETIQ meta-system represents the first step toward "Cognitive DevOps," where the platform is not just a tool but an intelligent team member.<sup>18</sup> Future developments will likely focus on "Episodic Memory" for agents, allowing them to learn from every past incident and deployment failure to improve their future decision-making.<sup>30</sup>

Moreover, as models like Kimi K2.5 continue to optimize their MoE architectures, the cost of running massive swarms will decrease, making it viable for even small startups to deploy "swarms of 1,000 agents" for tasks like exhaustive security testing or real-time cost optimization. The goal of the "finest" DevOps platform must be to build a self-optimizing, self-healing digital ecosystem where the human role shifts from reactive firefighting to high-level strategic guidance and creative architecture.

In conclusion, by synthesizing the parallel throughput of Kimi's Agent Swarm with the recursive reliability of the POETIQ meta-layer, a DevOps platform can achieve unprecedented levels of efficiency, safety, and insight. The path forward lies in embracing the "Prompt as an Interface" philosophy, treating every operational task as an iterative reasoning loop, and

maintaining a steadfast commitment to transparency through a robust "Observability Conscience." This architectural blueprint provides a comprehensive path toward defining the future of autonomous software delivery.

## Works cited

1. Kimi K2.5: Complete Guide to Moonshot's AI Model - Codecademy, accessed on February 10, 2026,  
<https://www.codecademy.com/article/kimi-k-2-5-complete-guide-to-moonshots-ai-model>
2. Kimi K2.5 API | Together AI, accessed on February 10, 2026,  
<https://www.together.ai/models/kimi-k2-5>
3. Kimi K2.5 Tech Blog: Visual Agentic Intelligence - Kimi AI, accessed on February 10, 2026, <https://www.kimi.com/blog/kimi-k2-5.html>
4. Moonshot AI's Kimi K2.5 Expands What Open-Weight Models Can Do - Alwire - HPCwire, accessed on February 10, 2026,  
<https://www.hpcwire.com/aiwire/2026/01/30/moonshot-ais-kimi-k2-5-expands-what-open-weight-models-can-do/>
5. Kimi K2.5 in 2026: The Ultimate Guide to Open-Source Visual Agentic Intelligence, accessed on February 10, 2026,  
<https://dev.to/czmilo/kimi-k25-in-2026-the-ultimate-guide-to-open-source-visual-agentic-intelligence-18od>
6. Traversing the Frontier of Superintelligence - Poetiq, accessed on February 10, 2026, [https://poetiq.ai/posts/arcagi\\_announcement/](https://poetiq.ai/posts/arcagi_announcement/)
7. Kimi K2.5 and Agent Swarm: A Guide With Four Practical Examples - DataCamp, accessed on February 10, 2026,  
<https://www.datacamp.com/tutorial/kimi-k2-agent-swarm-guide>
8. One Hundred Agents, One Command, Kimi K2.5 Just Rewrote the Rules of Automation, accessed on February 10, 2026,  
<https://medium.com/@cognidownunder/one-hundred-agents-one-command-kimi-k2-5-just-rewrote-the-rules-of-automation-0e9db50bc694>
9. NEW Kimi Swarm Update is INSANE!, accessed on February 10, 2026,  
[https://www.youtube.com/watch?v=HQdiZO\\_OVRE](https://www.youtube.com/watch?v=HQdiZO_OVRE)
10. [2602.02276] Kimi K2.5: Visual Agentic Intelligence - arXiv, accessed on February 10, 2026, <https://arxiv.org/abs/2602.02276>
11. Poetiq Did It!!! Poetiq Has Beaten the Human Baseline on Arc-AGI 2 (<60%) | "Poetiq's approach of building intelligence on top of any model allowed us to integrate the newly released Gemini 3 and GPT-5.1 models within hours of their release to achieve the SOTA-results presented here." : r/ - Reddit, accessed on February 10, 2026,  
[https://www.reddit.com/r/accelerate/comments/1p2grr3/poetiq\\_did\\_it\\_poitiq\\_has\\_beaten\\_the\\_human/](https://www.reddit.com/r/accelerate/comments/1p2grr3/poetiq_did_it_poitiq_has_beaten_the_human/)
12. Poetiq Achieves SOTA on ARC-AGI 2 Public Eval : r/singularity - Reddit, accessed on February 10, 2026,  
[https://www.reddit.com/r/singularity/comments/1pu5mhk/poetiq\\_achieves\\_sota\\_o](https://www.reddit.com/r/singularity/comments/1pu5mhk/poetiq_achieves_sota_o)

n\_arcagi\_2\_public\_eval/

13. ARC AGI 2 is solved by poetiq! : r/accelerate - Reddit, accessed on February 10, 2026,  
[https://www.reddit.com/r/accelerate/comments/1pu7llk/arc\\_agi\\_2\\_is\\_solved\\_by\\_poetiq/](https://www.reddit.com/r/accelerate/comments/1pu7llk/arc_agi_2_is_solved_by_poetiq/)
14. Multi-Agent Code-Orchestrated Generation for Reliable Infrastructure-as-Code - arXiv, accessed on February 10, 2026, <https://arxiv.org/html/2510.03902v1>
15. Self-Auditing Deep Learning Pipelines for Automated Compliance Validation with Explainability, Traceability, and Regulatory Assurance - ResearchGate, accessed on February 10, 2026,  
[https://www.researchgate.net/publication/399185737\\_Self-Auditing\\_Deep\\_Learning\\_Pipelines\\_for\\_Automated\\_Compliance\\_Validation\\_with\\_Explainability\\_Traceability\\_and\\_Regulatory\\_Assurance](https://www.researchgate.net/publication/399185737_Self-Auditing_Deep_Learning_Pipelines_for_Automated_Compliance_Validation_with_Explainability_Traceability_and_Regulatory_Assurance)
16. What is an AI SRE? - Neubird, accessed on February 10, 2026,  
<https://neubird.ai/glossary/what-is-an-ai-sre/>
17. Overview of Azure SRE Agent Preview | Microsoft Learn, accessed on February 10, 2026, <https://learn.microsoft.com/en-us/azure/sre-agent/overview>
18. What Are AI Agentic Assistants in SRE and Ops, and Why Do They Matter Now? - Medium, accessed on February 10, 2026,  
<https://medium.com/@ad.shaikh2003/what-are-ai-agentic-assistants-in-sre-and-ops-and-why-do-they-matter-now-7ed5f6ac5a56>
19. Bridging the SRE gap: Towards autonomous observability and RCA - Thoughtworks, accessed on February 10, 2026,  
<https://www.thoughtworks.com/en-us/insights/blog/generative-ai/bridging-the-SRE-gap-towards-autonomous-observability-and-RCA>
20. Observe Introduces AI SRE and o11y.ai Agents, Accelerating Developer Productivity While Cutting Enterprise Observability Costs, accessed on February 10, 2026,  
<https://www.observeinc.com/news-pr/observe-introduces-ai-sre-and-o11y-ai-agents-accelerating-developer-productivity-while-cutting-enterprise-observability-costs>
21. Operationalizing Agentic AI in Enterprise Delivery: Observability ..., accessed on February 10, 2026, <https://engrxiv.org/preprint/download/5657/9463/7893>
22. MAIF: Enforcing AI Trust and Provenance with an Artifact-Centric Agentic Paradigm - arXiv, accessed on February 10, 2026,  
<https://arxiv.org/html/2511.15097v1>
23. Day 45: Compare and contrast Traditional ML vs Deep Learning vs Agentic AI (Part 3), accessed on February 10, 2026,  
<https://luxananda.medium.com/day-45-compare-and-contrast-traditional-ml-vs-deep-learning-vs-agentic-ai-part-3-d418bb2ca95f>
24. Automate CircleCI Monitoring and Triage with Edge Delta's AI Teammates - Edgedelta, accessed on February 10, 2026,  
<https://edgedelta.com/company/blog/automate-circleci-monitoring-and-triage-with-edge-deltas-ai-teammates>
25. LangGraph vs AutoGen: Multi-Agent AI Framework Comparison - Leanware,

- accessed on February 10, 2026,  
<https://www.leanware.co/insights/auto-gen-vs-langgraph-comparison>
26. AutoGen vs LangGraph: Comparing Multi-Agent AI Frameworks - TrueFoundry, accessed on February 10, 2026,  
<https://www.truefoundry.com/blog/autogen-vs-langgraph>
27. LangGraph vs AutoGen: Comparing AI Agent Frameworks - PromptLayer Blog, accessed on February 10, 2026,  
<https://blog.promptlayer.com/langgraph-vs-autogen/>
28. Top AI Agent Orchestration Platforms in 2026 - Redis, accessed on February 10, 2026, <https://redis.io/blog/ai-agent-orchestration-platforms/>
29. Agentic AI: The Complete Guide to Autonomous Intelligence and the Future of AI Agents - Services Ground, accessed on February 10, 2026,  
<https://servicesground.com/blog/agentic-ai/>
30. SiriusS: Self-improving Multi-agent Systems via Bootstrapped Reasoning - arXiv, accessed on February 10, 2026, <https://arxiv.org/html/2502.04780v1>
31. Top 25 GenAI Patterns in Production: Agentic and Non-Agentic - Medium, accessed on February 10, 2026,  
<https://medium.com/@raj-srivastava/top-25-genai-patterns-in-production-agentic-and-non-agentic-2e1b3cccd8517>