



10X FASTER TROUBLESHOOTING:

THE AI SRE PLAYBOOK

 **OBSERVE**
BY SNOWFLAKE

INTRODUCTION

Despite the exponential growth in observability data, organizations still face the same challenges during production incidents. Engineers context-switch between tools, monitor accumulating dashboards and alerts, and often rely on a small number of experts for resolution.

Major incidents continue to require large teams and extended hours to resolve. For many organizations, detection has improved, but investigation remains the dominant contributor to MTTR.

Google's research on debugging distributed systems shows that investigation is inherently cyclical. Their 2019 study, published in ACM Queue ("Debugging Incidents in Google's Distributed Systems"), describes how engineers debug production issues at scale. The canonical debugging journey is not a linear path from alert to resolution. It's a cyclical, non-sequential process in which engineers repeatedly form hypotheses, query telemetry, discard dead ends, and start again. Google's SRE Handbook frames it more formally. Troubleshooting is "an application of the hypothetico-deductive method"—engineers iteratively hypothesize potential causes and test those hypotheses against available data.

Even highly skilled engineers are slowed by the structure of the investigation process. Each cycle demands context gathering, signal correlation, and synthesis across fragmented data. This often happens under significant time pressure. Increased complexity expands the hypothesis space and prolongs investigation.

AI can significantly accelerate the investigation loop. However, meaningful gains appear only when built on the right architectural foundation. Over the past year, a wave of AI SRE offerings has entered the market. AI systems are now designed to augment site reliability engineering workflows, particularly investigation and root cause analysis.

Many current implementations are limited by observability platforms that, despite translating natural language into queries, lack unified context. As a result, they cannot effectively reason across signals, understand system relationships, or scale operations efficiently.

This ebook defines what truly effective AI-driven investigation looks like. It provides a practical framework for distinguishing between AI SRE tools that achieve real operational improvements and those that merely offer workflow optimizations without improving accuracy, speed, or cost.

SECTION 1

THE PROMISE OF AI SRE: SCALING INVESTIGATION IN COMPLEX SYSTEMS

Why Investigation Is Harder Than Ever

Three forces are simultaneously compounding the investigation problem.

- **Data volume is increasing faster than existing systems were designed to handle.** As system complexity accelerates, telemetry output grows accordingly. Containerization, microservices, AI-generated code, and agentic workflows have multiplied the observability data teams must analyze when troubleshooting.
- **System dependencies have become more opaque due to modern architectures.** These often create deep and undocumented dependency chains. Understanding these relationships requires institutional knowledge, which is often incompletely documented in runbooks or modeled in observability tools.
- **Expertise is often concentrated among a few individuals.** In most organizations, only a small number of senior engineers possess the experience and familiarity required for complex investigations. When that expertise is unavailable, investigation cycles tend to lengthen.

What Investigation Actually Looks Like

The outline of an incident investigation is deceptively simple: detect the problem, scope the blast radius, query telemetry, form hypotheses, and resolve. In practice, each of those steps involves substantial manual work.

- **Trace investigation requires navigating to a tracing UI,** formulating queries, analyzing span hierarchies, checking for errors across services, and writing up findings. Even in moderately complex systems, a single trace can require tens of minutes to fully analyze.
- **Log analysis is where investigation time quietly compounds.** An engineer tracking down an unexpected access pattern change might start by orienting to the data structure. Then, they run a series of exploratory queries to establish a baseline. Next, they build trend analysis across time windows, break down results by role, user, or region, and finally, deep-dive into whatever looks anomalous. Each step feels small. The full sequence routinely takes 30 to 45 minutes. This is for a single line of inquiry.
- **Latency investigations triggered by an API gateway alert are even more involved.** Engineers review alert details and consult runbooks. They run multiple queries across different signal types and correlate findings across services. Synthesizing recommendations can become time-consuming. Latency investigations often stretch beyond an hour as engineers iterate through competing hypotheses and cross-service signals.

And this is only one thread of investigation. Incidents usually involve several parallel threads, multiple hypotheses, and dead ends that force engineers to start over.

The cumulative cost is significant. In some organizations, particularly those operating complex distributed systems, incident response and its aftermath can consume as much as 30–50% of engineering capacity.

SECTION 2

THE PROMISE AND REALITY OF MOST AI FOR OBSERVABILITY

The Rise of AI SRE

AI-powered investigation tools have quickly become common in the observability market. The integration of LLM capabilities with telemetry platforms has made the AI SRE a standard offering. Many mainstream observability tools now include conversational AI that promises to query data, correlate signals, and identify root causes.

The drivers behind this wave are real. System complexity has outpaced human capacity to investigate incidents manually. Microservices, multi-cloud architectures, and AI-generated code create complex dependency graphs. No single engineer can hold all of these in their head. Teams need help and AI is the obvious candidate to provide it.

Despite the rapid adoption of AI features, real-world results often fall short of claims of faster triage, automated root-cause analysis, and natural-language investigation. The effectiveness of these features differs widely based on the technology used.

The Need to Go Deeper Than Chat

While AI SREs are now commonplace, the ability to ask questions in plain language alone does not provide differentiated outcomes.

The most important questions are more complex:

- **How accurate is the answer?** Does the AI SRE have access to enough context to give a correct response, or does it generate answers without sufficient grounding? Did the underlying observability tool retain enough data for the AI SRE to answer the question?
- **How fast do I get responses?** Can the AI reason across large volumes of telemetry in real time, or does it choke on the scale of data involved in a real incident?
- **Does it work without exploding costs?** Can the AI query the underlying observability tools sustainably, without escalating observability spend?
- **Does it connect infrastructure to business?** Can you troubleshoot not just with infrastructure context, but also with business context? Can you ask questions about customers, shopping carts, transactions, and other business concepts?

In many cases, AI SRE capabilities are limited by structure. They often rely on observability platforms that weren't built for AI and don't scale efficiently.

What Happens Without Architecture Optimized for Performance and Cost

When AI operates on top of observability architectures not designed for unified, large-scale reasoning, predictable failure modes emerge. These architectures make AI less reliable.

- **Unreliable answers happen when data is fragmented.** When telemetry is siloed by different models, schemas, or retention boundaries, the AI SRE sees only a partial system view. Missing signals must be inferred. Human experts must then validate conclusions. This erodes intended efficiency gains.
- **Correlation is shallow or manual,** lacking a unified data model that captures relationships among services, deployments, users, and business entities. Often, correlation is limited to temporal proximity and surface-level identifiers. It becomes difficult to connect infrastructure behavior to customer impact and business outcomes.
- **High costs constrain investigation.** Observability platforms are frequently criticized for high data and query costs, leading organizations to limit retention or restrict analysis. These constraints become more acute at scale. When AI must execute broad, cross-signal queries across large datasets, the economic burden can make comprehensive investigation impractical.

The result is an AI that functions more as a summarizer than an investigator. It can describe what dashboards already show, but it lacks the depth of context needed to uncover what they omit.

SECTION 3

WHAT MAKES AI SRES ACTUALLY WORK

Architecture Is the Differentiator

Effective AI investigation depends on multiple architectural components that must work together. None is sufficient on its own.

1 Unified low-cost storage for all telemetry.

AI-driven troubleshooting requires unified access to logs, metrics, traces, and events. When telemetry is fragmented by data models, retention policies, or query engines, cross-signal reasoning becomes slower, more expensive, and less reliable. A unified storage layer allows AI systems to reason efficiently across signal types. This removes the need to assemble partial results from multiple backends.

Cost efficiency is equally critical here. When storing or querying, comprehensive telemetry becomes prohibitively expensive. Teams often respond by sampling data, dropping high-cardinality fields, or shortening retention periods, creating blind spots for AI-driven investigation. Architectures built on scalable, cost-efficient object storage with strong compression capabilities enable the retention and analysis of full-fidelity telemetry at scale.

2 A context graph that models semantic relationships.

Raw telemetry alone is insufficient for effective AI-driven investigation. Beyond logs, metrics, and traces, the system must understand topology and relationships, including service dependencies, pod-to-node mappings, and deployment changes. This contextual layer transforms disconnected signals into a structured representation of the system's operation.

Beyond infrastructure context, an effective AI SRE should also incorporate business context: users, transactions, shopping carts, shipments, and other business-level entities. When operational and business contexts are brought together, the AI SRE becomes a very powerful tool, directly connecting infrastructure behavior to real-time business outcomes.

A context graph, or knowledge graph, represents these entities and their relationships with explicit semantic meaning. During investigation, the AI can traverse this graph, moving from affected customers or services through their dependencies, infrastructure, recent deployments, and downstream impacts, enabling structured reasoning across the entire environment.

3 Tight integration between the AI SRE and the data and semantic foundations.

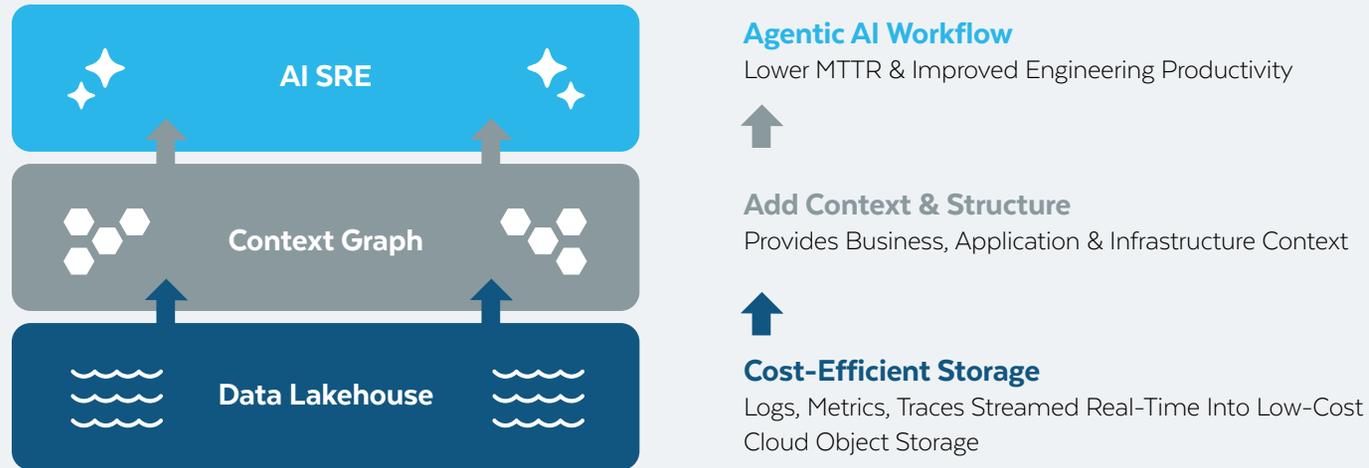
An AI SRE that is merely layered on top of existing observability tools inherits the architectural constraints of those platforms. Without deep integration into the underlying data and semantic models, the system cannot fully optimize how telemetry is accessed, correlated, or reasoned over.

These limitations often surface as unreliable answers, inconsistent query performance, timeouts during complex analysis, or escalating costs as the AI SRE repeatedly queries external observability systems.

By contrast, when AI is natively integrated with an observability platform designed to provide cost-efficient scale and unified context, it can optimize queries, data access, and reasoning paths holistically. The result is more predictable performance, higher accuracy, and greater cost efficiency.

A Modern Architecture for Observability at Scale

Three Related, Dependent, Components



What AI-Optimized Architecture Enables in Practice

With a well-architected foundation, AI can handle the most time-consuming tasks during investigations:

- **Searching and analyzing logs across distributed systems:**
The AI SRE can query across log data at scale, surface the most relevant entries, and recognize patterns that may be missed by manual review.
- **Investigating error patterns and anomalies:**
The AI SRE can correlate error spikes across services, identify common attributes, and suggest likely causal chains, tasks that usually require expert knowledge.
- **Querying metrics and tracing data in context:**
The AI SRE queries across signal types to present a unified view, eliminating the need for manual correlation across multiple interfaces.
- **Generating hypotheses with confidence levels:**
Well-designed AI surfaces multiple hypotheses ranked by likelihood, grounded in actual telemetry, mirroring the approach of expert investigators.
- **Correlating deployments and CI/CD changes:** The AI SRE can automatically determine whether recent deployments, configuration changes, or CI/CD commits correlate with the onset of an issue, streamlining a typically tedious investigation step.

Where AI Delivers the Biggest Gains

Not all investigation tasks benefit equally from AI. The highest-impact use cases share a common characteristic: they involve synthesizing large volumes of data across multiple sources under time pressure.

- **Exploratory investigations**, where the engineer is unsure where to begin, benefit most. The AI can autonomously query multiple data sources, follow relationship chains, and surface relevant signals without requiring precise initial queries.
- **Log analysis involving pattern matching across large datasets** is another high-impact area. AI can quickly scan for anomalies and maintain context across large datasets, outperforming manual review.
- **Incident triage**, particularly in the initial minutes after an alert, benefits from AI's ability to quickly summarize the issue, identify affected services and users, and suggest investigation paths when time and context are most critical.

SECTION 4

REAL-WORLD RESULTS

Accelerating Troubleshooting and Reducing Toil

The investigation loop- detect, scope, query, hypothesize, resolve- remains with AI, but is significantly accelerated. Many investigation tasks that previously required tens of minutes to hours can now be completed in a fraction of the time. AI makes the process accessible to all users, not just a select few experts.

The Productivity Math

Analysis of thousands of AI-assisted investigation sessions shows consistent productivity gains. Using a methodology adapted from Anthropic’s framework for measuring AI productivity, we compared estimated manual investigation time with AI-assisted completion time across a wide range of tasks, including scenarios where AI provided limited benefit.

Across measured interactions:

- Productivity gains clustered in the **3-10x range**
- **30% of interactions** showed an improvement of more than 5x
- **5% of interactions** exceeded 10x improvement

The largest gains occurred in investigations involving large-scale data synthesis across multiple signal types.

Representative Tasks with Time Savings

Task	Human Time	AI-Assisted Time	Multiplier
Trace investigation	15 min	4 min	3.8x
Production issue investigation	85 min	25 min	3.4x
CloudTrail log analysis	35 min	6 min	5.8x
API gateway latency alert	90 min	9 min	10x
QA environment performance	30 min	2 min	15x

What Teams Are Seeing

Productivity has improved across roles, team sizes, and use cases, leading to faster incident triage and more efficient investigations.

Foursquare is leveraging AI SRE to improve incident response. “Observe’s AI SRE and MCP Server have the potential to transform how we investigate incidents and reduce the time engineers spend on resolving issues by providing faster, more contextual insight into system behavior.”

At Topgolf, the benefits extend beyond reactive troubleshooting to proactive reliability.

“It’s a major unlock for our SRE practice, empowering teams to proactively detect system issues and minimize disruptions, resulting in improved player experiences and greater system stability.”

These results are consistent across teams. Other customer feedback highlights that AI SRE enables faster information gathering and reduces manual review time.

- In one case, a SaaS provider for the automotive industry reduced incident investigation time from over 3 hours to minutes, significantly improving the customer experience.
- Another team saw fewer escalations and faster support ticket resolution after adoption.

Overall, teams report less manual effort and faster resolution times.

From Hours To Minutes

These improvements reduce investigation time and, in many cases, overall MTTR. Many investigation workflows that previously required hours of manual work can now be completed in minutes. This progress results from compressing the entire investigation, including faster context gathering, hypothesis testing, and broader team access to investigation capabilities.

CONCLUSION

Current challenges with investigation and troubleshooting stem not from a shortage of dashboards, alerting, or tooling, but from a structural mismatch between the volume and complexity of available data and engineers' capacity to synthesize it under time pressure.

AI has the potential to close that gap, but only when it operates on a well-designed architectural foundation: cost-efficient telemetry storage; a context graph that models technical and business relationships; and an AI SRE deeply integrated with both. When these elements are present, investigation becomes faster, more consistent, and accessible to a broader range of engineers.

For engineering leaders, the strategic question is no longer whether AI will be part of observability, but whether your architecture enables AI to deliver sustained and measurable productivity gains.

SOURCES:

"Debugging Incidents in Google's Distributed Systems," ACM Queue, 2020. Google SRE Handbook, Chapter 12: "Effective Troubleshooting."
Productivity analysis methodology adapted from Anthropic's framework for measuring AI-assisted task completion.



OBSERVE
BY SNOWFLAKE

observeinc.com

