



Why AI Slows Down When It Reaches Production

The next productivity frontier
February 2026

Contents

Key Takeaways p 2

The Production Gap p 4

Why Pilots Break in Production p 5

The Productivity Paradox p 7

Ownership, Governance, and Risk p 8

Delivery Capacity as the Bottleneck p 9

Comparative Deployment Outcomes p 10

Implications for CTOs p 11

Executive Summary

Enterprise AI adoption has reached a structural bottleneck. While experimentation is widespread, sustained production impact remains rare. This paper examines why.

Key Takeaways

AI does not stall in production because the technology fails.

It stalls because organizations are not built to operate systems that change continuously. AI introduces constant updates into delivery, governance, and cost structures that were designed for infrequent, discrete releases.

The biggest failure point is ownership, not models.

In pilots, responsibility is implicit. In production, it must be explicit. When no single owner is accountable for reliability, cost, and compliance, AI systems degrade quickly and adoption stalls.

Productivity gains do not translate into faster delivery.

Individual developers may become 20–40% more productive with AI tools, yet most organizations see no acceleration in end-to-end delivery. The execution system absorbs the efficiency gain instead of converting it into throughput.

Governance and risk become decisive at scale.

Once AI enters production, decisions shift from innovation teams to legal, compliance, security, finance, and executive leadership. At this stage, AI is evaluated as infrastructure, judged on accountability, auditability, and operational risk rather than technical novelty.

Similar technology produces radically different outcomes.

Organizations deploying comparable models and platforms achieve very different results depending on execution structure. Where AI is treated as production infrastructure with clear accountability, it delivers sustained value. Where responsibility is fragmented, it remains marginal.

AI amplifies the delivery model it is placed into.

If the underlying execution model is not designed for continuous operation, AI will slow down until the organization adapts. At scale, AI success is an execution problem before it is a technology problem.

88%

of large enterprises use AI in at least one function, but only 33% have begun scaling AI across the organization; fewer than 6% report enterprise-level EBIT impact above 5% (McKinsey Global AI Survey 2025).

Developers report

20-40%

productivity gains from AI tools, yet fewer than 30% of organizations see faster end-to-end delivery cycles (MIT Sloan research).

Large-scale deployment studies show that organizations using comparable models achieve divergent results primarily due to differences in execution ownership and governance, not technical capability (academic deployment surveys).

Over

40%

of agentic AI initiatives are expected to be cancelled by 2027 due to cost escalation, unclear ownership, and insufficient risk controls (Gartner).

AI governance failures have already generated cumulative global losses exceeding

\$4.4 billion,

shifting oversight to legal, compliance, and executive risk functions (EY).

AI initiatives involving legal, security, and finance from the start are more than twice as likely to reach sustained production use than those driven solely by innovation teams (Stanford HAI / MIT Sloan).

1.

The Production Gap



In a large enterprise, an AI system can look production-ready right up until the moment it has to behave like one.

A fraud detection model passes offline validation. A support agent handles test conversations flawlessly. Then the system is wired into live applications, real data streams, and on-call rotations, and suddenly the problems begin. Latency spikes. Alerts trigger with no clear owner. Rollbacks become manual. What worked in isolation struggles once it is exposed to the realities of production.

Recent industry data confirms that this is not an isolated pattern. According to [McKinsey's 2025 Global AI Survey](#), 88% of large enterprises report regular use of AI in at least one business function, yet only 33% have begun scaling AI systems across the organization. More strikingly, fewer than 6% report measurable

enterprise-level EBIT impact above 5%. The gap between experimentation and production remains the dominant structural failure mode in enterprise AI adoption.

The problem is simpler than it sounds. AI introduces continuous change into organizations that are built to approve, release, and operate software in discrete steps. Once AI systems are expected to behave like production infrastructure, gaps in ownership, governance, and cost control become impossible to ignore.

The result is predictable. Reviews queue up. Incidents have no clear owner. Costs rise before anyone notices. To protect stability, teams slow delivery. What looks like a delivery problem is an execution and ownership problem.

2.

Why Pilots Break in Production

The first failure rarely comes from the model. It comes from ownership.

In pilots, no one is on call. In production, someone must be. When an AI system degrades at 2 a.m., the question is not whether the model works, but who is responsible for fixing it. In many organizations, that question has no clear answer.

The second failure is cost visibility. Pilots run on limited volumes and controlled datasets. Production does not. By the time real usage patterns emerge, spend has already escaped forecast and finance is reacting after the fact.

The third failure is governance latency. Review and approval processes designed for occasional system changes are suddenly asked to keep up with continuous AI-driven updates. They cannot, so delivery slows to preserve stability.

Governance failure is no longer theoretical. [Gartner](#) estimates that over 40% of agentic AI initiatives will be cancelled by 2027, primarily due to escalating operational costs, unclear ownership, and insufficient risk controls. Separately, EY reports that AI-related governance failures have already generated cumulative losses exceeding \$4.4 billion globally, shifting AI oversight firmly into the remit of legal, compliance, and executive risk functions.

This shift is visible in buying behavior. According to [Stanford HAI](#) and MIT Sloan research, AI initiatives that involve legal, security, and finance stakeholders from the start are more than twice as likely to reach sustained production use compared to initiatives sponsored exclusively by innovation or IT teams. Once AI enters production, the dominant success factor becomes decision defensibility, not technical novelty.

These issues are invisible in proof of concept. They only surface once AI systems are expected to behave like real production infrastructure.

This helps explain why, despite widespread experimentation, only a minority of organizations have deployed AI systems broadly in production environments.

Enterprise survey data illustrates how sharply AI adoption diverges once organizations attempt to move from pilots to production-scale impact.

At this stage, the nature of the decision changes. AI adoption is no longer driven by early adopters or innovation teams, but by risk owners. Legal, compliance, security, finance, and executive leadership become central to the decision process. AI decisions start to resemble infrastructure decisions, judged less on capability and more on accountability, auditability, and long-term operational risk.

Organizational Stage	Percentage of Respondents	Enterprise-Level EBIT Impact
Experimentation/Piloting	~67%	Minimal (< 5%)
Early Scaling	~33	Limited (39% report any impact)
High Performers (5%+ EBIT)	~6%	Significant (5%+)

Source: McKinsey Global AI Survey 2025; enterprise respondents categorized by AI maturity and reported EBIT impact.

3.

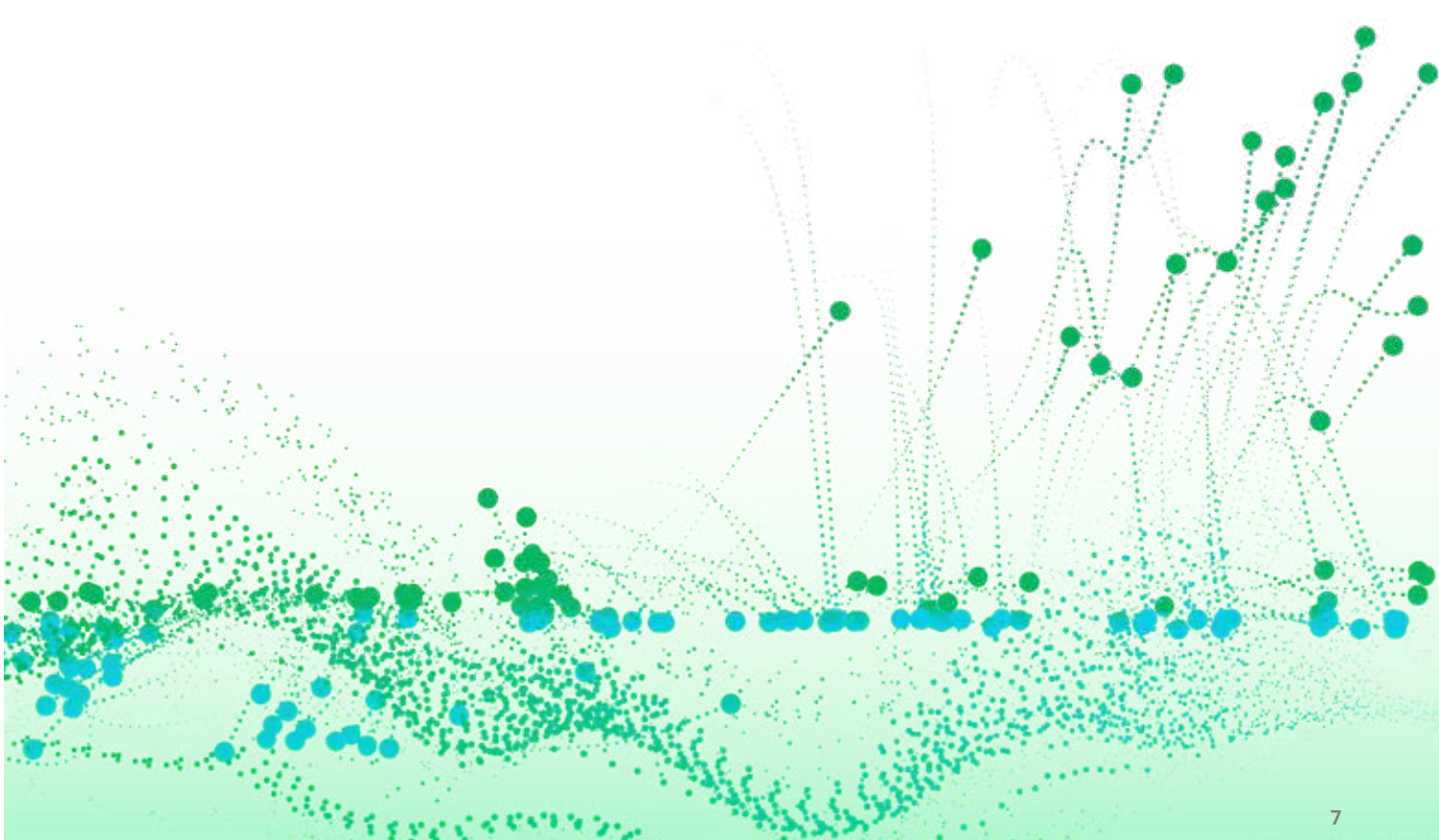
The Productivity Paradox

In most production environments, most of the engineering capacity is already consumed by maintenance and operational work. AI initiatives therefore compete for scarce delivery of bandwidth rather than creating new capacity, which is why early gains rarely translate into faster execution at scale.

This explains a paradox observed consistently in enterprise studies. [MIT Sloan](#) research shows that while individual developers using AI tools report productivity gains of 20–40%,

fewer than 30% of organizations observe any acceleration in end-to-end delivery cycles. The efficiency gain is absorbed by the system rather than converted into throughput, because execution constraints sit downstream of the individual contributor.

Most CTOs have seen this pattern before. Systems look production-ready on their own, then break down once they meet real roadmaps, real governance, and real operational constraints.



4.

Ownership, Governance, and Risk

At this stage, the nature of the decision changes. AI adoption is now driven by risk owners. Legal, compliance, security, finance, and executive leadership become central to the decision process. AI decisions start to resemble infrastructure decisions, judged less on capability and more on accountability, auditability, and long-term operational risk.

Once AI systems move beyond isolated pilots, the question shifts from whether the technology works to whether the organization can defend, operate, and sustain it over time. Responsibility, cost exposure, and governance are no longer secondary considerations; they become primary decision criteria. As a result, AI initiatives are increasingly evaluated through the same lens as other mission-critical systems, where unclear ownership or fragmented accountability represents unacceptable risk.

This transition explains why many AI programs slow down precisely at the point where they are expected to scale. The limiting factor is no longer technical feasibility, but the organization's ability to assign clear responsibility, enforce governance consistently, and absorb AI into existing operational and risk management structures.

5.

Delivery Capacity as the Bottleneck

The constraint usually shows up before the model does.

In most large organizations, engineering teams are already operating near capacity. Roadmaps are full. On-call rotations are tight. Maintenance and incident work consume the majority of available time. AI initiatives do not arrive in a greenfield environment; they arrive on top of an already saturated system.

When AI is added without changing how work flows through that system, its effects remain local. Individual tasks get faster, but reviews, integration, testing, deployment, and incident response do not. The system absorbs the gain instead of converting it into throughput.

This is why many CTOs see the same pattern: developers report higher efficiency, but delivery dates do not move. In practice, the primary constraint is not the model itself, but the execution system and ownership structure around it.

6.

Comparative Deployment Outcomes

This contrast is not anecdotal. Large-scale empirical studies of machine learning deployment, including multi-year surveys of real production systems, consistently show that technical parity does not translate into operational parity. Across organizations deploying similar models under comparable conditions, differences in governance, accountability, and execution ownership emerge as the primary drivers of divergent outcomes in production. This pattern has been systematically documented in [Challenges in Deploying Machine Learning: A Survey of Case Studies](#) and its extended academic version published through the [University of Sheffield](#).

Consider two large enterprises rolling out AI for customer support. Both use the same models, the same cloud provider, and comparable budgets. On paper, the technical choices are nearly identical. In practice, the outcomes diverge quickly.

In the first organization, the system is framed as an experiment. An innovation team builds it. IT supports it intermittently. Legal reviews it late in the process. Responsibility is distributed across functions, which in practice means it is unclear who owns the system end to end.

When responses are inconsistent, it is dismissed as a pilot issue.

When compliance asks for traceability, it is deferred. When costs begin to rise, finance asks who is accountable, and no clear answer emerges. Usage remains limited, partly because the system is unreliable, and partly because no one is responsible for making it reliable.

In the second organization, the technology is similar but the framing is different. The system is treated as production infrastructure from the start. Ownership is explicit. When outputs degrade, there is a named owner responsible for fixing them. When compliance requests traceability, the mechanisms are already in place. When costs increase, accountability is clear.

The difference comes down to ownership. When AI systems are owned as part of core execution, they become reliable enough to be used and to produce business value. When they are not, they remain marginal.

At scale, AI breaks when it is introduced into execution systems that were never designed to run it.

In organizations where ownership is diffuse, costs are opaque, and operational responsibility is deferred, AI predictably stalls.

In organizations that treat AI as production infrastructure from day one, the same technology becomes reliable enough to matter.

7.

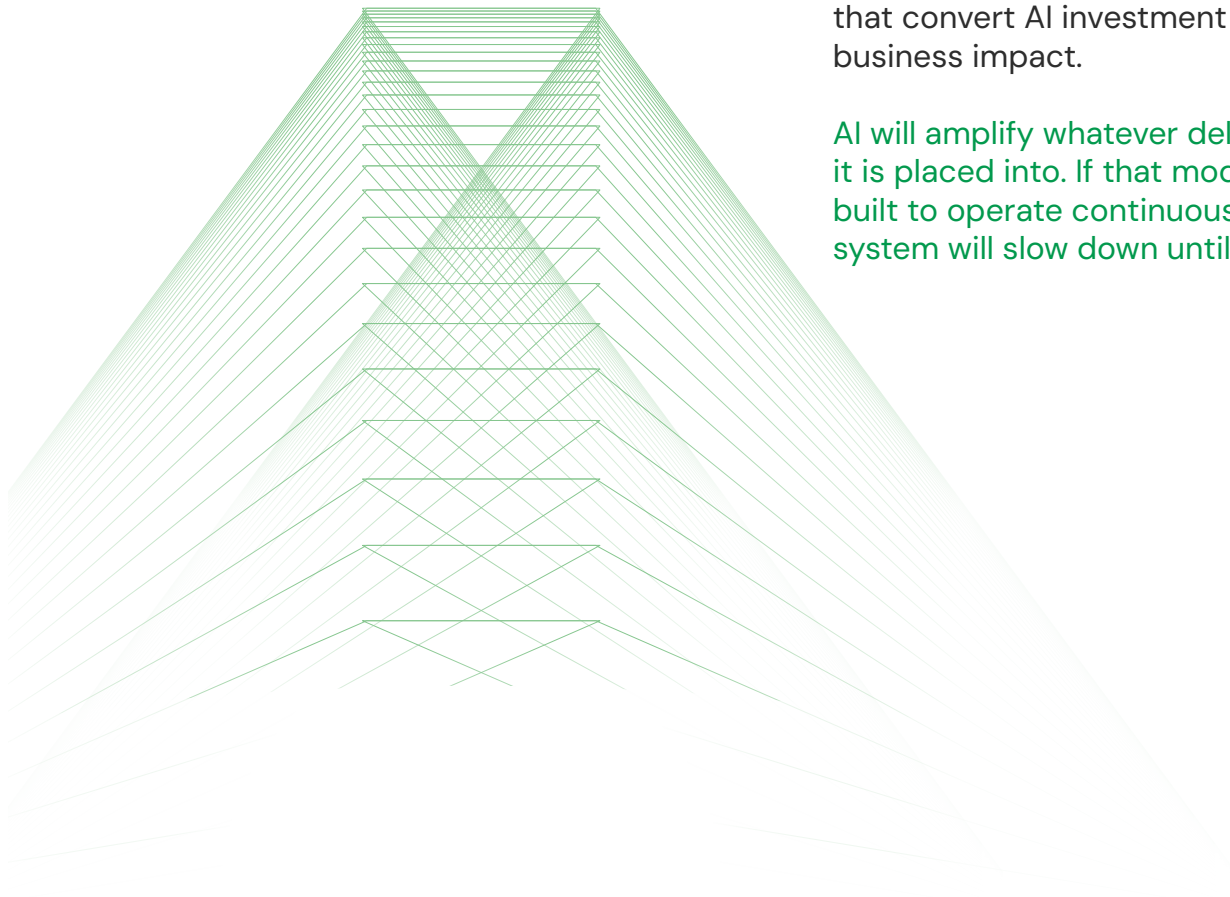
Implications for CTOs

For CTOs, this comes down to execution. At scale, AI does not fail because the technology is immature. It fails because responsibility is fragmented.

When ownership is unclear, costs become opaque, governance slows delivery, and operational risk accumulates without a single point of accountability. The same technology produces radically different outcomes depending on whether execution accountability is clearly defined end to end or diffused across disconnected teams and vendors.

Across all major enterprise studies, one pattern is consistent. Organizations that treat AI as experimental technology struggle to scale. Organizations that treat AI as production infrastructure, with explicit execution accountability, shared governance, and operational rigor comparable to other mission-critical systems, are the ones that convert AI investment into durable business impact.

AI will amplify whatever delivery model it is placed into. If that model is not built to operate continuously, the system will slow down until it is.





Bucharest, Romania

Green Court, 4C Gara Herastrau Street,
Building B, 3rd Floor, District 2

Phone: +40 374 400 731

www.rinf.tech