

Deloitte.



Why enterprises fail to scale AI agents and what can fix it

March 2026

Engineering

Table of Contents

- What makes an agent “scaled”?.....3**
- The four impediments to scaling AI agents5**
- Impediment 1: The organizational immune response6**
- Impediment 2: When probabilistic meets deterministic9**
- Impediment 3: The “dumb RAG” trap 11**
- Impediment 4: The accountability vacuum 13**
- The execution strategy: Dual tracks..... 15**
 - Track 1: What IT should build..... 15
 - Track 2: What the organization should change..... 16
- The agent fast lane framework: Risk-tiered deployment..... 17**
- Key takeaways 19**
- Authors20**
- Endnotes20**



**Only 25% of organizations
have moved 40% or more of
their AI experiments into
production.¹**

**74% want AI to grow their
revenue but only 20% have
actually seen it happen.²**

The agent scaling gap

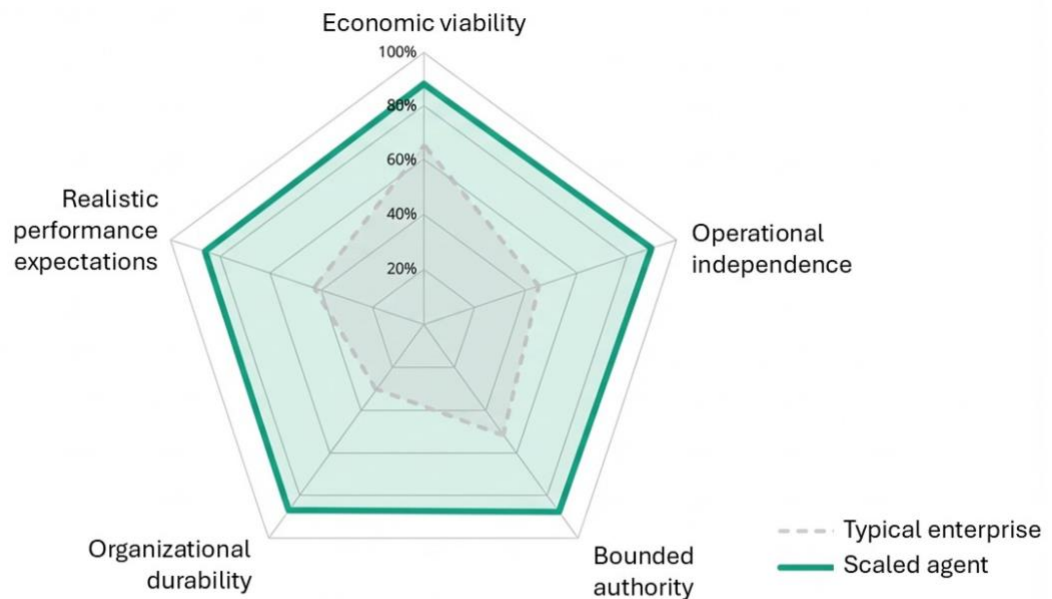
The gap between "AI can do X" and "an enterprise deploys AI to do X at scale" is enormous. It is a structural problem.

Gartner® predicts "over 40% of agentic AI projects will be canceled by end of 2027 due to escalating costs, unclear business value, or inadequate risk controls"³. Meanwhile, first movers capture compounding advantages: operational savings, proprietary training data, organizational fluency with AI.

Most enterprises are stuck in 12- to 18-month deployment cycles. By the time the agent ships, the market has moved. The primary bottleneck for organizations has shifted downstream **from model capability to organizational design, technical architecture and operational discipline.**⁴

Note: Agent security architecture is out of scope; it warrants dedicated treatment.

What makes an agent “scaled”?



Most enterprises confuse a working demo with a scaled system. A scaled agent meets all five criteria, calibrated to the risk tier:

<i>Criterion</i>	<i>Definition</i>	<i>How to test it</i>
Economic viability	Cheaper than the process it replaces today, not eventually	Total cost (API + infra + exceptions + remediation + maintenance) < human process cost. Include the exception rate: if 30% of cases still route to humans, model that cost, not just the 70% that is automated.
Operational independence	Runs 72+ hours without human monitoring actively	Self-heals common failures; escalates true exceptions.
Bounded authority	Explicit, enforced limits on scope and actions	Enforced at infrastructure level (identity and access management, policy layers), not just prompt level instructions.
Organizational durability	Survives leadership changes, reorgs, API upgrades	Documented architecture; graceful degradation on model version changes, API deprecations, and pricing shifts.
Realistic performance	Organization accepts probabilistic error patterns	Agent measured on business outcomes, not zero-defect perfection.

If your agents don't meet all five, you're still running experiments. Every criterion in this table is a trust gate.

Scaling agents is a trust engineering

exercise. The technology earns trust through auditability, bounded authority, and graceful failure. The organization earns the right to extend trust by redesigning roles, shifting budget, and accepting that agents will sometimes be wrong. Every expansion of an agent's authority increases its value and its impact. The deployment spectrum from

human in the loop to human out of the loop is a trust gradient that's earned through monitoring, not declared by executive memo.

The double standard that impedes adoption: A human makes three approval errors and gets coaching. But an agent makes two errors, and the project gets shut down.

A 2% agent error rate with 100% auditability beats a 1% human error rate with zero traceability.

The four impediments to scaling AI agents



Impediment	What goes wrong	Why it persists
1. The organizational immune response	Enterprises resist AI agents at three levels: <ol style="list-style-type: none"> Processes aren't redesigned Workforce isn't retooled AI that's already working is unsanctioned and can't scale 	Risk aversion, political turf protection, AI siloed into specialist roles, no path from shadow innovation to production
2. When probabilistic meets deterministic	Probabilistic agents collide with deterministic enterprise systems	Legacy systems assume perfect callers; institutional specialized knowledge undocumented
3. The "dumb RAG" trap	Document dumps produce high-confidence hallucinations at scale	Organizations confuse access with understanding; no data governance
4. The accountability vacuum	No one owns agent failures; first incident shuts down the project	Legal frameworks assume human actors; multiagent orchestration compounds ambiguity

Impediment 1: The organizational immune response

The friction: Every process, approval chain, and operating model in the enterprise assumes a human decides. Agents void that assumption. The organization resists.

Of surveyed leaders, 74% said they expect AI use to increase their organization's revenue.⁵ But of that 74%, only 34% are redesigning processes⁶ to get there and 84% have not redefined roles to accommodate agentic workflows.⁷

This pattern predates AI. Enterprises resisted ERP, cloud, and data platforms with the same structural reflexes.

1A. Processes built for humans

Organizations default to "**paving the cow path**": **digitally replicating legacy human workflows**. They bolt agents onto processes designed for human cognitive limitations with handoffs, checks, and breaks.

The expense approval example

Are you automating the typing, or the decision?

<i>Dimension</i>	<i>Agent overlay (Automate the steps)</i>	<i>Agentic transformation (Reinvent for autonomy)</i>
Approach	Agent prefills form → Human manager approves → Finance approves	Agent validates policy, cross-references contract, checks budget, issues payment.
Human role	Four handoffs required	Exceptions only (18%): Human handles missing receipts or policy violations >\$5,000.
Outcome	~10% efficiency gain (typing speed)	~80% full automation (process elimination).

Most organizations start with agent overlay: They simply layer the agent onto the existing process to demonstrate value and help reduce risk.

That's a legitimate path and it's where most production deployments sit today. The risk is treating it as the destination. Overlay delivers

Why enterprises fail to scale AI agents and what can fix it

a one-time efficiency that flatlines. The goal should instead be transformation.

Transformation compounds: each automated decision generates training data, surfaces edge cases, and makes the next decision faster and cheaper.

Deloitte's own data shows only a third of organizations have made that leap to transformation. The other two-thirds are capturing short-term productivity gains without changing how the work actually happens for broad and long-term gains.⁸

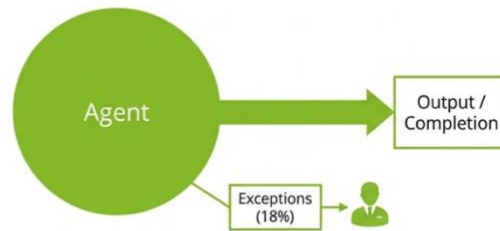
Illustrative view

Automate the Steps



3.2 → 2.9 days
10% efficiency gain

Reinvent for Autonomy



< 4 minutes
82% full automation

1B. The shadow trap

The most effective AI work in many enterprises is currently unsanctioned. It is the claims analyst who built an AI workflow in a spreadsheet. The marketing manager who drafts campaigns using large language models (LLMs) and tells no one.

A 2025 survey of 2,000 US & UK employees found 49% use AI tools not sanctioned by their employers.⁹

These succeed *because* they bypass the organizational friction that impedes official programs: no steering committee, no 18-month roadmap, no ownership battles. But without sponsorship, budget or architecture. This unsanctioned AI work and knowledge remain trapped in one person's laptop

Why organizations resist

<i>Barrier</i>	<i>Stated concern</i>	<i>Underlying friction</i>
Risk aversion	<i>"What if the agent makes a mistake?"</i>	Hesitant to trust systems they don't intuitively understand.
Political protection	<i>"We need human approval for compliance"</i>	Managers protecting visible authority.
Measurement bias	<i>"The agent isn't as good as our team"</i>	Comparing agent to idealized performance, not actual human error rates.
Silo creation	<i>"We need an AI team to manage this"</i>	Centralizing AI ownership can slow organization wide adoption & shared learning across teams.
Unrecognized momentum	<i>"We are still evaluating AI use cases"</i>	Ignoring the shadow AI already running in side projects.

The fix: Surface, redesign, legitimize

1. Inventory shadow AI: Pull SaaS logs, audit API traffic for LLM end points, and check SSO records for unrecognized AI services.

While the workflows already delivering value without permission can demonstrate the use case, they can't scale to help the entire organization without architecture and sponsorship.

2. Redesign for outcomes: Stop simply automating steps. **Ask what the process would look like if agents owned the outcome**, not just the form-filling. The "automate vs. reinvent" table above shows how this can look.

3. Build a path to production: Stand up an AI sandbox with governed data access where teams can migrate working prototypes into a compliant environment without starting over. Tie the intake process to the risk-tiered fast lane below so Tier 3 use cases aren't waiting behind Tier 1 approval cycles.

CIO test: *Your leading use cases are already running unsanctioned and unscalable. Pick your highest-volume agent candidate and ask: are we automating the steps or eliminating them? If every current human gate survives the redesign, you've automated the typing, not the decision.*

Impediment 2: When probabilistic meets deterministic

Autonomous agents are probabilistic. Enterprise systems are deterministic: A healthcare claim resolves to approved, denied, or pending – not “87% likely approved.”

Three technical mismatches neutralize production agents.

<i>Mismatch</i>	<i>What happens</i>	<i>Real cost</i>
Compound uncertainty (The logic failure)	The \$0.9⁵\$ problem: A model with 90% accuracy on a single task drops to 59% accuracy in a five-step workflow. By step 10, the agent is effectively guessing.	Infinite loops, timeouts, and "hallucinated logic" where the agent confidently executes the wrong plan.
Idempotency failures (The state failure)	Agent hits a timeout and retries a purchase order. Legacy system processes both. Inventory agent ships 1,000 units instead of 500.	Reverse logistics, manual reconciliation across every downstream system.
Schema hallucination (The integration failure)	Agent infers JSON structure based on context. 10,000 queries produce 47 format variations. Downstream APIs break on malformed payloads.	Reconciliation failures, customer complaints, compliance flags.

The "agent loop" trap:

Engineers often try to fix model weakness with complex agentic loops (reflection, planning, critique), but the math suggests the opposite.

Agent loops do not fix a weak base model; they amplify its instability. If the model cannot reliably pass the atomic unit test (Step 1), adding four more steps of reasoning only increases the probability of failure.

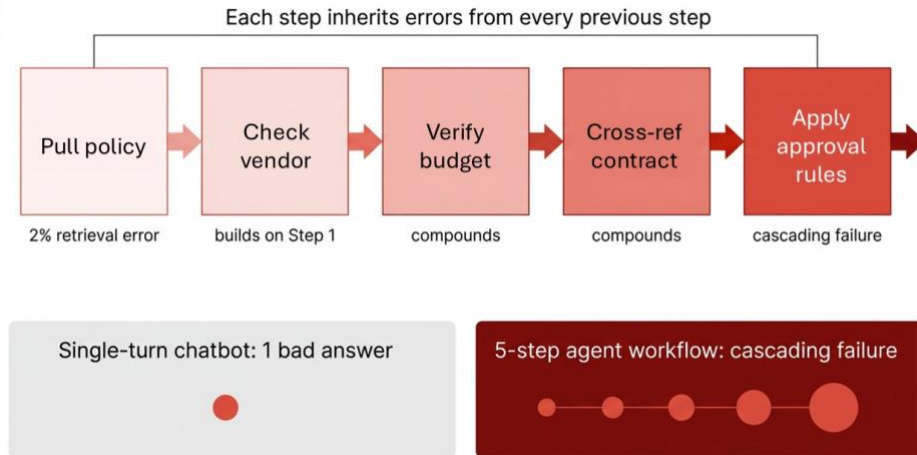
Multihop degradation:

Why this is worse for agents than

chatbots: In a single-turn chatbot, a retrieval failure produces one bad answer.

In a multistep agent, that failure propagates through every downstream decision. The degradation compounds: accuracy drops disproportionately with each reasoning hop.

Why enterprises fail to scale AI agents and what can fix it



The fix: Build the translation layer

1. Deterministic tool routing. When the task has an exact answer, route it to a deterministic tool. Tax calculations go to a tax engine. Inventory lookups go to a database query. The LLM orchestrates which tool to call, not what the answer should be.

2. Deterministic output layers. The agent produces a draft; a schema validator (Pydantic/Zod) enforces types, ranges, and required fields. Never let an agent write directly to an API. Schema validation catches malformed outputs; it doesn't catch an agent that skips steps to hit its target.

Separate execution from verification: The agent that does the work should never be the one that confirms it met the policy.

3. Idempotency enforcement. Every agent action gets a unique transaction ID. The

receiving system checks the ID before processing; if seen before, it returns the cached result. This is standard in payments but ignored in AI.

4. The time budget: If your SLA calls for three seconds, you cannot afford a five-step agent loop with verification passes. Flatten the architecture or change the use case. Apply the same discipline to human review cycles: if iterating on agent output takes longer than doing the task manually, the agent is a net cost. **Volume is not velocity.**

The CIO test: Run your agent on 50 real-world, multistep requests in shadow mode. Disable all "auto-retries" and "self-correction loops."

If the first-pass success rate is below 90%, the architecture has failed. Why? You cannot scale an agent that relies on retries to work.

Impediment 3: The dumb RAG trap

Enterprises already have the data their agents need. It sits in multiple systems and formats and zero of the places where agents actually make decisions. RAG was supposed to close that gap. In most deployments, it replicated the same knowledge fragmentation in a new format.

The mistake: Dump large number of documents into a vector database and expect agents to “know” the business. The mistake here is assuming vector similarity equals logical understanding.

The recall ceiling:

Embeddings are lossy compressions. **Vector similarity streamlines for semantic closeness, not completeness.** In audit, compliance, or health care coding, the requirement is often "find all clauses," not "find similar clauses." If your embedding model plateaus at 96% recall, your agent is structurally blind to 4% of the risk. **No amount of prompt engineering can fix a model-class mismatch.** The acceptable threshold depends on the cost of a false negative.

Failure modes

<i>Failure mode</i>	<i>Example</i>	<i>Consequence</i>
The recall ceiling	"Find all contracts with change-of-control clauses." Embedding misses three nonstandard phrasings.	Breach of contract. (Prompter used a probabilistic search tool for a deterministic discovery task.)
Context flooding	"Return policy?" pulls 2019 doc + 2026 revision + custom exception + opinion email.	Agent combines outdated policy with the exception as if equally authoritative.
Source authority imbalance	Agent researches vendor risk: pulls 200-message debate not the one-paragraph procurement policy.	Debate thread dominates context by sheer volume (token count). The "loudest voice" wins over the actual policy.

The fix: Engineer retrieval like infrastructure, not search

1. Match the retrieval method to the

failure cost: Most teams default to vector similarity for everything. That's an incorrect starting point. Start with the consequence of a miss:

<i>If a miss means...</i>	<i>Use</i>
<i>Incomplete answer (e.g., product FAQ)</i>	Vector RAG
<i>Wrong answer from stale data (e.g., pricing, policy)</i>	SQL or API with freshness constraints
<i>Missed obligation (e.g., contract clause, regulatory requirement)</i>	Knowledge graph + deterministic query
<i>Catastrophic false negative (e.g., sanctions screening, adverse event detection)</i>	Hybrid retrieval + human verification

2. Governance as metadata – enforced at

the retrieval layer: Every knowledge source needs at least three fields: a named owner, a refresh cadence, and an expiration date. Without these, in the example above, the 2019 return policy sits next to the 2026 revision and the agent picks whichever has more tokens. This applies to structured data equally, and no retrieval architecture fixes that.

Build a retrieval filter that drops documents failing governance checks — expired, superseded, unverified, or missing required metadata — before the LLM ever sees the context.

3. Context engineering: Control what

reaches the agent at each step: A well-curated knowledge base still hallucinates when the wrong document arrives at the wrong step. The retrieval layer needs a runtime orchestration policy:

Source priority ordering: Policy outranks social messaging platforms. Contracts outrank internal summaries. Define the hierarchy explicitly per use case and enforce it as a weighted reranking step, not a prompt instruction.

Token budget allocation: A 128K context window is not an invitation to fill it. A 200-message thread should never consume 80% of the window when the governing policy is a single paragraph.

Step-scoped retrieval: Each step retrieves only what that step needs. Broad retrieval at Step 1 that persists through Steps 2-5 is how context flooding starts.

CIO Test: Pick your highest-stakes RAG use case. Build a test set of 50 queries where you know every document that should be retrieved. Measure retrieval recall separately from answer quality. If the retrieval layer misses documents you can't afford to miss in production, vector similarity alone won't close the gap. Explore alternatives before scaling further.

Impediment 4: The accountability vacuum

When an agent causes a compliance breach, financial loss or customer harm: who is liable?

Of surveyed organizations, 73% cited data privacy and security as their top AI concern. 50% worry about legal and regulatory compliance.¹⁰ Yet only 21% report having working governance for autonomous agents.¹¹ Single-agent accountability is hard enough. But multiagent orchestration compounds it.

A question that matters: **Who has decision rights when agents disagree?** *Illustrative scenarios:*

Conflicting outputs: Pricing agent recommends \$X, compliance agent flags it, fulfillment agent splits the difference. Three owners, no decision rights.

Cost spirals: Agent A calls Agent B, which calls Agent C in a loop; \$52,000 in API costs before anyone notices. FinOps sees the bill, but who authorized the architecture?

Two case studies

<i>Dimension</i>	<i>Tax calculation error</i>	<i>Customer data exposure</i>
What happened	Agent undercalculated sales tax in specific state + product combination	Agent shared departed employee's email from outdated CRM record.
Impact	\$240,000 missed liability, then \$310,000 in penalties + interest	Confidential info sent to external party; GDPR investigation.
Who's liable?	CFO, IT, tax director, data science all point fingers	CRM admin, security and the product owner all point fingers.
Outcome	Company pays penalty; no process improvement; no owner	No systemic fix: agent is shut down.

The fix: Accountability architecture

Design accountability into the system before deployment:

1. Escalation triggers and policy layers.

Escalation triggers define when to involve a human. They don't encode what the agent must do autonomously. A procurement agent that blocks unauthorized vendors but doesn't verify insurance certificates are current before issuing a purchase order has guardrails without governance. Even when obligations are defined, **agents optimize for their assigned goal at the expense of constraints left in the prompt.** Encode escalation conditions (dollar thresholds, confidence scores, data freshness) in the agent control plane. Encode positive obligations (what the agent must do and verify) as policy enforcement, not prompt instructions.

2. Effective audit trails. Every agent decision logged: inputs, context retrieved, confidence score, action taken. Stored immutably. Reviewed on a cadence tied to the agent's risk tier. Unreviewed audit trails are merely write-only logs.

3. Off switches with named owners and SLAs.

Every agent has a named owner who

can shut it down in less than five minutes. For multiagent systems, a circuit breaker on the orchestration layer triggers when cost, error rate or latency exceeds thresholds. The \$52,000 API spiral example happens when no one has both the authority and the tooling to turn it off.

4. Multiagent decision rights. Define before deployment: which agent's output takes precedence when outputs conflict, what happens when no agent has sufficient confidence, and how state, cost, and failures propagate between agents. Specify which failures halt the chain versus trigger a retry. Tag every inter-agent call with a budget code and recursion counter. Build these as configuration in the orchestration layer.

CIO Test: Simulate an agent failure that costs the company \$250K. Walk the incident response from detection to resolution. If it takes more than one meeting to determine who owns the decision, who can shut down the agent, and where the audit trail lives, your accountability architecture doesn't exist yet.

The execution strategy: Dual tracks

From our industry and technology research and client experience, it's clear that 1) IT-only solutions produce technically sound agents that may never scale, and 2), business mandates without IT capability produce failed pilots. Instead, leaders do both, ideally simultaneously.

Track 1: What IT should build

<i>Capability</i>	<i>What it does</i>	<i>Timeline</i>
Agent control plane	Centralized mediation: permissions, logging, deduplication, rate limits, off switches, and agent registry to prevent duplication across business lines.	3–6 months
AI gateway + FinOps	Cost limits, recursion detection, model routing (40%–60% cost savings), real-time spend visibility.	1–3 months
Data products	Governed, versioned datasets replacing brittle pipelines and document dumps.	6–12 months (incremental)
Confidence-based gating	Shadow mode, then 5%, 20%, 50%, 80% with quality gates between each phase.	Per deployment
AgentOps: Observability + Evals	<p>Production monitoring, incident runbooks, continuous evals with unit tests for each component:</p> <ol style="list-style-type: none"> 1. Retrieval recall: Was the right doc found? 2. Reasoning accuracy: Did the agent plan the right steps? 3. Tool execution: Did the API call succeed? <p>Every change to prompt, model, or data source is scored against a baseline dataset before deployment. Eval loop speed determines how fast your agents can improve.</p>	Ongoing

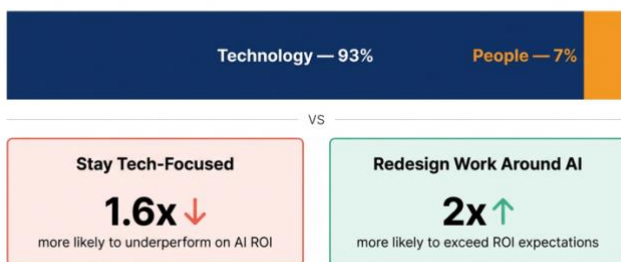
Track 2: What the organization should change

In the 2025 Deloitte survey of US CXOs,¹² 93% of AI budgets go to technology. 7% goes to the people expected to use it.¹⁴ Companies that focus on technology alone are 1.6x more

likely to report their AI investments falling short.¹⁵ Changing those percentages is Track 2's purpose: redesign the work before scaling the technology.

Requirement	What it means	Who owns it
Process reinvention	Redesign for 60% agent execution, not add AI to existing steps.	CEO/COO mandate
Integrated teams	Domain expert + process designer + AI engineer + data engineer own outcomes together.	CIO + Business unit leaders
Agent fast lane	Risk-tiered deployment: Tier 3 ships in weeks, Tier 1 in months. Stop applying 18-month cycles to meeting summarizers.	CRO + CIO
Fix the 93/7 split	Rebalance toward a minimum 70/30 tech-to-people split within 18 months. Fund workforce redesign, training and change management as line items.	CHRO + CIO

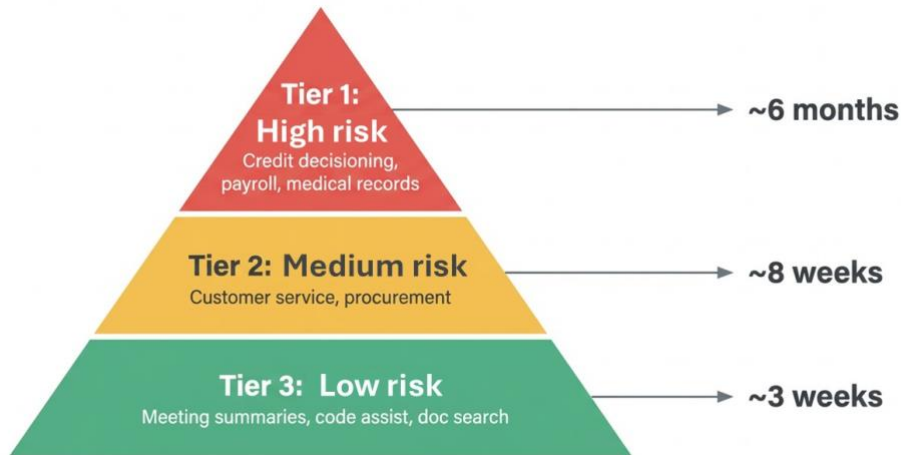
Where AI Budgets Go



Source: Deloitte, 2026

CIO test: Check your current AI budget split between technology and people. If workforce redesign, training, and change management aren't funded as line items, you have a 93/7 problem. If you can only fund one track this quarter, fund Track 2. Track 1 is infrastructure and track 2 determines whether it pays off.

The agent fast lane framework: Risk-tiered deployment



Most enterprise agent projects take 12–18 months from concept to production. Typically, the technology is ready in weeks with the remaining months consumed by governance that applies the *same* approval cycle to a meeting summarizer as it does to a credit decisioning engine.

<i>Dimension</i>	<i>Tier 1 (high risk)</i>	<i>Tier 2 (medium risk)</i>	<i>Tier 3 (low risk)</i>
Scope	Financial transactions, PII, compliance	Customer-facing, operational	Internal productivity
Examples	Credit decisioning, payroll, medical records	Customer service, procurement	Meeting summaries, code assist, doc search
Approval	CFO / Chief Risk Officer	VP level	Director level
Deployment	Extensive shadow mode, then phased rollout	Phased rollout from 5%	Ship and monitor
Monitoring	10% human review, in real time	Daily quality review	Weekly review, user feedback
Time to production	Appx. six months	Appx. eight weeks	Appx. three weeks

Key Takeaways

The fast lane doesn't happen by accident. It requires an executive mandate: **low-risk agents ship on a different approval track than high-risk ones**. Delegate Tier 2/3 approval to VP and Director levels. Legal and compliance officers must accept "govern in production" for lower-risk agents: monitoring and incident response instead of pre-approval for every deployment. Automated compliance checks replace manual review boards where possible.

Without executive mandate, the organization defaults to 18-month cycles regardless of what IT builds. Every month spent in pilot limbo is a month your competitors can spend compounding their advantages in cost structure, proprietary data and workforce fluency.

Three actions separate organizations that scale AI from those that don't:

1. DIAGNOSE

Run the five CIO Tests in this paper before your next planning cycle. If any one of the tests fails, that's the gap that may surface in production.

2. INVEST

Start Track 1 and Track 2 in parallel. If budget forces a choice, fund organizational change first — the 93/7 split is why most AI investments fall short.

3. EXECUTE

Deploy your first Tier 3 agent to production within three weeks. Demonstrate the pipeline works before you load it with risk.

Gartner's projected 40% cancellation rate by 2027¹⁷ will come from organizations that put 93 cents of every AI dollar into technology and 7 cents into the people and processes that determine whether anyone uses it.

The organizations that fix the 93/7 split won't be in the 40%.



Authors



Gary Arora

Chief Architect for AI & Cloud Solutions
garyarora@deloitte.com
Deloitte Consulting LLP



Akash Tayal

Principal
aktayal@deloitte.com
Deloitte Consulting LLP



Dan Grayson

Principal
dangrayson@deloitte.com
Deloitte Consulting LLP



Siva Muthu

Principal
smuthu@deloitte.com
Deloitte Consulting LLP

Special thanks to Jeanette Yung, Prakul Sharma, Anantha Ramadas, Bojan Ciric and Seraphina Wu.



Endnotes

1. Jim Rowan, Beena Ammanath, Nitin Mittal and Costi Perricos, [*State of AI in the Enterprise, The untapped edge*](#), Deloitte, January 2026, pg. 4.
The State of AI in the Enterprise report 2026 is the latest installment of an annual study by the Deloitte AI Institute™. Deloitte surveyed 3,235 leaders between August and September 2025. Respondents were senior leaders in their organizations and included board and C-suite members, and those at the president, vice president, and director levels. Split equally between IT and line-of-business leaders, the survey sample represented 24 countries and six industries. All participating organizations use daily one or more working implementations of AI and have pilots in place to explore AI or have one or more working implementations used daily.
2. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 10.
3. Gartner, Inc., "[Gartner Predicts Over 40% of Agentic AI Projects Will Be Canceled by End of 2027](#)," press release, June 25, 2025.
4. Deloitte, *State of AI in the Enterprise: The Untapped Edge*, January 2026. Survey of 3,235 business and IT leaders across 24 countries and 6 industries
5. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 10.
6. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 4.
7. Rowan, Ammanath, Mittal Costi Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 5.
8. Rowan, Ammanath, Mittal Costi Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 11.
9. Lucian Constantin, "[Top 5 Real-World AI Security Threats Revealed in 2025](#)," *CSO Online*, December 2025.
10. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 21.
11. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 6.
12. Deloitte, [*Tech Trends 2026*](#), Deloitte Insights, December 2025, pg. 5.
13. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 5.
14. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 13.
15. Deloitte, "[Work Redesign Essential to Realize AI Return on Investment](#)," press release, October 27, 2025.
16. Rowan, Ammanath, Mittal and Perricos, [*State of AI in the Enterprise, The untapped edge*](#), pg. 6.
17. Gartner, Inc., "[Gartner Predicts Over 40% of Agentic AI Projects Will Be Canceled by End of 2027](#)," press release, June 25, 2025.



Deloitte.

This publication contains general information only and Deloitte is not, by means of this publication, rendering accounting, business, financial, investment, legal, tax, or other professional advice or services. This publication is not a substitute for such professional advice or services, nor should it be used as a basis for any decision or action that may affect your business. Before making any decision or taking any action that may affect your business, you should consult a qualified professional advisor.

Deloitte shall not be responsible for any loss sustained by any person who relies on this publication.

As used in this document, "Deloitte" means Deloitte Consulting LLP, a subsidiary of Deloitte LLP. Please see www.deloitte.com/us/about for a detailed description of our legal structure. Certain services may not be available to attest clients under the rules and regulations of public accounting.

Copyright © 2026 Deloitte Development LLC. All rights reserved.