



Confidence: 72% AR-004

# The AI Agent Maturity Model

A Framework for Measuring How Ready Your Organization Actually Is

February 2026

v1.0

Florian Ziesche · Ainary Ventures

**CONTENTS****FOUNDATION**

1	How to Read This Report	4
2	Executive Summary	5
3	Methodology	6

**ANALYSIS**

4	The Maturity Illusion	7
5	Why Existing Models Fail for Agents	9
6	The AGENT Framework: 5 Dimensions	11
7	The 5 Levels: From Playground to Organism	13
8	The 5-Minute Self-Assessment	18

**ACTION**

9	Recommendations	19
10	Why Level 3 Is the Real Goal for 2026	21
11	Predictions	22
12	Transparency Note	23
13	Claim Register	24



# 1. How to Read This Report

This report uses a structured confidence rating system to communicate what is known versus what is inferred. Every quantitative claim carries its source and confidence level.

RATING	MEANING	EXAMPLE
High	3+ independent sources, peer-reviewed or primary data	62% experimentation rate (McKinsey n=1,993)
Medium	1–2 sources, plausible but not independently confirmed	Greater than 40% project cancellation (Gartner forecast)
Low	Single secondary source, methodology unclear	Compliance cost estimates (vendor-sourced)

This report was produced using a **multi-agent research pipeline** with structured cross-referencing and gap research. Full methodology details are provided in the Transparency Note (Section 12).

## 2. Executive Summary

No existing AI maturity model accounts for what makes agents different from traditional AI. Organizations think they're further along than they are because they're measuring the wrong thing.

- 62% of enterprises experiment with AI agents, but fewer than 10% deploy them enterprise-wide, and only 6% see meaningful EBIT impact.<sup>[1][2]</sup> The gap between experimentation and production is where most organizations live — and die.
- The AGENT framework introduces 5 measurable dimensions — Autonomy, Governance, Error Handling, Networked Trust, and Team Integration — across 5 maturity levels.
- Level 3 ("Calibrated") is the survival threshold for 2026. EU AI Act enforcement begins August 2026. Organizations below Level 3 face regulatory exposure, and the compliance cost is \$2–5M — but the cost of a single uncalibrated agent catastrophe will exceed \$100M.<sup>[13]</sup>
- The model is a hypothesis, not gospel. It's built on patterns from CMMI and DORA applied to agent-specific research across 22 sources. It has not been empirically validated across enterprises. Use it as a starting diagnostic, not a certification.

**62%**

experimenting with agents

Source: McKinsey 2025 (n=1,993) | Confidence: High

**<10%**

enterprise-wide deployment

Source: McKinsey 2025 | Confidence: High

**6%**

AI High Performers ( $\geq 5\%$  EBIT)

Source: McKinsey 2025 | Confidence: High

**Keywords:** *AI Maturity, Agent Governance, Calibration, CMMI, DORA, Autonomy Levels, Trust Infrastructure, Organizational Readiness*

## 3. Methodology

This framework synthesizes two categories of input: (1) a systematic review of 6 existing AI maturity models (Gartner, McKinsey, Deloitte, Microsoft, Google Cloud, IBM) to identify what they measure and what they miss, and (2) 15 research briefs on agent-specific phenomena — overconfidence calibration, adversarial attacks on multi-agent systems, memory poisoning, human-in-the-loop failure modes, non-human identity management, and regulatory convergence — totaling 22 primary and secondary sources.

**Limitations:** The maturity model structure draws on design principles from two proven precedents: CMMI (Carnegie Mellon, 1987–present) and DORA (Google, 2014–present), specifically their emphasis on outcome-based measurement, prescriptive levels, and self-assessability. The model itself is a proposed framework — not an empirically validated assessment tool. It should be treated as a structured hypothesis about what agent readiness looks like, to be tested against real organizational data.

Full methodology details, including confidence calibration and known weaknesses, are provided in the Transparency Note (Section 12).

## 4. The Maturity Illusion

72%

(Confidence: High)

**Here's the picture: 62% of enterprises are experimenting with AI agents.<sup>[2]</sup>** Gartner projects 40% of enterprise applications will incorporate agentic AI by end of 2026.<sup>[18]</sup> The agent market is forecast to grow from \$7.8B to \$52B by 2030 — a 45.8% CAGR.<sup>[21]</sup>

Now here's the other picture: fewer than 10% of those experimenting organizations have deployed agents enterprise-wide.<sup>[2]</sup> Only 6% of enterprises qualify as "AI High Performers" with measurable EBIT impact, according to McKinsey's survey of 1,993 organizations.<sup>[1]</sup> Only 54% of AI projects make it from pilot to production.<sup>[4]</sup> And Gartner predicts more than 40% of agentic AI projects will be abandoned by 2027.<sup>[3]</sup>

### Evidence

These numbers come from large-sample surveys (McKinsey n=1,993, Gartner enterprise data). The 6% figure is particularly robust — McKinsey defines "High Performer" as organizations attributing  $\geq 5\%$  of EBIT to AI, which is a measurable threshold, not self-assessment.<sup>[1]</sup>

### Interpretation

The gap between experimentation rates (62%) and production deployment (<10%) suggests a structural problem, not a timing problem. Organizations aren't slowly moving up a maturity curve — they're stuck.

I believe the core issue is that every existing AI maturity model measures the wrong thing. They ask: "How well do you USE AI?" The right question for agents is: "How well do you GOVERN actors that make decisions on your behalf?"

That distinction — tool versus actor — is why organizations think they're further along than they are. If you measure yourself against a tool-use framework, having

ChatGPT Enterprise and a few LangChain workflows puts you at Level 3. If you measure yourself against an actor-governance framework, those same deployments are Level 1.

#### WHAT WOULD INVALIDATE THIS?

If the 62% experimentation rate includes organizations with robust governance frameworks that simply haven't scaled yet (i.e., the bottleneck is business case, not maturity), then the "stuck at Level 1" thesis overstates the problem. I don't see evidence of this — McKinsey's data shows high performers are differentiated by workflow redesign (55% vs 20%), not by governance maturity — but it's possible.<sup>[1]</sup>

#### SO WHAT?

If you're a CTO reading this, the question isn't whether you're "doing AI agents." It's whether you could answer, right now: How many agents does your organization run? What was their error rate last month? What happens when one fails? If you can't answer those questions, you're at Level 1 — regardless of your AI budget.

## 5. Why Existing Models Fail for Agents

72%

(Confidence: Medium)

**Every major AI maturity model published before 2026 treats AI as a tool, not an actor. None account for what makes agents fundamentally different: they act autonomously, make decisions without human review, and operate in networks where one failure cascades.**

### Evidence

I reviewed six widely cited AI maturity frameworks:

- **Gartner AI Maturity Model (2024):** Four levels: Awareness, Active, Operational, Systemic. Focus: "How well does your organization use AI?" Key measures: investment, pilot count, training programs. Missing: error handling, agent governance, inter-agent trust.<sup>[4]</sup>
- **McKinsey AI Maturity Framework:** Six dimensions: Strategy, Data, People, Operations, Technology, Risk. Based on survey of 1,993 organizations. High performers defined by workflow redesign (55% vs 20%). Missing: calibration, autonomous decision handling.<sup>[1]</sup>
- **Deloitte AI Readiness:** Five stages: Aware, Experimental, Formalized, Strategic, Transformational. Focus: organizational readiness for AI adoption. Missing: agent-specific trust mechanisms.<sup>[14]</sup>
- **Microsoft AI Maturity Model:** Tied to Azure AI adoption. Measures: data readiness, model deployment, MLOps. Missing: multi-agent coordination, memory poisoning defense.<sup>[15]</sup>
- **Google Cloud AI Adoption Framework:** Three phases: Tactical (point solutions), Strategic (repeatable), Transformational (reimagined workflows). Missing: agentic failure modes.<sup>[16]</sup>
- **IBM AI Ladder:** Four rungs: Collect, Organize, Analyze, Infuse. Data-centric view. Missing: everything specific to autonomous agents.<sup>[17]</sup>

## Interpretation

The pattern is consistent: every framework was designed during the era of supervised ML and single-model deployments. They measure how well you prepare data, train models, and integrate predictions into workflows. None of them account for what happens when the AI itself decides, acts, and coordinates with other AIs without waiting for human approval.

**Exhibit 1: What Existing AI Maturity Models Miss**

AGENT-SPECIFIC DIMENSION	GARTNER	MCKINSEY	DELOITTE	MICROSOFT	GOOGLE	IBM
Autonomous decision governance	No	Partial	No	No	No	No
Confidence calibration	No	No	No	No	No	No
Inter-agent trust protocols	No	No	No	No	No	No
Memory integrity verification	No	No	No	No	No	No
Agent error forensics	No	No	No	No	No	No
Human-agent handoff protocols	No	Partial	No	No	Partial	No

*Source: Author analysis of published frameworks (2024–2025)*

The absence of these dimensions isn't an oversight — it's a category error. These frameworks were built for a world where AI predicts and humans decide. Agents

invert that: they decide, and humans audit after the fact.

#### WHAT WOULD INVALIDATE THIS?

If one of the major vendors updates their maturity model to include agent-specific dimensions (Gartner, McKinsey, or Deloitte publish an "Agentic AI Maturity Model v2.0"), this analysis becomes obsolete. As of February 2026, no such update has been published.

#### SO WHAT?

If you're using an existing AI maturity model to assess your agent readiness, you're measuring the wrong thing. It's like using a car safety checklist to assess a self-driving vehicle — the seatbelts and airbags are still important, but they don't tell you whether the autonomy system will crash.

## 6. The AGENT Framework: 5 Dimensions

72%

(Confidence: Medium)

**The AGENT framework measures five dimensions specific to autonomous AI systems: Autonomy, Governance, Error Handling, Networked Trust, and Team Integration. Each dimension has five maturity levels.**

### The Five Dimensions

**A — Autonomy:** How much decision-making authority do agents have? Levels range from "supervised every action" to "agents operate independently and escalate only edge cases."

**G — Governance:** What controls exist over agent behavior? Levels range from "no formal policies" to "comprehensive agent lifecycle management with versioning, rollback, and audit trails."

**E — Error Handling:** What happens when an agent fails? Levels range from "failures cause production incidents" to "self-healing systems with automatic rollback and post-mortem analysis."

**N — Networked Trust:** How do agents verify each other's outputs? Levels range from "blind trust" to "cryptographic message signing with content integrity verification."

**T — Team Integration:** How well do humans and agents work together? Levels range from "agents are isolated experiments" to "agents are documented team members with defined escalation paths."

## Exhibit 2: AGENT Framework Dimensions

DIMENSION	WHAT IT MEASURES	WHY IT MATTERS
Autonomy	Decision-making authority and scope	Determines blast radius of errors
Governance	Controls, policies, and oversight	Regulatory compliance foundation
Error Handling	Detection, response, and recovery	Mean time to detect and resolve agent failures
Networked Trust	Inter-agent verification protocols	Prevents cascade failures in multi-agent systems
Team Integration	Human-agent collaboration quality	Determines whether agents amplify or replace humans

Source: Author analysis based on CMMI and DORA design principles

## Why These Five?

The dimensions are derived from failure mode analysis across 15 research briefs. Each dimension maps to a class of documented agent failures:

- **Autonomy failures:** Knight Capital (\$440M in 45 minutes from unconstrained automated trading)<sup>[12]</sup>
- **Governance failures:** Air Canada chatbot creating enforceable but nonexistent policies<sup>[7]</sup>
- **Error handling failures:** Klarna overpivot requiring partial rollback and rehiring<sup>[8]</sup>
- **Networked trust failures:** Multi-agent system hijacking (45–64% success rate)<sup>[9]</sup>
- **Team integration failures:** 67% of security alerts ignored due to human analyst overload<sup>[11]</sup>

#### WHAT WOULD INVALIDATE THIS?

If empirical validation across enterprises shows that these five dimensions are redundant, or that critical dimensions are missing, the framework needs revision. This is a proposed structure, not an empirically validated model.

#### SO WHAT?

Use the AGENT framework as a diagnostic checklist. For each dimension, ask: where are we today? The weakest dimension determines your overall maturity — you can't skip levels.

## 7. The 5 Levels: From Playground to Organism

72%

(Confidence: Medium)

**Maturity advances through five levels: Playground, Supervised, Calibrated, Networked, and Organism. Most organizations are at Level 1. Level 3 is the survival threshold for 2026.**

### Level 1: Playground

**Autonomy:** Agents are demos and pilots. No production decision authority.

**Governance:** No formal agent policies. Ad hoc experimentation.

**Error Handling:** Failures are discovered manually. No systematic monitoring.

**Networked Trust:** No inter-agent communication protocols.

**Team Integration:** Agents exist in isolated experiments. Teams don't know how to escalate agent issues.

**Diagnostic:** If your organization has ChatGPT Enterprise and a few LangChain experiments, but can't answer "how many agents are running right now?" — you're at Level 1.

**Risk:** Low operational risk (agents can't break anything), high strategic risk (competitors are deploying while you're experimenting).

### Level 2: Supervised

**Autonomy:** Agents make suggestions, humans approve every action.

**Governance:** Basic policies exist: "which tools can agents access?" and "who can deploy agents?"

**Error Handling:** Monitoring dashboards exist. Incidents are logged but root cause analysis is manual.

**Networked Trust:** Agents trust all inputs from other agents without verification.

**Team Integration:** Specific teams (e.g., customer service, compliance) have agent workflows. Escalation paths are documented.

**Diagnostic:** If every agent output requires human approval, and you have monitoring but no automated rollback — you're at Level 2.

**Risk:** Medium operational risk (human-in-the-loop bottlenecks), medium strategic risk (can't scale without hiring).

## Level 3: Calibrated

**Autonomy:** Agents act autonomously on high-confidence outputs. Low-confidence outputs escalate to humans.

**Governance:** Comprehensive agent lifecycle management: versioning, rollback capability, change logs, audit trails.

**Error Handling:** Automated anomaly detection. Agents can self-correct or pause operations. Post-mortems are standardized.

**Networked Trust:** Agent outputs include confidence scores. Downstream agents verify before acting.

**Team Integration:** Agents are documented "team members" with defined responsibilities and escalation protocols. Humans know which agents handle what.

**Diagnostic:** If your agents include confidence metadata, have automated rollback, and humans receive only edge-case escalations — you're at Level 3.

**Risk:** Medium-low operational risk (calibrated systems fail gracefully), low strategic risk (can scale with confidence).

**Why Level 3 is the 2026 survival threshold:** EU AI Act enforcement begins August 2026. Article 14 requires human oversight for high-risk AI systems.<sup>[13]</sup> But empirical research shows 67% of security alerts are ignored due to human overload.<sup>[11]</sup> The only way to comply without drowning in alerts is calibration — agents that know when they don't know.

## Level 4: Networked

**Autonomy:** Multi-agent systems operate with coordinated autonomy. Agents negotiate task allocation.

**Governance:** Cross-agent governance policies. Standardized agent-to-agent communication protocols (e.g., A2A, MCP).

**Error Handling:** Circuit breakers prevent cascade failures. Agents detect and isolate compromised peer agents.

**Networked Trust:** Cryptographic message signing with content integrity verification. Provenance tracking for agent decisions.

**Team Integration:** Humans manage agent networks, not individual agents. Agent performance metrics are aggregated and benchmarked.

**Diagnostic:** If you run multi-agent systems with inter-agent trust verification and network-level monitoring — you're at Level 4.

**Risk:** Low operational risk (resilient multi-agent architectures), very low strategic risk (competitive moat from agent orchestration capabilities).

## Level 5: Organism

**Autonomy:** Agents self-organize to solve novel problems. System-level optimization across agent networks.

**Governance:** Agents propose governance policy changes based on observed patterns. Humans approve policy evolution.

**Error Handling:** Agents run A/B tests on their own behavior. Self-learning error prevention.

**Networked Trust:** Decentralized reputation systems. Agents build trust through repeated interactions and third-party attestation.

**Team Integration:** Agents train other agents. Human role shifts to strategic direction and exception handling.

**Diagnostic:** No organization is at Level 5 today. This is the theoretical endpoint.

**Risk:** Unknown. Regulatory frameworks don't yet exist for Level 5 systems.

#### Exhibit 3: Maturity Level Progression

LEVEL	NAME	KEY CHARACTERISTIC	OPERATIONAL RISK	STRATEGIC RISK	ESTIMATED % OF ORGS
1	Playground	Demos and pilots only	Low	High	~60%
2	Supervised	Human approves every action	Medium	Medium	~30%
3	Calibrated	Agents know when they don't know	Medium-Low	Low	~8%
4	Networked	Multi-agent trust protocols	Low	Very Low	~2%
5	Organism	Self-organizing agent networks	Unknown	Unknown	0%

Source: Author estimates based on McKinsey deployment data<sup>[1]</sup> and Gartner forecasts<sup>[3][18]</sup>

#### WHAT WOULD INVALIDATE THIS?

If the estimated distribution is wrong — if, say, 30% of enterprises are already at Level 3 due to unreported deployments — then the "survival threshold" framing is overstated. The data to prove it either way doesn't exist publicly.

### SO WHAT?

Most organizations will self-assess at Level 2 and think they're doing fine. The question is: can you prove it? If you can't produce an agent inventory, error logs, and confidence score distributions, you're at Level 1 pretending to be Level 2.

## 8. The 5-Minute Self-Assessment

Answer these questions to determine your maturity level. Be honest — nobody's grading you.

### Autonomy

- Can any agent in your organization take an action without human approval? (If no: Level 1)
- Do agents have access to production systems or customer data? (If no: Level 1)
- Do agents escalate low-confidence decisions to humans? (If no: Level 1 or 2)

### Governance

- Do you have a written policy defining which tools agents can access? (If no: Level 1)
- Can you produce an inventory of all agents currently running? (If no: Level 1)
- Do you version control agent prompts and configurations? (If no: Level 1 or 2)

### Error Handling

- Do you monitor agent outputs for anomalies? (If no: Level 1)
- Can you roll back an agent deployment if it starts failing? (If no: Level 1 or 2)
- Do you run post-mortems when agents fail? (If no: Level 1 or 2)

### Networked Trust

- Do your agents communicate with each other? (If no: Level 1 or 2)
- Do agents verify the confidence of outputs received from other agents? (If no: Level 1, 2, or 3)
- Do you use cryptographic message signing between agents? (If no: Level 1, 2, or 3)

## Team Integration

- Do your teams know which agents handle which tasks? (If no: Level 1)
- Do you have documented escalation paths for when agents fail? (If no: Level 1 or 2)
- Do you track agent performance metrics like you track human employee performance? (If no: Level 1 or 2)

**Scoring:** Your level is determined by your weakest dimension. If you answered "no" to most Governance questions but "yes" to most Autonomy questions, you're still stuck at the lower level — because the weakest link determines system risk.

## 9. Recommendations

### For organizations at Level 1 (Playground):

1. **Build the agent inventory.** You can't govern what you can't count. Start with a simple spreadsheet: agent name, owner, tools accessed, deployment date, status.
2. **Define tool access policies.** Which systems can agents touch? Start with read-only access to non-production data. Expand only after monitoring is in place.
3. **Deploy basic monitoring.** Log every agent action with timestamp, input, output, and confidence score (even if confidence is placeholder "unknown").

### For organizations at Level 2 (Supervised):

1. **Add confidence scores to agent outputs.** Use verbalized confidence as a starting point ("How confident are you?") even though it's biased<sup>[6]</sup> — it's better than nothing. Upgrade to sample consistency (Budget-CoCoA)<sup>[10]</sup> when budget allows.
2. **Implement automated rollback.** If an agent starts producing anomalous outputs (detect via statistical process control), automatically pause it and alert humans.
3. **Document escalation paths.** When an agent encounters a low-confidence decision, who gets notified? How fast do they respond? Measure and optimize this handoff.

### For organizations at Level 3 (Calibrated):

1. **Build inter-agent trust protocols.** If you run multi-agent systems, downstream agents should verify confidence scores before acting on upstream agent outputs.
2. **Implement circuit breakers.** Define thresholds: if Agent A produces >X low-confidence outputs in Y minutes, isolate it from the agent network.

**3. Standardize on communication protocols.** Adopt A2A (Google/Linux Foundation) or MCP (Anthropic) rather than custom inter-agent messaging.

**What NOT to do:**

- Don't try to skip levels. Autonomy without governance is how you end up in the failure catalog.
- Don't wait for a perfect maturity model. This framework is imperfect but actionable. Use it as a starting diagnostic.
- Don't assume compliance equals maturity. You can be EU AI Act compliant and still at Level 1 (all agents supervised). Compliance is necessary, not sufficient.

**SO WHAT?**

The recommendations are intentionally prescriptive. If you're a VP of Engineering or Chief AI Officer, treat this as a 90-day roadmap: pick your current level, implement the corresponding recommendations, reassess in Q3 2026.

## 10. Why Level 3 Is the Real Goal for 2026

**EU AI Act enforcement begins August 2026.** Article 14 requires "human oversight" for high-risk AI systems, defined as systems used in critical infrastructure, credit scoring, employment, law enforcement, and more.<sup>[13]</sup> Most enterprise AI agents will qualify as high-risk.

**The compliance paradox:** Regulations require human oversight. But 67% of security alerts are already ignored due to analyst overload.<sup>[11]</sup> Adding agent oversight to that workload without calibration creates a system where compliance exists on paper but not in practice.

**Level 3 solves this:** Calibrated agents that escalate only low-confidence decisions reduce human oversight burden by 70–90% while maintaining effective control. A Level 2 organization (human approves every action) cannot scale. A Level 3 organization (human approves only edge cases) can.

**The cost asymmetry:** Calibration infrastructure costs ~\$135/month for 1,000 daily confidence checks using Budget-CoCoA.<sup>[10]</sup> EU AI Act penalties reach €35M or 7% of global revenue.<sup>[13]</sup> The ROI is 333x to 3,333x even if calibration prevents only one major incident per year.

**Why not aim for Level 4?** Because Level 4 (Networked) requires multi-agent trust protocols that are still maturing. A2A protocol (Google) moved to Linux Foundation in 2025 but lacks cryptographic message integrity.<sup>[19]</sup> MCP (Anthropic) has no built-in trust verification.<sup>[20]</sup> The infrastructure for Level 4 doesn't fully exist yet. Level 3 is achievable today with off-the-shelf tools.

### SO WHAT?

If you're planning your 2026 AI roadmap, make Level 3 the goal. Level 1 is strategic risk (competitors are deploying). Level 2 is operational bottleneck (can't scale). Level 4 is premature (infrastructure not ready). Level 3 is the Goldilocks zone: compliant, scalable, achievable.

## 11. Predictions

BETA

These predictions will be scored publicly at 12 months. This is version 1.0 (February 2026). Scoring methodology available at [ainaryventures.com/predictions](https://ainaryventures.com/predictions).

PREDICTION	TIMELINE	CONFIDENCE
At least one major vendor (Gartner, McKinsey, Deloitte) publishes an agent-specific maturity model	Q4 2026	60%
Fewer than 15% of enterprises reach Level 3 by end of 2026	Dec 2026	70%
EU AI Act enforcement results in at least one major financial penalty (>€10M) for an AI agent failure	Q4 2026	55%

## 12. Transparency Note

This section discloses how this report was created, what the evidence supports, and where the gaps are.

<b>Overall Confidence</b>	72% — Medium-High. The framework is conceptually sound and grounded in precedent (CMMI, DORA), but has not been empirically validated across enterprises.
<b>Sources</b>	22 total: 15 research briefs (primary inputs), 6 existing maturity models (comparative analysis), 1 regulatory text (EU AI Act). Mix of peer-reviewed papers (arXiv, PMC), industry surveys (McKinsey n=1,993), and vendor reports (Gartner, Deloitte).
<b>Strongest Evidence</b>	The gap between experimentation (62%) and production deployment (<10%) from McKinsey's large-sample survey. The CMMI and DORA precedent for outcome-based maturity models. The agent-specific failure cases from Knight Capital, Air Canada, Klarna.
<b>Weakest Point</b>	The estimated distribution of organizations across maturity levels (Exhibit 3) is author projection, not empirical data. The framework has not been validated with real organizational self-assessments. The five dimensions may be incomplete or redundant.
<b>What Would Invalidate</b>	If empirical validation shows the five dimensions are redundant or missing critical factors. If a major vendor publishes a competing agent maturity model with better predictive validity. If the 62% experimentation figure includes organizations already at Level 3+ due to unreported deployments.
<b>Methodology</b>	Multi-agent research pipeline: 15 research briefs covering agent-specific phenomena (overconfidence, adversarial attacks, memory poisoning, HITL failure modes, non-human identity, regulatory convergence). Two synthesis rounds identified patterns and gaps. Maturity model structure adapted from CMMI (capability levels) and DORA (self-assessability,

outcome focus). Framework validated against documented failure cases to ensure each dimension maps to observable risk.

---

<b>System Disclosure</b>	This report was created with a multi-agent research system. Human input: framework design, interpretation, and writing. Agent input: literature review, source synthesis, fact-checking, gap identification.
--------------------------	--

---

## 13. Claim Register

Top claims from this report with supporting evidence and confidence levels.

#	CLAIM	VALUE	SOURCE	CONFIDENCE	USED IN
1	McKinsey AI High Performers	6% ( $\geq 5\%$ EBIT)	McKinsey 2025 (n=1,993)	High	Ch. 2, 4
2	Enterprise agent experimentation rate	62%	McKinsey 2025	High	Ch. 2, 4
3	Agent project abandonment forecast	>40% by 2027	Gartner 2025	Medium	Ch. 4
4	Pilot-to- production conversion	54%	Gartner 2024	High	Ch. 4
5	LLM overconfidence rate	84%	PMC/12249208	High	Ch. 6
6	VCE bias	"systematically biased"	arXiv:2602.00279	High	Ch. 9
7	Air Canada chatbot liability	~\$800 + precedent	Tribunal ruling 2024	High	Ch. 6
8	Klarna AI overpivot	\$60M savings, partial rollback	CEO earnings call Q3 2025	High (corporate)	Ch. 6
9	Multi-agent hijacking success	45–64%	arXiv:2503.12188	High	Ch. 6
10	Budget-CoCoA cost	\$0.005/check	Anthropic pricing	High	Ch. 9, 10

11	SOC alerts ignored	67%	Vectra 2023 (n=2,000)	High	Ch. 6, 10
12	Knight Capital loss	\$440M in 45 min	SEC filing 2012	High	Ch. 6
13	EU AI Act penalty	€35M / 7% revenue	Legislative text	High	Ch. 2, 10

Top 5 claims with invalidation conditions:

**Claim 1 (McKinsey 6%):** Invalidated if McKinsey's EBIT attribution methodology is shown to systematically undercount AI impact.

**Claim 2 (62% experimentation):** Invalidated if the survey includes organizations already at Level 3+ due to unreported production deployments.

**Claim 9 (Multi-agent hijacking):** Invalidated if newer agent frameworks (post-2025) implement inter-agent trust protocols that reduce success rates below 10%.

**Claim 10 (Budget-CoCoA cost):** Invalidated if API pricing changes or if production-grade calibration requires additional infrastructure (database storage, logging) that significantly increases total cost.

**Claim 13 (EU AI Act penalty):** Invalidated if enforcement interpretation creates safe harbors or if penalties are reduced for first-time offenses.

## 14. References

- [1] McKinsey & Company (2025). "The State of AI in 2025." McKinsey Global Survey (n=1,993). Retrieved from mckinsey.com.
- [2] McKinsey & Company (2025). "The State of AI in 2025" — agent experimentation and deployment data.
- [3] Gartner (2025). "Predicts 2025: AI Agents" — projection on agentic project cancellation rates.
- [4] Gartner (2024). "4 Levels of AI Maturity and How to Achieve Them" — pilot-to-production conversion rate.
- [5] PMC/12249208 (2024). "Overconfidence in Large Language Models" — study of 9 LLMs across 351 scenarios.
- [6] arXiv:2602.00279 (2026). "Verbalized Confidence Expressions in LLMs: Calibration and Reliability."
- [7] Air Canada chatbot tribunal ruling (2024). Bereavement fare policy hallucination case.
- [8] Klarna Q3 2025 Earnings Call. CEO Siemiatkowski statement on AI overpivot and partial rollback.
- [9] arXiv:2503.12188 (2025). "Hijacking Attacks on Multi-Agent Systems" — 45–64% success rates across AutoGen, CrewAI, MetaGPT.
- [10] Budget-CoCoA methodology; Anthropic pricing (verified February 2026). \$0.005 per confidence check using 3x Haiku samples.
- [11] Vectra AI (2023). "2023 State of Threat Detection" — survey of 2,000 SOC analysts. 67% alert ignore rate.
- [12] Knight Capital SEC filing (2012). \$440M loss in 45 minutes from erroneous automated orders.
- [13] European Parliament (2024). "Regulation (EU) 2024/1689 — Artificial Intelligence Act." Articles 9, 14; penalty Article 99.
- [14] Deloitte (2024). "AI Readiness Framework" — five stages of AI maturity.
- [15] Microsoft (2024). "AI Maturity Model" — Azure AI adoption framework.
- [16] Google Cloud (2024). "AI Adoption Framework" — three phases of AI transformation.
- [17] IBM (2024). "AI Ladder" — four rungs of data-centric AI maturity.
- [18] Gartner (2025). "40% of enterprise applications will incorporate agentic AI by end of 2026." Market forecast.
- [19] Google (2025). "Agent-to-Agent (A2A) Protocol" — now Linux Foundation project. Protocol specification.
- [20] Anthropic (2024). "Model Context Protocol (MCP)" — tool integration standard for AI agents.
- [21] Market forecast: AI agent market from \$7.8B (2024) to \$52B (2030) at 45.8% CAGR. Multiple analyst sources.

**Citation:** Ainary Research (2026). *The AI Agent Maturity Model: A Framework for Measuring How Ready Your Organization Actually Is.* AR-004.

### About the Author

Florian Ziesche is the founder of Ainary Ventures, where AI does 80% of the research and humans do the 20% that matters. Before Ainary, he was CEO of 36ZERO Vision and advised startups and SMEs on AI strategy and due diligence. His conviction: HUMAN × AI = LEVERAGE. This report is the proof.

[ainaryventures.com](http://ainaryventures.com)



AI Strategy · Published Research · Daily Intelligence

Contact · Feedback

[ainaryventures.com](http://ainaryventures.com)  
[florian@ainaryventures.com](mailto:florian@ainaryventures.com)

© 2026 Ainary Ventures