

# The Defensible ROI

*How leading enterprises turn metered AI activity into a recomputable, board-ready ledger — and finally close the gap between what AI costs and what AI returns.*

---

## \$3.07M

net verified AI value, annualized — the worked example this paper recomputes in full: a 6.0× gross return on \$620K of metered spend, payback in roughly eight weeks, quality grade A-.

## ILLUSTRATIVE · THE MERIDIAN CASE

**\$3.07M** net verified value, annualized

A 5,000-person firm surfaced \$3.69M of value on \$620K of metered AI spend — a 6.0x gross return, payback in roughly eight weeks, and a board one-liner the CFO's team could re-derive from their own rate card.

## QUALITY GRADE

A- (72% ledger- and meter-backed)

## PAYBACK

~8 weeks

## MODELED UPSIDE

to ~\$4.5M

## EXECUTIVE SUMMARY

*Provenance beats precision.***The finding in three sentences**

Enterprise AI has a measurement problem that has nothing to do with the technology. Most firms cannot tell a CFO, with audit-grade evidence, what their AI returned last quarter — and the people who can are about to define the next category of enterprise software.

Three findings follow from a survey of the public research and from a worked example of the framework the AgentX platform implements:

- 1. Adoption is not value.** Two years into the enterprise AI buildout, the share of organizations with established ROI is stuck in the single digits — 8% in Q1 2026, 7% in Q2 (KPMG Global AI Pulse). The cost side is real and large; the return side is unauditably for most firms.
- 2. The four existing categories miss the join.** FinOps sees the cost. GRC sees the existence and risk. Observability sees the execution. None of them sees the return. AgentX sits on the join — and turns the same evidence that audits the work into the arithmetic that proves the value.
- 3. Defensibility is a primitive, not a feature.** The Credibility Ladder — a four-tier ranking that maps every dollar to where its evidence came from, and that labels every estimate honestly — lets a firm present a board one-liner the CFO can re-derive without trusting the vendor. That primitive is the system of record for AI value.

The rest of this paper walks the evidence behind each finding, then runs the framework end-to-end on a worked example — every number recomputable, every assumption disclosed, every haircut visible. In this category, transparency is not a marketing posture. It is the product.

*Designed for AI-forward mid-market and lower-enterprise organizations (~500–5,000 employees) with active AI deployments in at least two functions. Not a regulated-vertical primer; the optional Regulated Industries module is a separate artifact.*

## PART ONE

# The context: capability has moved to the edge

*“If you discover yourself in partnership with the A.I. system, you are uniquely vulnerable to all of the failures of that A.I. system.”*

JACK CLARK, CO-FOUNDER, ANTHROPIC · VIA THE NEW YORK TIMES

## 1.1 The AI Sprawl Thesis

Ninety-five percent of organizations report having an AI strategy. The share that can prove what theirs returns is stuck in the single digits — 8% in the first quarter of 2026, 7% in the second (KPMG Global AI Pulse). That gap — between strategy and evidence, between intent and proof — is the single most important fact about enterprise AI in 2026. Two years into the buildout, the needle is not moving.

The shift behind the gap is structural. In the prior era, AI capability was scarce and centralized — owned by a small group of specialists in research or platform engineering. That has changed. Today, every desk, team, and function can stand up its own copilots, agents, and embedded models. Capability is no longer scarce; it is ambient.

*“Ultimately, there is no agentic future without trust and no trust without governance that keeps pace.”*

STEVE CHASE, GLOBAL HEAD OF AI AND DIGITAL INNOVATION, KPMG INTERNATIONAL · Q2 2026

These builds happen in isolation. Each team selects its own models, writes its own prompts, instruments its own evaluations, and ships its own integrations. No shared inventory. No common standards. No cross-pollination. The institution accumulates redundant *effort* (the same capability rebuilt in five places), redundant *spend* (compute, tokens, and validation labor duplicated), and redundant *risk* (every new build widens the surface; nothing connects).

### WHY IT MATTERS NOW

Three forces converge in 2026 to make the sprawl expensive enough to matter: regulatory frameworks (EU AI Act and sector-specific guidance) demand an aggregate, documented view of AI; the duplicated spend is real, large, and ownerless on the P&L; and leadership can no longer answer “where do we stand on AI?” to a board, an allocator, or an examiner — because no single source of truth exists.

## 1.2 The Measurement Gap

The dominant metric for enterprise AI today is adoption: seats, prompts per day, weekly active users, satisfaction scores. They are useful diagnostics; they are not value. A team that processes ten thousand prompts a day has not produced ten thousand outcomes — and most of the outcomes that do exist are unmeasured, unattributed, and unaudited.

The spend behind those diagnostics is real and accelerating. Gartner forecasts \$2.59 trillion of worldwide AI spending in 2026 — up 47% year over year — on the way to \$3.49 trillion in 2027 (Gartner, May 2026). And Gartner's own lead analyst attaches the caveat that matters: “CIOs face

challenges in proving the value from AI investments and demonstrate tangible business outcomes,” notes John-David Lovelock, Distinguished VP Analyst. The question is not whether enterprises will spend at that scale. The question is whether any of it will be defensible when an allocator or an examiner asks where it went.

The finance function has now said this out loud. In EY's June 2026 survey of CFOs, 47% said their teams cannot effectively measure AI value, and 71% said traditional financial metrics are insufficient for evaluating AI investments (EY, via CFO Dive, June 2026). The gap is precise. Vanity metrics — adoption rate, daily prompts, user satisfaction, uptime — are what most vendors lead with. Value metrics — revenue impact, cost reduction, risk avoided, cycle-time compression — are what CFOs ask for and almost never get: each a *business* metric, not an *AI* metric.

The gap is also one of perception. Wall Street Journal reporting in 2026 found a 38 -percentage -point split: roughly 40% of workers say AI saves them zero hours a week, against just 2% of the C-suite who say the same (WSJ, 2026). When the people closest to the work see no return and the people farthest from it believe it is large, the disagreement is itself the data.

The value side is just as stark. McKinsey's State of AI 2025 — the most comprehensive read available — found that while 88% of organizations now use AI, only about 6% qualify as high performers — defined as attributing at least 5% of EBIT to their AI use. Adoption is nearly universal; captured value is rare and concentrated (McKinsey State of AI 2025, n=1,993 across 105 nations).

Two further findings round out the picture, each from an independent measurement:

- **95% see no measurable return.** An MIT Project NANDA study that drew wide attention — and methodological criticism — reported 95% of enterprise GenAI pilots showing no measurable P&L impact. The finding does not stand alone: KPMG finds established ROI stuck in the single digits (8% in Q1 2026, 7% in Q2), and BCG finds 74% of companies have yet to show tangible value from AI. Three independent measurements, one conclusion — a skeptic cannot dismiss all three.
- **The ungoverned share is already being breached.** 13% of organizations reported breaches of their AI models or applications; of those, 97% lacked proper AI access controls — and breaches linked to shadow AI added up to \$670K to the average breach cost (IBM Cost of a Data Breach, 2025).

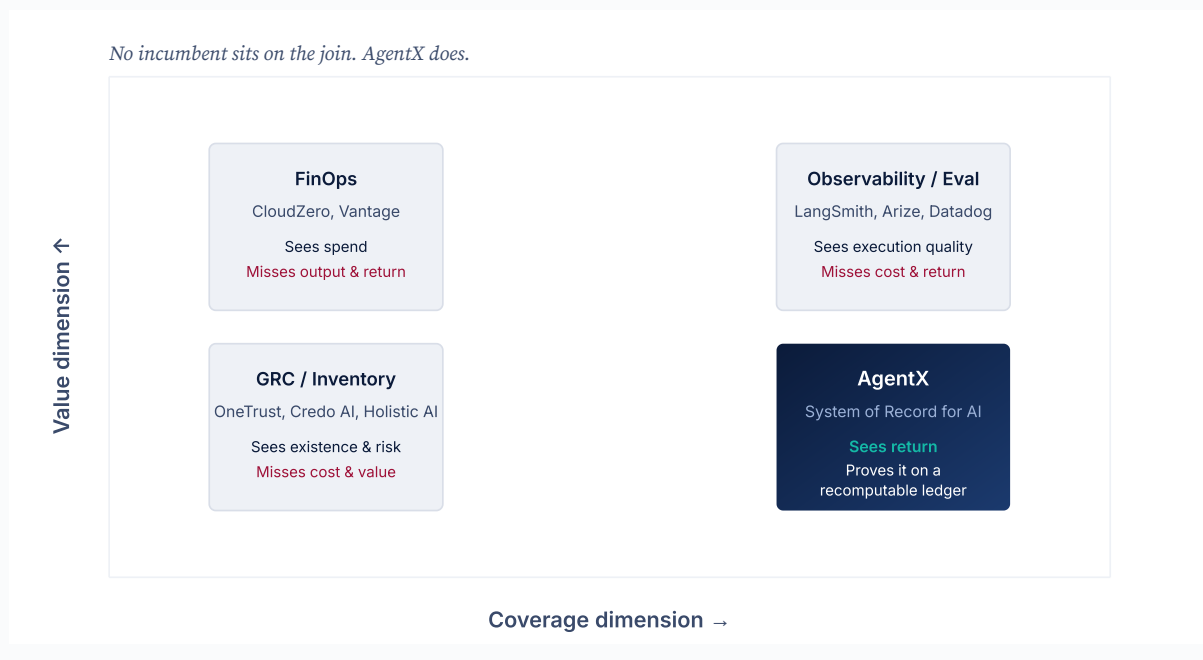
These are not problems the next vendor in any of the four existing categories will solve. They are symptoms of a missing system of record.

### 1.3 Why Existing Approaches Fail

Four categories of enterprise software already touch AI. Each sees a different slice; none sees the return.

The pattern repeats across all four: each incumbent sees one slice of the problem, integrates well into its own slice, and has no economic incentive to extend into the next. The buyer's experience is a portfolio of dashboards that do not add up to an answer. The CFO's question — “what is our AI returning, on a ledger my auditor will accept?” — has no incumbent owner.

## Four categories, one gap



**FIG 1** – Mapping existing categories along two axes — what dimension of AI activity they cover (horizontal) and what dimension of value they surface (vertical) — exposes the gap. The bottom-right quadrant is the join between coverage and value, and the only product category positioned there is AgentX.

One further finding points to where that return actually lives — and the advisory literature is converging on it. PwC's AI Business Survey estimates that only about 20% of the value an organization captures from AI comes from the technology itself; the remaining 80% comes from redesigning the work around it — the workflow, the org, the incentives, the review cadence. Existing tooling markets itself on the 20%. The value lives in the 80%: the redesign that turns a deployed model into a measured outcome — and the work the next Launchpad engagement is scoped to do.

That distinction — the technology versus the redesign around it — is why a system of record is the right shape for the answer. A point tool that deepens one slice will not produce a board-ready ledger. A consulting engagement that produces a static binder will not survive the next quarter. What is needed is a living artifact that the executive and the worker both trust — one that translates metered activity into a recomputable number and that translates workflow redesign into a Tier 2 dollar.

## PART TWO

# The data: provenance beats precision

## 2.1 The Credibility Ladder

A board does not reject an AI value claim because the number is small. It rejects it because the number is unsourced. The first principle of credible AI measurement is therefore not accuracy — it is *provenance*: where did the evidence for each dollar come from?

The AgentX framework answers that question with **the Credibility Ladder**. Every dollar claimed is ranked on one of four tiers by the strength of its evidence, and the headline is built only from dollars the customer's own finance team could re-derive without trusting the vendor.

TIER	EVIDENCE TYPE	SURVIVES THE BOARD BECAUSE...	HEADLINE?
<b>1 — Ledger-verified</b>	Money that already moved or stopped: cancelled licenses, held backfills, cut overtime	It is in the general ledger; the CFO can pull it	Yes (the floor)
<b>2 — Metered x rate card</b>	System-counted outputs x the customer's own finance-signed unit cost	Counts come from logs (no self-report); rate countersigned by FP&A	Yes
<b>3 — Modeled with counterfactual</b>	Throughput, win-rate, and cycle-time gains, with an attribution haircut applied	Stated assumption + disclosed haircut; quarantined from the floor	Labeled upside only
<b>4 — Strategic / optionality</b>	Risk reduction, capability build, retention	Honestly un-dollarized	Narrative only

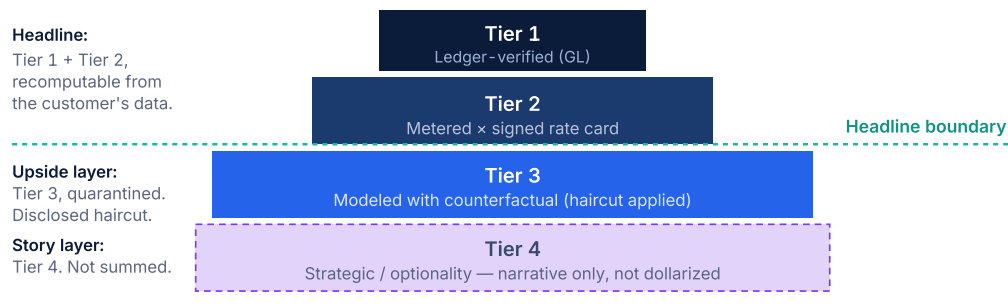
The discipline this enforces is the one that makes the number defensible: Tier 1 and Tier 2 dollars are inputs the customer can audit; Tier 3 dollars are presented as upside with their haircut visible; Tier 4 is told as a story, not a number. A single undifferentiated total that mixes all four is the pattern AgentX refuses to produce.

**OPERATIONAL DEFINITIONS**

**Tier 1** requires a general-ledger entry that closes inside the measurement window. **Tier 2** requires a system-counted event linked to a unit rate countersigned by FP&A in the same fiscal quarter. **Promotion** from Tier 3 to Tier 2 requires at least one quarter of acceptance-sample data with a confidence interval below ±10%.

The framework is closest in spirit to Forrester's Total Economic Impact methodology — but with a strictly enforced separation between hard and estimated dollars, an immutable record of every promotion, and an explicit refusal to blend Tier 4 narrative into the headline.

## The Credibility Ladder — four tiers, one headline boundary



*The headline rests on Tier 1. Everything below it is upside at a stated haircut, or story.*

*Above the dashed line: the floor. Below: labeled upside, then narrative only.*

**FIG 2** – A single undifferentiated total destroys the trust the framework rests on. The headline is the top two tiers only; modeled upside is shown separately at a disclosed haircut; Tier 4 is told as a story, not a number.

## 2.2 How the numbers are built

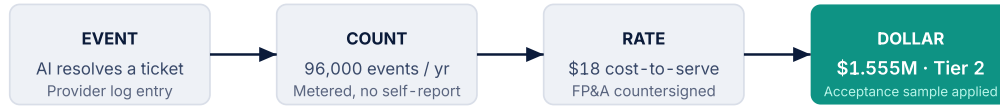
Three mechanisms do the work that makes the ladder operational.

**The unit-rate library.** The single most important move. Capture once, at onboarding, from the customer's own FP&A: cost-to-serve per ticket, loaded dollars per hour by role, cost per contact. Counter-sign and version. Dollar conversion is then arithmetic on metered counts, not estimates per event. This is also why output × unit-cost beats time-saved: time-saved asks the worker "how long would this have taken?" (noisy, inflatable, per-event self-report). Output × unit-cost moves the estimate from the worker to finance, and from per-event to a stable average finance already budgets against.

**Acceptance sampling.** Instead of asking every user about every task, randomly sample  $N$  metered events; a human verifies "AI did this, it was usable." Produces a verified acceptance rate with a confidence interval (e.g. 90% ± 3%,  $n=384$ ) applied as a population haircut. That is statistics, not anecdote.

**Conservative defaults.** Attribution defaults low (e.g. 50% for cycle-time gains, 40% for sales-response win-rate effects). Only the customer can ease the haircut upward with their own data; the framework cannot soften it below the conservative default.

## From a logged event to a defensible dollar



Three human inputs in the entire chain — the rate, the attribution, the counterfactual — all captured once and signed.

**FIG 3** – The transformation from a logged event to a defensible dollar has three human inputs and four automatic steps. Once the rate is signed, every additional event is arithmetic.

### 2.3 The Meridian case

The following worked example uses figures published in AgentX's methodology documentation. Meridian is a 5,000-employee firm running AI in support, operations, sales-assist, and knowledge work. All figures are annualized from a 90-day observation window, per the AgentX methodology. Acceptance rates carry a 95% confidence interval of  $\pm 0.03$  ( $n=384$  sampled events per category).

A. COST DISPLACEMENT	TIER	SOURCE	ANNUALIZED (\$K)	
Cancelled 2 SaaS tools AI replaced	Tier 1	Vendor invoices	140	
3 attriting support reps not backfilled (3 × \$65K loaded)	Tier 1	Payroll	195	
Overtime / contractor hours cut (4,000 × \$40)	Tier 1	Timesheets	160	
<i>Subtotal — incontrovertible</i>			<b>495</b>	
B. LABOR-OUTPUT SUBSTITUTION (METERED × RATE × ACCEPTANCE)	COUNT	UNIT COST	VERIFIED ACCEPTANCE	\$K
Support tickets resolved w/o human	96,000	\$18 cost-to-serve	0.90	1,555
Self-serve deflection	120,000	\$6	0.95	684
Knowledge drafts accepted	40,000	\$30 (0.60 hr × \$50 loaded)	0.80	960
<i>Subtotal</i>				<b>3,199</b>

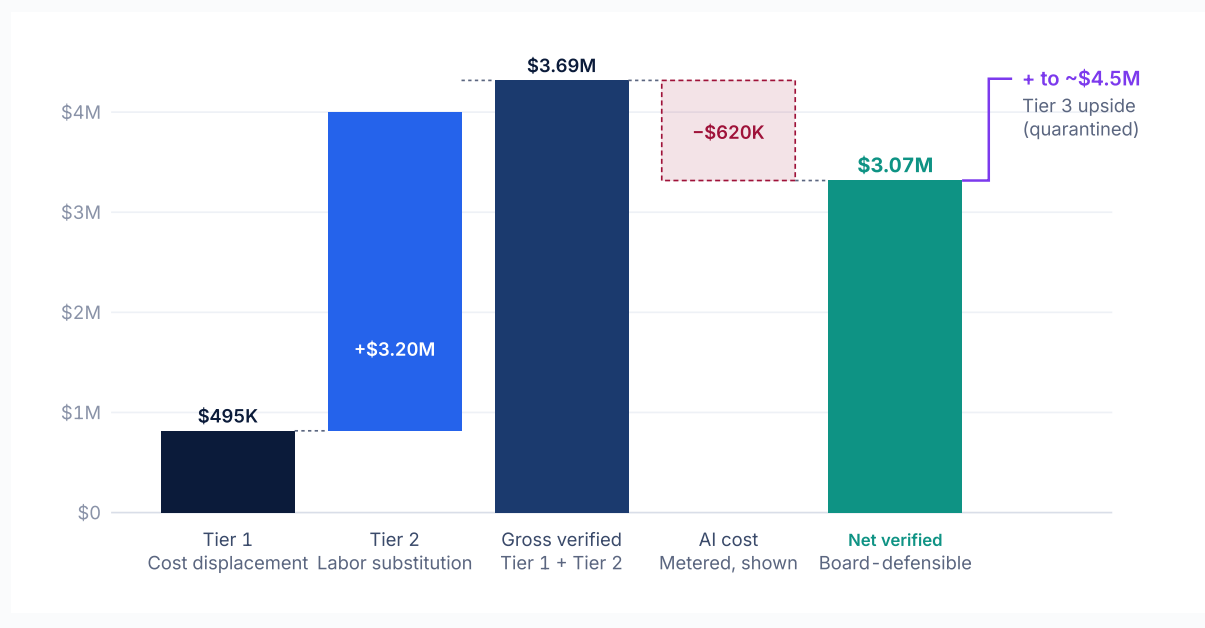
Below the headline, in the upside layer:

C. THROUGHPUT & CYCLE-TIME	TIER	METHOD	\$K
Sales response 4h → 5min on 20,000 leads (modeled)	Tier 3	+1.5pt win-rate × deal value, 40% attribution haircut	~1,400

*Labeled upside only. Not added to headline.*

D. AI COST (THE DENOMINATOR)	SOURCE	\$K
Platform + tokens + amortized implementation	Metered	620
<b>Gross verified (A + B)</b>		<b>3,694</b>
<b>Less AI cost (D)</b>		<b>(620)</b>
<b>Net verified value (annualized)</b>		<b>3,074</b>
<b>Gross return on \$620K spend</b>		<b>6.0×</b>
<b>Approximate payback</b>		<b>~8 weeks</b>
<b>Quality grade</b>	<b>A- (72% ledger- and meter- backed)</b>	
<b>Modeled upside (Tier 3, quarantined)</b>		<b>to ~\$4.5M</b>

### How the \$3.07M net is built



**FIG 4** – The waterfall is internally consistent: \$495K + \$3.20M = \$3.69M gross verified, less \$620K metered AI cost = \$3.07M net. Tier 3 modeled upside sits above as a bracket, never added to the headline. Source: AgentX value framework, worked example.

### HONESTY IN THE COST COLUMN

The single largest trust signal in any AI ROI claim is the AI cost itself. Showing the denominator — and subtracting it — is what separates a real number from a vendor slide.

Meridian is a worked example. The same shape of claim has now been made by a real institution, in public: Commerzbank projects **€300M in benefits** against €140M of AI investment — an implied ROI of roughly **120%** — material enough to account for roughly a quarter of the bank’s guided profit growth through 2028 (as reported by Bloomberg).

**~120%**    **€300M**                      **€140M**                      **~25%**  
 IMPLIED ROI      PROJECTED BENEFITS      AI INVESTMENT      OF GUIDED PROFIT GROWTH

A bank does not publish a number like that casually; it publishes it because it can defend the arithmetic. That discipline — benefits a CFO will put beside spend, in public — is exactly what the Credibility Ladder operationalizes for firms without a bank’s measurement staff.

## 2.4 The Day-1 / Day-30 / Day-90 progression

The framework is designed to produce something defensible on day one — and progressively more defensible as evidence accumulates.



**Day 1.** Metered spend, metered usage, metered output counts, and a provisional rate. Annualized projection; payback estimate; grade marked *Provisional*. Independently useful — most organizations do not know their AI spend or who is using it — with zero baselines required.

**Day 30.** The customer's FP&A countersigns the rate card. The first event-sample audit establishes an acceptance rate with a confidence interval. The earliest Tier-1 deltas (a cancelled tool, a held backfill) appear on the ledger. A real grade is on the number.

**Day 90.** Tier-1 deltas realized, Tier-2 metered-and-sampled, Tier-3 modeled upside quarantined, the headline and grade and QoQ trend on a fully recomputable backing ledger. This is the artifact that survives the board challenge.

## 2.5 The Governance Multiplier

The Meridian case measures one firm's floor. A separate body of survey evidence measures the ceiling — and it is consistent enough to deserve a name. Call it **the Governance Multiplier**: in every major survey that quantifies the link, organizations that govern their AI — accountability at the top, full visibility into cost, real-time oversight of the work — are consistently *more likely* — up to several times more likely — to realize value than organizations that do not. Three independent samples, three different methods, one direction.

**EY — real-time oversight.** The EY Responsible AI Pulse (October 2025, n=975 C-suite respondents across 21 countries) found that organizations with real-time AI monitoring were **34% more likely to see improvements in revenue growth** and **65% more likely to see improved cost savings** than organizations without them (EY, 2025).

**KPMG — accountability and cost visibility.** KPMG's Q2 2026 pulse (n=2,145) isolates two levers, and the split is stark. Organizations with clearly defined AI accountability at the leadership level are **3.5x more likely to report established ROI** than organizations without it (14% versus 4%); organizations with full visibility into AI operating costs are **5x more likely** (15% versus 3%) — in KPMG's own words, "five times more likely to achieve ROI." Note the unit: a greater *likelihood* of proven ROI, not a multiple on the return itself (KPMG Global AI Pulse, Q2 2026).

**McKinsey — the maturity gap.** McKinsey's State of AI 2025 supplies the third leg: **71% of high performers have mature AI risk-management practices in place, versus 38% of everyone else** — a 33-point governance-maturity gap between the roughly 6% of organizations that qualify as high performers (attributing at least 5% of EBIT to AI) and the rest (McKinsey, 2025, n=1,993). The multipliers are large precisely because the discipline behind them is still scarce.

### The Governance Multiplier — three independent samples, one conclusion

EY RESPONSIBLE AI PULSE · OCT  
2025 · N=975

**34% / 65%**

more likely to see improvements in revenue growth / in cost savings, with real-time monitoring

KPMG GLOBAL AI PULSE · Q2 2026 ·  
N=2,145

**3.5x / 5x**

more likely to achieve established ROI, with leadership accountability (14% vs 4%) / full cost visibility (15% vs 3%)

MCKINSEY STATE OF AI · NOV 2025 ·  
N=1,993

**71% vs 38%**

high performers vs everyone else with mature AI risk-management practices — a 33-point maturity gap

*Every source that quantifies the governance-to-value link finds the same direction: governance is not a tax on AI value; it is the multiplier on it.*

**FIG 6** – Three surveys, three methods — oversight-practice splits (EY), established-ROI rate splits (KPMG), maturity-gap analysis (McKinsey) — on three independent samples. All multipliers are stated as likelihoods of realized value, never as multiples on the return itself.

## THE HONEST CAVEAT — AND THE LIMITS OF THE CLAIM

These multipliers are cross-sectional correlations from independent self-reported samples, not causal estimates: organizations that govern well also measure well, and no one has run the randomized trial. What survives scrutiny is the consistency — three samples, three methods, one direction — plus one outcome-based datum pointing the same way (IBM's breach-cost data, measured in dollars rather than sentiment). Nor does the framework reach everywhere: a firm running AI in a single cost center can reasonably keep bolting ROI math onto its FinOps tooling, and frontier research deployments or AI embedded invisibly inside vendor products resist event-level metering altogether. Stating this plainly is not a concession; it is the argument. A correlation this consistent only becomes a defensible number for *your* firm through an evidence layer that meters, grades, and recomputes it on your own data.

Read together, the methodology points one way: these multipliers are not the model's effect — they are the effect of the operating model around it. Accountability, cost visibility, and real-time oversight are exactly the plumbing a system of record installs. The EY and KPMG numbers describe what that plumbing is worth once it exists; the Meridian floor is what it produces in the first ninety days. The floor and the ceiling are the same story told at two altitudes.

## PART THREE

# AgentX: the system of record for AI

## 3.1 The category claim

AgentX is the AI governance system that audits the value it helped you create — recomputable from your own finance data, defensible to your board.

Read that sentence the way a board would. The asymmetry is structural: the system that meters the work is the only system that can audit the value of the work — and no incumbent fuses the two. And the proof standard is the highest one on offer: not "trust our dashboard," but "recompute it from your own finance data." Governance that pays for itself in evidence is not a tax on AI value. It is the multiplier on it — the 3–5× greater likelihood of established ROI that KPMG measures, in the direction every survey of the link keeps finding.

## 3.2 The trust primitive

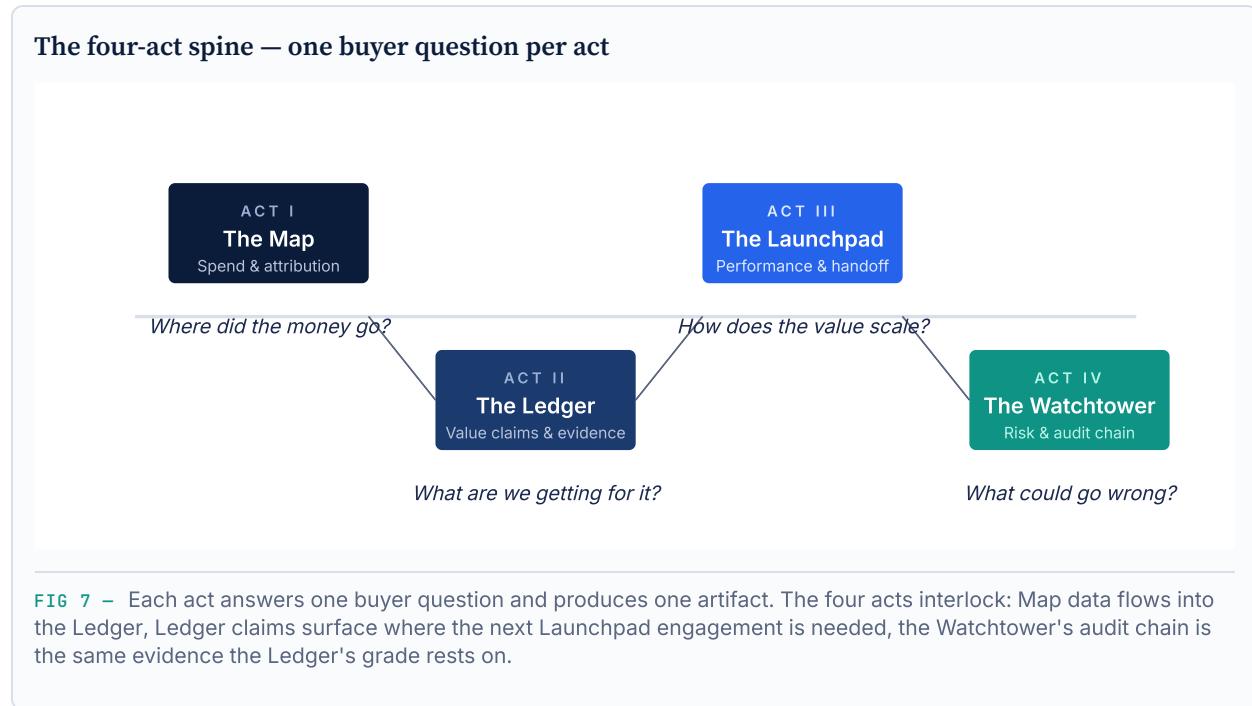
The mechanism that survives the board is older than any vendor in the category. It is the same primitive that turns a journal into a ledger: every mutation logged, no entry posted without separation of duties, every total recomputable from the underlying entries.

Applied to AI value, the primitive has three components. **Hardened Value** is the customer's own numbers, promoted up the Credibility Ladder, on an immutable record the customer owns — defensible. **Assisted Estimate** is a product-suggested metric or peer benchmark offered to help a customer who has no numbers yet — always labelled as an estimate, never folded into a hardened total. **Separation of duties** means no user can promote their own claim past validation; the rate card is countersigned by FP&A, not authored by the team selling the product. The labelling line between Hardened and Assisted is the line on which the entire defensibility claim rests.

This is not a feature. It is the primitive the system runs on. Once an audit committee has agreed that *this* is the system by which AI value gets promoted from "someone said it" to "audit-signed," switching it out is a real decision — the same kind of decision that is rarely revisited once a general ledger is in place.

### 3.3 The four-act spine

The buyer's journey through AgentX is structured as four acts. Each answers a single buyer question and produces a single artifact. AgentX is not a replacement for observability tooling or for governance inventory — it is the system of record that turns observability and inventory data into a board artifact. Most customers keep their existing LangSmith, OneTrust, or CloudZero deployments; AgentX sits above them as the join.



**Act I — The Map (Governance).** Where did the money go? Spend, usage, and attribution across every model, copilot, and agent the firm runs — provider billing, cloud audit logs, and identity/workspace signals, intersected with a fingerprint channel the customer owns end-to-end. The Map makes the sprawl legible.

**Act II — The Ledger (Value capture).** What are we getting for it? Value claims on the Credibility Ladder, every dollar mapped to where its evidence came from. This is the wedge; this is where the board one-liner lives.

**Act III — The Launchpad (Design handoff).** How does the value scale? Where the Ledger surfaces a gap, the next agentic-solution engagement is scoped. Design engagements feed real requirements back into the product.

**Act IV — The Watchtower (Risk & audit).** What could go wrong? Shadow-AI breaches, model drift, access-control failures, and the immutable audit chain that proves governance. The market is candid about why this act cannot wait: 74% of organizations plan to deploy agentic AI within two years, while only 21% report mature oversight frameworks (FT/Deloitte, June 2026). "We need to do more work not only on their capabilities, but also on controls and governance," as ACCA's Alistair Brisbourne puts it. Risk findings are also a leading indicator of ROI: an incident erodes the value claim before the value materializes.

### 3.4 The board one-liner

#### THE HEADLINE, IN ONE LINE

AgentX: \$3.1M net verified AI value (annualized) — 6.0× gross return on \$620K spend, paid back in ~8 weeks. Quality grade A- (72% ledger- and meter-backed, 28% modeled). Modeled upside to ~\$4.5M.

The artifact AgentX produces at the end of a 90-day onboarding is a single page. That one line is its headline; everything beneath it is the backing ledger, built for the analyst sent to break it. The claim earns its grade the only way that matters: someone tried to break it, and could not.

*Every dollar of value is ranked by where its evidence came from. The customer's own GL entries and metered counts become the floor; modeled upside sits in a separate, hair-cut bracket; strategic narrative is told as a story, never a number. The rate card is countersigned by FP&A, never by the team selling the product. The headline is built only from what survives an audit.*

#### METHODOLOGY IN ONE PARAGRAPH

Meridian's \$3.07M net verified is the floor — the number a CFO can defend today. The ceiling sits higher: the Governance Multiplier establishes that organizations that govern AI well — clear accountability at the top, full visibility into cost — are 3–5× more likely to achieve established ROI (KPMG), with EY's real-time-oversight split pointing the same direction. The Meridian floor and that evidenced ceiling are the two ends of the spectrum a system of record is built to make visible at once — the audited present beside the evidenced upside, on the same page.

### A one-page recomputable value ledger — illustrative

DRIVER	TIER	COUNT	RATE + SOURCE	ATTRIBUTION	ANNUALIZED \$	EVIDENCE
Cancelled SaaS tools (2)	1	—	Vendor invoices	100%	\$140K	GL entry
Held backfills (support)	1	3 FTEs	Payroll, \$65K loaded	100%	\$195K	HR roster
Cut overtime / contractors	1	4,000 hrs	Timesheets, \$40	100%	\$160K	Time system
Tickets resolved w/o human	2	96,000	\$18 cost-to-serve (FP&A v3)	0.90 acceptance	\$1,555K	Provider logs + sample
Self-serve deflection	2	120,000	\$6 per resolution	0.95 acceptance	\$684K	Provider logs + sample
Knowledge drafts accepted	2	40,000	\$30 (0.60 hr × \$50 loaded)	0.80 acceptance	\$960K	Provider logs + sample
<b>Gross verified (Tier 1 + 2)</b>					<b>\$3,694K</b>	
AI cost (platform + tokens)	1	—	Metered	100%	(\$620K)	Provider billing
<b>Net verified value (annualized)</b>					<b>\$3,074K</b>	
Sales response cycle-time	3	20,000 leads	+1.5pt win-rate × deal value	40% haircut	~\$1,400K	Modeled; quarantined

FIG 8 – Every row is sourced. Every total is recomputable. The customer's CFO can pull the rate card, the customer's auditor can pull the underlying logs. The headline is for the boardroom; the table is for the analyst sent to break it.

## PART FOUR

### Implications

The mechanism is in place. The question is what a CFO, a Head of AI, and a board do with it.

#### 4.1 What this means for CFOs

The risk of the next AI budget cycle is not under-investment; it is mis-investment that cannot be defended. A board that asks “what is our AI returning?” and receives a vendor ROI deck will, increasingly, demand the arithmetic behind it. The CFO who cannot produce the arithmetic is the CFO who loses the next AI budget cycle.

The arithmetic is producible — and the bar it must clear is not sentiment but cost of capital. As Pankaj Mhatre put it this year, writing for CFA Institute's *Enterprising Investor*: “If incremental returns on AI-related investment fail to exceed the firm's weighted average cost of capital (WACC), shareholder value is not created regardless of technological sophistication.” A recomputable ledger is the only instrument that lets a CFO test AI spend against WACC the way every other capital

allocation is tested. It is expensive to build once. Every quarter's answer after that is nearly free — and 97% of finance chiefs say their boards now expect regular AI readouts (OneStream, via CFO Dive).

## 4.2 What this means for Heads of AI

The Head of AI's job has quietly changed. "Leaders are no longer content to run pilots. They want proof. GenAI is being held to the same standards as other major investments" (Wharton, 2025). Three partial answers — an adoption dashboard, a vendor ROI deck, a compliance binder — no longer survive that room. The move is from three partial answers to one defensible one, and it has to happen without an inventory project, without a multi-quarter consulting engagement, and without surrendering ownership of the rate card. That is the job AgentX was built to do.

### SELF-ASSESSMENT · FIVE QUESTIONS YOUR BOARD WILL ASK

1. What did AI cost us last quarter — metered, not estimated?
2. What did it return — on evidence our auditor would accept?
3. Who signed the rates behind that number — finance, or the vendor?
4. Can we recompute the headline on demand, from our own data?
5. Which of our value claims expire, and when?

Score yourself in five minutes at [agentxscorecard.com](https://agentxscorecard.com). Fewer than four confident answers is the gap this paper is about.

## 4.3 The 90-day decision

### Before and after: a 90-day cycle



**FIG 9** – The transition is not from “no answer” to “an answer.” It is from “a number you cannot defend” to “a number your auditor can re-derive.” The cost is the same conversation you would have had about AI ROI anyway — the difference is that the conversation now produces an artifact.

## 4.4 About AgentX

AgentX is the System of Record for enterprise AI — the platform that turns the metered activity of every model, copilot, and agent a firm runs into a recomputable, board-defensible ledger of value, governance, and risk. Built product-first; the consulting engagements exist to accelerate lighthouse deployments and feed real requirements back into the platform.

### THE 14-DAY SANDBOX

No demo gate. Sign in, register two or three agents, and see the value → governance → performance loop on synthetic data before committing. The same workflow runs on your data on day one of a production deployment.

Within a year, every AI line item in your budget will be challenged at least once — by a board member, an allocator, or an auditor. Firms that answer with adoption charts will re-litigate their AI budget every quarter. Firms that answer with a recomputable ledger will spend that meeting deciding where to scale. Provenance beats precision. Open a sandbox today; put a provisional, honestly-graded number in front of your CFO in fourteen days — and a board-ready one in ninety.

For the methodology in full — including the worked Meridian case, the operational definitions for each tier, and the statistical treatment of acceptance sampling — request the AgentX Value Framework.

---

## Appendix · Glossary

<b>System of Record (SoR)</b>	The canonical, single, living source of truth for every model and AI tool an institution runs. Delivered as SaaS.
<b>Defensible ROI</b>	An AI value claim that the customer's own finance team can re-derive without trusting the vendor. The framework, not the number, is the asset.
<b>Hardened Value</b>	The customer's own numbers promoted up the Credibility Ladder on an immutable record the customer owns. Defensible.
<b>Assisted Estimate</b>	A product-suggested metric or peer benchmark offered to help a customer who has no numbers yet. Always labelled as an estimate; never folded into a hardened total.
<b>Separation of duties (SoD)</b>	The discipline that no user can promote their own claim past validation. The rate card is countersigned by FP&A; the customer cannot be both author and auditor.
<b>The Credibility Ladder</b>	Four tiers of evidence (ledger-verified / metered × rate card / modeled with counterfactual / strategic). The headline is built from tiers 1–2 only; tier 3 is labeled upside; tier 4 is narrative.
<b>Fingerprint channel</b>	The attribution path AgentX owns end-to-end. It is the primary path; the three-signal stack (provider billing, cloud logs, identity / workspace) is the safety net.
<b>Quality grade</b>	A bond-rating-style classification of a value claim's evidence quality. The Meridian case carries an A– at 72% ledger- and meter-backed.

### ABOUT THIS PAPER

The Defensible ROI — An AgentX whitepaper, July 2026. Methodology references the AgentX Value Framework; worked example from the Meridian case. All figures are illustrative.

### CONTINUE

[Read the methodology](#) →  
[See a sandbox demo](#) →  
[Forward to your CFO](#) →

### CONTACT

[research@agentx.example](mailto:research@agentx.example)  
[agentxscorecard.com](https://agentxscorecard.com)