

### THE 80% MARGIN IS GONE. NOW WHAT?

Every SaaS company is adding AI. Most of them are losing money on every query. Traditional B2B SaaS ran at **80% to 90% gross margins** because the marginal cost of an additional user was close to zero. AI-first companies are operating at **50% to 60% margins**, with some early-stage companies as low as 25%. The pricing models haven't caught up. That's the trap.

But here's what most of the conversation gets wrong. Not all SaaS companies are the same. Not all products need an AI component. Some products have multiple AI components maturing in real time, each with a different cost structure, usage pattern, and value delivery mechanism. Treating them as one category and debating 'seat vs. usage vs. outcome' at the company level is the wrong conversation. **The real architecture happens at the workload level.**

**The problem statement is simple:** not all AI workloads do the same thing or cost the same. How are you going to price them without understanding the cost-to-serve variants? A single product might contain standard inference, RAG-based retrieval, agentic orchestration, and batch processing, each with fundamentally different economics. Pricing them with one model is like charging the same rate for electricity, water, and gas because they all come through pipes.

### THE STRUCTURAL SHIFT: SOFTWARE AS LABOR

The industry has moved from an Ownership Era (perpetual licenses) through an Access Era (seat-based subscriptions) into what Andreessen Horowitz calls the **Value Era**: customers pay for a job to be done. Software is becoming labor. Bessemer reports buyers now evaluate AI as a productive teammate capable of independent execution. If your AI handles 45% of support tickets autonomously, the customer needs **fewer seats, not more**. For the vendor to grow revenue under a per-seat model, they'd need the AI to fail. Seat-based pricing dropped from **21% to 15%** of the market in a single year, with **-40% lower gross margins** and **2.3x higher churn** for companies sticking to per-seat pricing for AI products.

### THE MARGIN WEDGE

GitHub Copilot launched at \$10/month with unlimited usage. Compute costs ran as high as **\$80 per user/month** for heavy users, with average losses of \$20/user. By mid-2025, Microsoft introduced \$0.04/request pricing beyond caps. Replit scaled from \$2M to \$144M ARR but only achieved positive gross margins by moving to usage-based models. Clay accelerated from \$2M to \$37M ARR by iterating pricing twice per year (seat → hybrid → pure usage). Even as inference costs drop 80-90% per year, user consumption accelerates faster. This is the **Jevons Paradox** applied to AI compute.

[EXHIBIT: The Margin Wedge: Traditional SaaS vs AI-First SaaS]

Metric	Traditional B2B SaaS	AI-First SaaS
Gross Margins	80% – 90%	50% – 60% (early stage: ~25%)
Marginal Cost/User	Near-zero	20–40% of revenue (variable per query)
Infrastructure % of Revenue	8–12%	25–40%
Revenue Scaling	Sublinear cost growth	Costs scale proportionally to usage
Power User Risk	Low (more usage = more retention)	High (more usage = more COGS)
Pricing Model Risk	Seat-based works: marginal cost is zero	-40% margins, 2.3x churn if seat-based persists

### NOT ALL AI WORKLOADS ARE CREATED EQUAL

A SaaS product might contain six or more distinct AI workload types, each with a different cost driver, scaling pattern, and pricing implication. **Inference now accounts for 55%+ of total AI cloud infrastructure spend** (\$20.6B of \$37.5B), surpassing training for the first time. Deloitte estimates inference will reach two-thirds of all AI compute by end of 2026. But inference is not one thing. Real-time chat, batch processing, RAG retrieval, and agentic orchestration all fall under 'inference' with wildly different cost profiles.

The table below ranks workload types by estimated share of AI infrastructure spend, includes the typical user action that triggers each workload, provides a cost estimation formula, and narrows the unit cost to a mid-market benchmark rather than the full range.

[EXHIBIT: AI Workload Taxonomy: Cost Drivers, Spend Share, Typical Queries, and Cost Estimation]

Workload Type	Est. Share of AI Spend	Typical User Action	Cost Driver	Scaling Pattern	Cost Estimation Formula	Mid-Market Benchmark	Passed Through?
<b>Real-time Inference</b>	~40–45%	User prompts chatbot, asks a question, requests a draft	Tokens (in + out), GPU time, model tier	Non-linear, spiky per user session	Queries/mo × avg tokens/query × \$/token	~\$0.002–\$0.05 per query (GPT-4o-mini to GPT-4o class)	Yes: per-token, per-call, per-session, or credits
<b>Agentic Orchestration</b>	~5–8% (fastest growing)	"Research competitors and write a report"; multi-step autonomous workflow	Multiple LLM calls, tool use, state mgmt, verification loops	Highly non-linear; 10–100+ steps per query	Sessions/mo × avg steps/session × avg \$/step	\$0.50–\$5.00 per session (frontier models)	Yes: per-session, per-workflow, per-outcome
<b>RAG / Vector Search</b>	~8–12% (often bundled with inference)	"What does our policy say about X?"; knowledge base lookup	Vector DB ops, embedding retrieval, generation tokens	Linear per query, spiky at bulk ingest	Queries/mo × (retrieval cost + generation tokens × \$/token)	\$0.01–\$0.10 per query; \$0.40–\$0.70 per resolved workflow	Yes: bundled into AI query credits
<b>Training / Fine-tuning</b>	~25–30% (declining from dominant)	Vendor customizes model for domain (no direct user trigger)	GPU hours, data prep, ML engineering	Lumpy, infrequent (quarterly or on-demand)	GPU hours × \$/hour × iterations per cycle	\$5K–\$50K per iteration; updates \$500–\$5K	Absorbed or professional services fee
<b>Batch Inference</b>	~5–8%	Nightly document processing, bulk categorization, content moderation	Tokens (batch-optimized), lower GPU priority	Quasi-linear, step-wise at batch windows	Documents/mo × avg tokens/doc × batch \$/token	30–70% cheaper than real-time (~\$0.001–\$0.02/query)	Per-document, per-job, or bundled credits
<b>Embedding / Indexing</b>	~2–3%	New documents uploaded to knowledge base; corpus re-indexing	GPU/CPU for embedding model, vector storage	Linear w/ data volume, spiky at re-index	Documents × tokens/doc × embedding \$/token + storage/mo	\$20–\$120 per 1B tokens; near-zero per incremental page	Absorbed or per-document fee
<b>Human-in-the-Loop</b>	~2–3% (labor component)	"Review and approve this AI-generated contract"; moderation queue	Human labor + AI inference for pre-labeling	Linear, predictable per item	Items/mo × (AI cost/item + human min/item × \$/min)	\$0.01–\$1.00 per item (AI); \$50–\$80/hr (human overhead)	Per-review seat or submission volume
<b>Multi-Model Routing</b>	N/A (cost reducer)	No user trigger; infrastructure routes simple queries to cheap models automatically	Orchestration logic (minimal)	Reduces inference cost by up to 85%	Savings = baseline inference × routing efficiency %	80% of queries to small models; 20% to frontier	Absorbed (margin protection lever)

**Key insight:** the true cost of a resolved AI task is often **10 to 50 times higher** than the posted per-call price when vector search, memory, concurrency, and moderation are included. A \$0.01 model call becomes a \$0.40 to \$0.70 workflow. OpenAI burned roughly **\$8.7 billion on Azure inference** in the first three quarters of 2025. That's not training. That's serving outputs. And 80% of engineering teams miss their AI infrastructure cost forecasts by more than 25%.

## THE ROI DIVIDE: SOFT VS. HARD

Bessemer draws a sharp line between **Soft ROI** (copilots that advise, hard to measure, high churn risk) and **Hard ROI** (agents that execute, concrete metrics, premium pricing power). Intercom Fin: \$0.99/resolution, 15% to 45% ticket share in 5 months. Chargeflow: 25% of recovered chargebacks. HighRadius: zero upfront fees, pay only on measurable outcomes. If your AI delivers measurable outcomes, your pricing should capture a share of that value.

[EXHIBIT: Soft ROI Trap vs Hard ROI Moat]

Soft ROI (Copilots)	Hard ROI (Agents / Services)
AI advises; human executes	AI executes; human supervises
Value: "better emails," "faster drafts" → hard to prove at renewal	Value: tickets resolved, revenue recovered, hours replaced → auditable
McKinsey: 14% issue resolution increase, 9% handling time reduction	Intercom Fin: 15% → 45% in 5 months at \$0.99/resolution; Chargeflow: 25% of recovered \$
Ex: Grammarly, Notion AI, GitHub Copilot (original flat-fee)	Ex: Intercom Fin, Chargeflow, EvenUp, HighRadius (zero upfront, pay on outcome)

## MATCHING THE CHARGE METRIC TO THE WORKLOAD

The model isn't something you pick from a menu. It's something you match to the level of autonomy your product delivers and the cost structure of each workload underneath it. Simon-Kucher's Price-Value Matrix: don't force outcome pricing on a copilot, don't sell an autonomous agent per-seat. **65% of vendors now use a hybrid approach; 85% of SaaS leaders adopted usage or hybrid pricing by 2025.** Zuora's COMPASS Framework overlays Scope of Work against Attributability to determine the correct metric.

[EXHIBIT: The Charge Metric Spectrum]

Model	Best For	Strengths	Risks	Real-World Example
Seat-Based	AI as feature enhancement; human does the work	Predictable revenue; familiar to procurement	Fails to capture non-linear AI value; -40% margins for AI products	Grammarly, ClickUp
Consumption / Token	Technical buyers; high-volume API/inference	Perfect cost alignment; protects margins	"Taxi-meter effect" discourages adoption	OpenAI API, Snowflake, Algolia
Credit / Hybrid	Balancing predictability with usage flexibility	Budget certainty + upside capture; dominant model (65%)	"Double conversion" overhead: credits → tasks → value	Adobe Firefly, GitHub Copilot Pro+, Clay
Outcome-Based / MASC	Full automation; AI closes the loop	Ultimate value alignment; vendor bets on own tech	Vendor absorbs cost variance; definition alignment friction	Intercom Fin (\$0.99/resolution), HighRadius (zero upfront)
Hourly / Digital Gig	AI as contract worker; direct labor comparison	Transparent ROI vs human (\$50–\$80/hr)	Unfamiliar; requires trust in agent reliability	Retool (hourly agent pricing)

## PACKAGING: FENCE FOR WILLINGNESS-TO-PAY, NOT FEATURES

Traditional SaaS packaging drew tier boundaries based on feature checklists: Basic gets 5 features, Pro gets 10, Enterprise gets 15. In AI, the fencing dimensions have shifted. **The question isn't what capabilities you can access. It's how well, how securely, and how independently the AI performs for you.** Every tier gets the AI. The fence determines the quality of execution, the level of trust, and the degree of autonomy. These aren't arbitrary upsell levers. Each fence maps directly to a different cost-to-serve: a frontier model costs 10x more to run than a basic one, a private instance costs more than shared infrastructure, and full autonomous execution triggers more compute steps than reactive prompting.

Some of these fencing dimensions are familiar from traditional SaaS (volume caps, support tiers). Others are genuinely new to AI because they carry real cost differences underneath that traditional features never had.

[EXHIBIT: AI Packaging Fence Dimensions]

Fence Dimension	How It Works	Why It Matters	New to AI?
Model Quality / Tier	Basic model (GPT-4o-mini class) at lower tiers; frontier model (GPT-4o, Claude Opus) at premium	Directly tied to inference cost: frontier models cost 10–50x more per token	Yes: no traditional SaaS equivalent
Autonomy Level	Reactive prompting at lower tiers; multi-step autonomous execution at premium	Most powerful fence: full autonomy triggers 10–100x more compute steps per task	Yes: gates highest-value capability where unit economics support it
Trust & Governance	Standard data handling at base; zero-retention, private instances, SSO, audit controls at premium	Enterprise procurement demands it; easiest fence to build (Anthropic, ServiceNow)	Expanded: data governance existed but AI amplifies sensitivity

<b>Specialization / Domain</b>	Generic agents at base; domain-specific agents (legal, finance, compliance) at premium	Real moat: proprietary domain data resists generic LLM replication	Yes: fine-tuned domain agents are a new cost and value layer
<b>Usage Volume / Throughput</b>	1,000 queries/mo at base; 50,000 at premium; unlimited at enterprise	Familiar lever but now directly tied to variable compute cost, not just access	Traditional: volume caps existed but cost was near-zero
<b>Speed / Latency Priority</b>	Standard queue at base; priority processing at premium	Real infrastructure cost difference: priority GPU allocation vs. batch queue	Expanded: SLA tiers existed but AI latency is GPU-bound
<b>Concurrency</b>	One agent at a time at base; parallel agents at premium	Each concurrent agent multiplies compute; directly impacts cost-to-serve	Yes: parallel AI execution is a new cost dimension
<b>Customization / Fine-tuning</b>	Generic model at base; model fine-tuned on customer's data at premium	High setup cost (\$5K-\$50K) justifies premium tier; creates lock-in	Yes: customer-specific model training is new

One critical friction point: transitioning legacy customers from unlimited tiers to consumption-metered ones. The 'taxi-meter effect' severely harms adoption when users link every action to a price increase. Smart vendors manage this with hybrid models that include generous base allowances and only charge overages past limits. Others subsidize initial AI usage for free, then spin features into stand-alone metered SKUs once the customer realizes value.

### THE COST VISIBILITY PROBLEM: WHEN DOES THE BUYER KNOW?

None of the pricing model or packaging discussion matters if the buyer can't predict their cost at the time of use. This is the gap most pricing articles ignore, and it's the dimension that determines whether customers adopt or throttle their usage.

*[EXHIBIT: Cost Visibility at Time of Use: Three Scenarios]*

Scenario	How It Works	Buyer Experience	Best For
<b>Tier-Selected (Known Upfront)</b>	Customer picks a plan. Plan determines model quality, security, autonomy. Monthly cost is fixed before any usage.	Full predictability. Buyers love it. But vendor absorbs all cost variance from power users.	Products where usage is relatively uniform; low agentic complexity
<b>Runtime-Determined (Unknown Until After)</b>	System routes queries based on complexity. Simple → cheap model. Complex → frontier. Cost determined after execution.	Buyer doesn't know cost until the bill arrives. Creates "taxi-meter effect" that kills adoption.	Technical buyers (API-first); high-volume infrastructure workloads
<b>Hybrid / Bounded (Floor Known, Ceiling Capped)</b>	Tier sets the ceiling of what you can access + base allowance. Overages are metered but visible in real-time with spend alerts and session caps.	Buyer knows their floor cost. Overages are bounded. Real-time visibility prevents surprises.	Dominant model (65% of vendors). Best balance of predictability and cost alignment.

**The unsolved problem: agentic workloads.** A customer sends what looks like a simple request: "research competitors and write a report." Underneath, the agent triggers 50+ steps: multiple LLM calls, tool use, web retrieval, verification loops, and state management. The customer had no idea at time of send what that would cost. This is why **guardrails, spend alerts, and session caps matter as much as the pricing model itself**. The model on paper means nothing if the buyer can't predict their cost at time of use. The vendors that solve cost visibility will win adoption. The ones that don't will watch customers throttle usage out of fear, regardless of how much value the AI delivers.

### THE STRATEGIC SCENARIO MATRIX (BAIN)

Scenario	Description	Pricing Implication	Mandate
<b>Core Strongholds</b>	Deep human judgment + proprietary data required	AI as productivity multiplier; price saved time at premium	Defend and enhance
<b>Open Doors</b>	Workflows accessible via open APIs; vulnerable to value siphoning	Launch proprietary agents to raise switching costs	Deepen integrations
<b>Gold Mines</b>	Vendor has exclusive domain data; defensible moat	End-to-end automation, outcome-based pricing	Go offensive
<b>Battlegrounds</b>	Highly repetitive, easily replicated tasks	Cannibalize own legacy SaaS before a disruptor does	Self-disrupt or die

## THE RENEWAL CLIFF IS COMING

Most AI deals closed in 2025 were subsidized. Those contracts are hitting renewal in 2026. The **double-cost transition gap** (AI agent + human salary during pilot) stalls enterprise sales. **85% of organizations misestimate AI project costs by more than 10%**; 80% of AI projects fail before production due to cost overruns. Definition alignment is the other landmine: outcome-based pricing only works where the vendor controls the outcome. The companies that survive will have defined success metrics before contract signature and built pricing structures with graduated adoption paths.

## THE BLUEPRINT

#	Move	What It Means
1	<b>Build Unit Economics Discipline</b>	Track inference, HITL, and API fees as direct COGS from day one. Model margins at current usage. The Jevons Paradox guarantees cheaper inference creates more usage, not more profit.
2	<b>Audit Workloads Individually</b>	Each AI component (inference, RAG, agentic, batch, HITL) needs its own commercial treatment: metering, guardrails, and pass-through vs. absorption. One model for the whole product is the mistake.
3	<b>Map Metric to Autonomy</b>	If AI augments, seat or consumption is fine. If it automates, charge for work completed. If it delivers financial outcomes, take a share. Match the metric to what the AI does.
4	<b>Cap the Downside</b>	Unlimited AI plans are financial liabilities. Hybrid thresholds, fair-use policies, soft caps. Multi-model routing cuts inference costs by up to 85%.
5	<b>Sell the Hard ROI</b>	Transition from 'software that helps you work' to 'digital labor that does the work.' Target labor budgets, not just software budgets. The TAM doubles when you do.
6	<b>Solve Cost Visibility</b>	If the buyer can't predict cost at time of use, they'll throttle adoption regardless of value. Build real-time spend alerts, session caps, and bounded overages. The pricing model on paper means nothing if cost is a surprise.

**The era of selling potential is over.** In 2026, software must earn its keep by delivering measurable work. Expect to spend **three dollars on change management for every dollar on technology**. The companies that win won't have the best AI models. They'll have figured out how to price and package what those models actually deliver.

## SOURCES & REFERENCES

1. The Economics of AI-First B2B SaaS in 2026 (Monetizely)
2. AI Is Driving A Shift Towards Outcome-Based Pricing, a16z (Dec 2024)
3. The AI Pricing and Monetization Playbook, Bessemer Venture Partners
4. Per-Seat Pricing Isn't Dead, but New Models Are Gaining Steam, Bain
5. From Seats to Consumption: Why SaaS Pricing Has Entered Its Hybrid Era
6. AI Agent Monetization: Lessons from the Real World, Stactize
7. How to Monetize Generative AI Features in SaaS, Simon-Kucher
8. Value Monetization in the Age of AI, Simon-Kucher
9. Evolving Models and Monetization Strategies in the New AI SaaS Era, McKinsey
10. Economic Potential of Generative AI, McKinsey
11. Zuora COMPASS Framework for Agentic AI Pricing
12. Bain SaaS Workflow Scenario Framework (via iMerge Advisors)
13. Deloitte: AI Compute Predictions 2026 (inference 2/3 of compute)
14. AI Inference Costs: 55% of Cloud Spending in 2026 (byteiota)
15. Inference Infrastructure Split: GTC 2026 (Rack2Cloud)
16. AI Cloud Infrastructure Spending 2026 (Market Clarity: \$37.5B total)
17. CloudZero: Guide to Inference Cost (\$0.01 call → \$0.40–\$0.70 workflow)
18. FinOps Foundation: Cost Estimation of AI Workloads
19. The SaaS Pricing Playbook: Models, Strategy & AI
20. OpenAI Azure inference spend: ~\$8.7B in 9 months (CloudZero/Microsoft)

21. Simon-Kucher: Agentic AI Price Metric Spectrum

22. Ibbaka Four-Layer Pricing Framework (HICSS)

---

### **About the Author**

**Massoud Ashrafi** is the Founder and Principal of Ashrafi Consulting, a pricing architecture and commercial strategy practice serving PE-backed and growth-stage technology companies. He advises on end-to-end pricing discipline including model design, packaging, governance, and monetization strategy. Previously: Amazon, Twilio, GoDaddy, PwC.

[ashraficonsulting.com](http://ashraficonsulting.com) | [linkedin.com/in/mashrafi](https://linkedin.com/in/mashrafi)