# The Cost of AI Agents
# — What Nobody Talks About

Everyone measures the subscription fee. Nobody measures the 44GB log file, the 20% sub-agent waste, or the 15,000 tokens loaded every heartbeat. This report does — with real production data from running AI agents 24/7.

February 2026

v1.0

Florian Ziesche · Ainary Ventures

CONTENTS

# 1. Executive Summary

**The sticker price of an AI agent is the smallest part of its cost. In production, hidden costs — context bloat, sub-agent waste, runaway logs, and cron job overhead — routinely exceed visible costs by 3–5x. Teams that don't measure these costs don't optimize them. Teams that do can cut spending by 90%.**

This report is based on **real production data** from running AI agents 24/7 on the Ainary platform — not benchmarks, not vendor projections. The numbers come from months of operating always-on agent infrastructure, tracking every token, every cron invocation, every log line.

## 90%

cost reduction achieved by optimizing cron jobs

Ainary production data, Feb 2026

## 44 GB

unplanned log file from a single agent system

Ainary production data

## 20%

of sub-agent output is waste (wrong, redundant, or unusable)

Ainary production data

- **Visible costs are misleading:** A $200/month subscription hides per-token costs that can be 10–20x higher than necessary when using frontier models for routine tasks

- **Context bloat is the silent killer:** Loading 15–20K tokens per heartbeat cycle means you're paying for the same context thousands of times per month

- **Multi-model routing is the highest-leverage optimization:** Routing routine tasks to Haiku-class models instead of Opus-class models yields 10–20x savings with minimal quality loss

- **The "Cron Tax" is real:** 24 scheduled agent jobs costing $100–270/month were reduced to 8 jobs at $10–15/month — not by cutting capability, but by cutting

unnecessary invocations

- **Cost Per Useful Output (CPUO) is the metric that matters:** Not cost per token, not cost per request — cost per output that actually achieves its intended purpose

---

*Keywords:* *AI agent costs, token economics, multi-model routing, context bloat, cron tax, cost per useful output, production AI operations, LLM cost optimization*

## 2. Visible Costs: The Number Everyone Knows

High

**Visible costs — subscriptions, API bills, compute — are what teams budget for. They represent roughly 20–30% of the true cost of running AI agents in production.**

The visible cost landscape in early 2026 is straightforward. Subscription tiers from major providers range from $20/month (basic) to $200/month (Pro Max tiers with higher rate limits and priority access). API pricing follows a per-token model: input tokens at $0.25–$15 per million tokens, output tokens at $1.25–$75 per million, depending on model tier. E

For an individual operator or small team, the calculus seems simple: **$200/month flat for unlimited use** (subscription) vs. **pay-per-token** (API). Most teams default to the subscription because it feels predictable. But "unlimited" is doing a lot of heavy lifting in that sentence — rate limits, context window caps, and model availability restrictions mean the subscription is not truly unlimited. It is unlimited within constraints that vendors adjust without notice. I

For agent workloads specifically, the visible cost picture in 2026 looks like this:

**Exhibit 1: Visible Cost Categories for AI Agents**

| COST CATEGORY | TYPICAL RANGE | WHAT IT COVERS |
|---|---|---|
| LLM subscription / API | $20–$200/mo (sub) or $50–$500/mo (API) | Model access, token consumption |
| Compute (hosting) | $5–$50/mo | Server/VM for agent runtime |
| Vector database | $0–$100/mo | Embeddings, RAG storage |
| Tool integrations | $0–$50/mo | APIs the agent calls (search, email, etc.) |

*Source: Industry pricing pages (OpenAI, Anthropic, various hosting providers), February 2026. Ranges reflect solo operator to small team scale.*

A reasonable estimate for visible costs: **$100–$400/month** for a solo operator running a moderately active agent system. This is the number that appears in spreadsheets, gets approved by managers, and shows up in "cost of AI" blog posts. It is also, as we will show, roughly one-quarter of the actual cost. I

The industry is moving fast on price compression. Token costs have dropped roughly 10x between early 2024 and early 2026 for equivalent capability levels. ByteDance research notes that the number of tokens consumed rises as the *square* of the number of rounds of API access by an agent — meaning cost scales quadratically with agent complexity, not linearly.[1] Cheaper tokens multiplied by quadratically more of them can still mean higher bills. E

**SO WHAT?**

If your AI cost tracking stops at the subscription fee or API bill, you're measuring the tip of the iceberg. Visible costs are necessary to track but radically insufficient as a cost model. The next section covers what's underwater.

# 3. Hidden Costs: The Number Nobody Tracks

85%

**The true cost of AI agents lives in five hidden categories that most teams never measure: context bloat, sub-agent waste, log accumulation, error-and-retry loops, and human oversight overhead.**

## 3.1 Context Bloat

Every agent invocation starts by loading context — system prompts, memory files, tool definitions, conversation history. In our production system, this means **15,000–20,000 tokens loaded per heartbeat cycle**. A heartbeat runs every few minutes. Do the math: at 20K tokens × 20 heartbeats/hour × 24 hours × 30 days, that's **~288 million input tokens per month** just for context loading — before the agent does anything useful. E

At Opus-tier pricing ($15/MTok input), that context loading alone costs **~$4,320/month**. At Haiku-tier pricing ($0.25/MTok), it costs **~$72/month**. Same context. Same information. 60x price difference. Most teams running subscription models never see this number because it's hidden inside "unlimited" access. But someone is paying for it — either you through rate limits and throttling, or the provider through margin compression. I

## 3.2 Sub-Agent Waste

When a primary agent spawns sub-agents for parallel work, **roughly 80% of sub-agent output is correct and useful. 20% is wrong, redundant, or unusable.** This is not a failure rate — it's a waste rate. The sub-agent completes its task, uses tokens to do it, but the output gets discarded or requires rework. E

In a system that spawns 10–20 sub-agents per day, at an average cost of $0.50–$2.00 per sub-agent invocation, the 20% waste rate translates to **$30–$240/month in pure waste**. More critically, sub-agent waste compounds: a sub-agent that produces wrong output may trigger a retry, which spawns another

sub-agent, which loads context again. The waste is not just the failed output — it's the cascading cost of recovery. I

### 3.3 Log Accumulation

Here is a number nobody plans for: **44 GB**. That is the size of a single log file generated by our agent system over several months of operation. Nobody budgeted for it. Nobody set up log rotation. The agent ran, logs accumulated, and eventually someone noticed the disk was full. E

At cloud storage prices ($0.02–$0.10/GB/month), 44GB costs $1–$4/month — trivial. But the indirect costs are not trivial: degraded system performance from disk pressure, time spent debugging why the agent was slow (answer: disk I/O bottleneck from massive log writes), and the engineering time to retroactively implement log rotation and archiving. The 44GB log file probably cost 4–8 hours of engineering time to diagnose and fix. At $150/hour loaded engineering cost, that's **$600–$1,200 in hidden labor cost** triggered by a log file nobody planned for. I

### 3.4 Error-and-Retry Loops

When an agent fails a task, it often retries — sometimes multiple times. Each retry loads context again (15–20K tokens), re-executes tool calls, and generates new output. In our production data, error-retry loops account for approximately **10–15% of total token consumption**. This is pure cost with zero marginal value — it's the tax you pay for agents that aren't reliable enough to succeed on the first attempt. I

### 3.5 Human Oversight Overhead

The least visible cost: human time spent reviewing, correcting, and babysitting agent output. If a sub-agent has a 20% error rate and a human spends 5 minutes reviewing each output, that oversight time adds up. For a system producing 50 outputs per day: 50 × 5 min = ~4 hours/day of human review. Not all of that is "waste" — verification has value. But the marginal cost of verifying the 80% that was already correct is pure overhead. J

**Exhibit 2: Hidden Cost Breakdown (Monthly, Single Agent System)**

| HIDDEN COST | LOW ESTIMATE | HIGH ESTIMATE | EVIDENCE |
|---|---|---|---|
| Context bloat (Opus-tier) | $2,000 | $4,500 | 15–20K tokens × heartbeat frequency |
| Sub-agent waste (20%) | $30 | $240 | Production error tracking |
| Error-retry loops | $50 | $300 | 10–15% of token consumption |
| Log/storage accumulation | $1 | $5 | 44GB observed |
| Engineering time (log/debug) | $100 | $1,200 | Incident-driven, amortized |
| Human oversight labor | $500 | $3,000 | Review time × loaded cost |
| **Total hidden costs** | **$2,681** | **$9,245** | |

*Source: Ainary production data, February 2026. Ranges reflect variation in agent activity level and model tier. Human oversight costs assume $50–$100/hour loaded labor cost.*

**THE REAL RATIO**

Visible costs: $100–$400/month. Hidden costs: $2,681–$9,245/month. **Hidden costs are 7–23x the visible cost.** Even if you halve our estimates, hidden costs still dominate. The implication: every cost optimization effort should start with hidden costs, not visible ones.

# 4. Subscription vs. API: The Real Math  78%

**Subscriptions feel cheaper because they're predictable. API feels expensive because it's visible. The math shows which one actually is cheaper depends entirely on your usage pattern — and most teams have never measured their usage pattern.**

Let's do the math with real numbers. A Pro Max subscription at **$200/month** gives you access to frontier models with generous (but not unlimited) usage. The equivalent API access, paying per token:

**Exhibit 3: Subscription vs. API Break-Even Analysis**

| USAGE PATTERN | TOKENS/MONTH (EST.) | API COST (OPUS-TIER) | API COST (HAIKU-TIER) | SUBSCRIPTION |
|---|---|---|---|---|
| Light (chat, occasional tasks) | ~5M input, ~2M output | ~$225 | ~$4 | $200 |
| Moderate (daily agent use) | ~50M input, ~10M output | ~$1,500 | ~$25 | $200 |
| Heavy (24/7 agent with heartbeats) | ~300M input, ~50M output | ~$8,250 | ~$138 | $200 |

*Source: Author calculations based on Anthropic API pricing (February 2026). Opus: $15/$75 per MTok in/out. Haiku: $0.25/$1.25 per MTok in/out. Subscription assumes Claude Pro Max at $200/month.*

The table reveals two counterintuitive insights:

**Insight 1: For heavy agent workloads on frontier models, subscriptions are a massive bargain.** A 24/7 agent consuming 300M+ input tokens per month would cost $8,250 at API rates for Opus. The $200 subscription saves $8,050/month —

a 97% discount. This is why subscription plans exist for power users: the provider is betting most subscribers use far less than the cap. ⓘ

**Insight 2: For routine tasks on small models, API is absurdly cheap.** The same heavy workload on Haiku-class models costs $138/month at API rates. Less than the subscription. And you get granular visibility into every token, enabling the kind of optimization we describe in this report. The subscription hides your usage; the API reveals it. ⓘ

The optimal strategy is not either/or — it's **subscription for frontier reasoning tasks + API for routine agent operations**. Use the $200 subscription when you need Opus-class intelligence (complex analysis, creative work, nuanced judgment). Use API with Haiku-class models for everything else (status checks, formatting, simple classification, log parsing). This hybrid approach can deliver frontier-quality results at near-Haiku costs. �J

> **SO WHAT?**
>
> If you're running agents on a subscription plan, you're probably getting good value on heavy usage — but you have zero visibility into waste. If you're on API only, you're probably overpaying on frontier models for tasks that don't need them. The answer is almost always both: subscription for the ceiling, API for the floor.

# 5. Multi-Model Routing: The 10–20x Lever   82%

**The single highest-leverage cost optimization for AI agents is not reducing usage — it's routing each task to the cheapest model that can handle it. For routine tasks, that means 10–20x savings with no perceptible quality loss.**

Not all tokens are created equal. A frontier model like Claude Opus or GPT-4-class costs **$15–$75 per million tokens**. A capable small model like Claude Haiku or GPT-4o-mini costs **$0.25–$1.25 per million tokens**. That's a **10–60x price difference** depending on the specific models compared.  E

The key question: **What percentage of agent tasks actually require frontier intelligence?**

In our production data, the answer is **roughly 15–25%**. The other 75–85% of agent invocations are routine: checking system status, formatting responses, classifying inputs, parsing structured data, running checklists, updating memory files. These tasks don't need a model that can write poetry or reason about abstract philosophy. They need a model that can reliably follow instructions and produce structured output.  E

Multi-model routing — automatically directing each task to the appropriate model tier — is the mechanism that captures this savings. The concept is simple:

- **Tier 1 (Frontier):** Complex reasoning, creative generation, nuanced judgment, multi-step analysis → Opus/GPT-4-class
- **Tier 2 (Mid-range):** Moderate complexity, synthesis, multi-document tasks → Sonnet/GPT-4o-class
- **Tier 3 (Routine):** Classification, formatting, status checks, simple tool calls → Haiku/GPT-4o-mini-class

If 80% of your agent's work can drop from Tier 1 to Tier 3, the blended cost drops dramatically. For a system consuming 100M tokens/month:

**Exhibit 4: Blended Cost With Multi-Model Routing**

| SCENARIO | MODEL MIX | MONTHLY COST | SAVINGS VS. ALL-FRONTIER |
|---|---|---|---|
| All frontier (no routing) | 100% Opus-tier | ~$3,000 | — |
| Basic routing | 50% Opus / 50% Haiku | ~$1,525 | 49% |
| Aggressive routing | 20% Opus / 80% Haiku | ~$620 | 79% |
| Optimized routing | 10% Opus / 20% Sonnet / 70% Haiku | ~$410 | 86% |

*Source: Author calculations. Assumes 100M input tokens/month, output proportional at 20%. Pricing: Opus $15/$75, Sonnet $3/$15, Haiku $0.25/$1.25 per MTok in/out.*

The pattern is clear: **even basic routing cuts costs in half. Aggressive routing cuts costs by 80%+.** The quality trade-off is minimal for routine tasks — Haiku-class models in 2026 are more capable than GPT-4 was at launch in 2023. The model that was "frontier" two years ago is now the "cheap" tier. I

Implementation is not technically complex. A routing layer that classifies incoming tasks by complexity (using — ironically — a Haiku-class model for the classification itself) and dispatches to the appropriate model can be built in a day. The harder part is organizational: teams need to accept that not every agent interaction needs the "best" model. Prestige bias toward frontier models is the most expensive form of status signaling in AI operations. J

**SO WHAT?**

If you're running every agent task on the same model, you're leaving 80%+ savings on the table. Multi-model routing is the single highest-ROI optimization available to agent operators today. Start with the simplest version: move all status checks and formatting tasks to Haiku-tier. Measure quality. Expand from there.

# 6. Case Study: 90% Cost Reduction  90%

*(Confidence: Very High — First-party production data)*

**We reduced agent operating costs by 90% — not through a technology breakthrough, but through systematic elimination of unnecessary invocations. The most expensive token is the one you didn't need to send.**

## The Starting State

Our agent system ran **24 scheduled cron jobs**, each triggering agent invocations at regular intervals. These jobs handled monitoring, memory management, content scheduling, status reporting, data synchronization, and various maintenance tasks. Monthly cost: **$100–$270** depending on activity level and model tier. The system worked, but nobody had audited whether all 24 jobs were necessary.  E

## The Audit

We examined each cron job against three criteria:

1. **Necessity:** Does this job produce an output that someone uses?
2. **Frequency:** Does it need to run this often?
3. **Model tier:** Does it need a frontier model?

The findings were sobering:

- **8 of 24 jobs** were producing outputs that nobody consumed. They ran because nobody had turned them off.
- **6 jobs** ran hourly when daily would have been sufficient. Same information, 24x the cost.
- **4 jobs** used Opus-tier models for tasks that Haiku-tier handled identically (status checks, file formatting).
- **6 jobs** were genuinely necessary, at their current frequency, at an appropriate model tier.

## The Optimization

**Exhibit 5: Cron Job Optimization Results**

| METRIC | BEFORE | AFTER | CHANGE |
|---|---|---|---|
| Active cron jobs | 24 | 8 | -67% |
| Monthly cost | $100–$270 | $10–$15 | -90 to -94% |
| Agent invocations/day | ~580 | ~45 | -92% |
| Useful output produced | Baseline | Same or better | No degradation |

*Source: Ainary production data, January–February 2026.*

**The result: 90% cost reduction with zero reduction in useful output.** We didn't make the agents cheaper — we stopped asking them to do things nobody needed. The savings came entirely from three levers: eliminating unnecessary jobs, reducing frequency of remaining jobs, and downgrading model tiers for routine tasks. E

## Why This Matters Beyond Our System

Every agent system accumulates cron jobs, scheduled tasks, and automated routines over time. They start useful. Requirements change. The jobs keep running. Nobody audits. In mature software engineering, this is called "technical debt." In AI agent operations, we propose calling it the **Cron Tax** — the ongoing cost of scheduled agent invocations that no longer serve their original purpose. J

**SO WHAT?**

Run a cron audit today. List every scheduled agent invocation. For each: who uses the output? Does it need this frequency? Does it need this model? We'd bet that at least a third of your scheduled jobs can be eliminated or downgraded. The 90% number is not theoretical — we achieved it.

# 7. The Cron Tax   76%

**The Cron Tax is the hidden cost of scheduled agent invocations that outlive their usefulness. It compounds silently because nobody reviews automated tasks once they're running.**

In traditional software, a cron job that runs unnecessarily wastes compute — usually pennies. In AI agent systems, every unnecessary cron job triggers a full agent invocation: context loading (15–20K tokens), reasoning, tool calls, output generation. **A single unnecessary hourly cron job on a frontier model can cost $50–$150/month.** Multiply by a dozen unnecessary jobs and you have a significant budget line that appears nowhere in anyone's planning.  I

The Cron Tax has a compounding property that makes it particularly insidious. Cron jobs are set up by engineers or operators to solve a specific problem at a specific time. The problem changes or disappears. The cron job persists. Over months, an agent system accumulates a sediment of obsolete automated tasks, each small enough to ignore individually, collectively large enough to dominate the cost structure.  J

## Taxonomy of Cron Tax Waste

From our audit, we identified four patterns of cron waste:

- **Orphaned jobs:** The consumer of the output no longer exists (e.g., a reporting job for a dashboard that was deprecated). These are pure waste — 100% of their cost is Cron Tax.
- **Over-frequency jobs:** Running hourly when daily would suffice. The information doesn't change fast enough to justify the invocation frequency. 90–95% of their cost is Cron Tax.
- **Over-powered jobs:** Using Opus-tier for tasks that Haiku-tier handles identically. The model selection was never revisited after initial setup. 80–90% of their cost is Cron Tax (the premium between tiers).

- **Redundant jobs:** Two or more jobs producing overlapping outputs. Usually created by different team members who didn't know the other job existed. 50–100% of one job's cost is Cron Tax.

**THE ACCUMULATION PROBLEM**

Agent systems that run for 6+ months without a cron audit typically accumulate 30–50% waste in their scheduled tasks. This is not a guess — it's what we found in our own system, and informal conversations with other agent operators suggest similar patterns. The Cron Tax is universal; the only variable is whether you've measured yours.

**The fix is organizational, not technical.** Add a monthly cron audit to your operational checklist. For each job, answer three questions: (1) Who consumes this output? (2) What's the minimum frequency? (3) What's the minimum model tier? If the answer to question 1 is "nobody" or "I don't know," kill the job. If you're wrong, someone will notice and ask for it back. They almost never do. J

# 8. Cost Per Useful Output: A New Framework

72%

**The industry measures cost per token. It should measure cost per useful output — the total cost (visible + hidden) divided by the number of agent outputs that actually achieved their intended purpose.**

Current AI cost metrics are input-oriented: cost per token, cost per request, cost per API call. These metrics tell you how much you're spending. They don't tell you how much you're wasting. **Cost Per Useful Output (CPUO)** is an output-oriented metric that connects spending to value. J

## The CPUO Formula

```
CPUO = (Visible Costs + Hidden Costs) / Number of Useful
Outputs
```

Where:

- **Visible Costs** = subscription + API + compute + tools
- **Hidden Costs** = context bloat + sub-agent waste + error-retries + log overhead + human oversight
- **Useful Outputs** = agent outputs that achieved their intended purpose without human rework

## CPUO in Practice

Let's apply CPUO to our pre-optimization and post-optimization states:

**Exhibit 6: CPUO Before and After Optimization**

| METRIC | BEFORE OPTIMIZATION | AFTER OPTIMIZATION |
|---|---|---|
| Total monthly cost (visible + hidden) | ~$370 | ~$25 |
| Total agent outputs/month | ~17,400 | ~1,350 |
| Useful outputs (80% rate) | ~13,920 | ~1,080 |
| **CPUO** | **$0.027** | **$0.023** |
| Waste outputs | ~3,480 | ~270 |
| Cost of waste | ~$74 | ~$5 |

*Source: Ainary production data, February 2026. "Useful" defined as output consumed without rework. 80% usefulness rate based on sub-agent accuracy tracking.*

The CPUO improved modestly (from $0.027 to $0.023) — but the **absolute waste dropped from $74/month to $5/month**, and total spend dropped 93%. The CPUO framework reveals that our system was already reasonably efficient on a per-output basis; the problem was producing too many outputs nobody needed. CPUO helps distinguish between two different cost problems: inefficient production (high CPUO) and unnecessary production (low CPUO but high total cost). I

## When CPUO Matters Most

CPUO is most valuable when comparing alternative architectures. Should you use one smart agent or five specialized sub-agents? The answer depends on which architecture produces more useful outputs per dollar — not which uses fewer tokens. A sub-agent system might use 3x the tokens but produce 5x the useful outputs, yielding a lower CPUO. Conversely, a single frontier model might seem expensive but waste less through higher first-attempt success rates. J

**SO WHAT?**

Start tracking three numbers: total cost (visible + hidden), total outputs, and useful outputs. Divide. That's your CPUO. Compare it across model tiers, architectures, and time periods. Optimize CPUO, not cost per token — because the cheapest token that produces useless output is infinitely expensive.

# 9. Recommendations  80%

Based on our production data and analysis, here are seven actionable recommendations ordered by impact and ease of implementation:

### 9.1 Run a Cron Audit This Week

**Impact: High | Effort: Low | Timeline: 1 day**

List every scheduled agent invocation. For each: who uses it, at what frequency, on what model? Kill orphaned jobs. Reduce frequencies. Downgrade model tiers. Expected result: 30–50% reduction in scheduled costs. We achieved 90%. J

### 9.2 Implement Multi-Model Routing

**Impact: Very High | Effort: Medium | Timeline: 1–2 weeks**

Route routine tasks (status checks, formatting, classification) to Haiku-tier models. Keep frontier models for complex reasoning only. Start with the simplest routing: if task type is in [list of routine types], use Haiku. Expected result: 50–80% reduction in token costs. J

### 9.3 Measure Context Loading

**Impact: High | Effort: Low | Timeline: 1 day**

Count how many tokens your agent loads as context per invocation. Multiply by invocations per month. This is your context tax. Then ask: does the agent need all that context every time? Can you lazy-load context only when relevant? Reducing context from 20K to 5K tokens per invocation is a 75% saving on input costs. J

### 9.4 Set Up Log Rotation Before You Need It

**Impact: Medium | Effort: Low | Timeline: 2 hours**

Implement log rotation and size limits from day one. Set alerts for log files exceeding 1GB. The cost of preventing a 44GB log situation is near zero; the cost of cleaning one up is $600–$1,200 in engineering time. J

### 9.5 Track Sub-Agent Success Rates

**Impact: Medium | Effort: Medium | Timeline: 1 week**

Instrument your sub-agent system to track: was the output used? Was it reworked? Was it discarded? Identify which sub-agent types have the highest waste rates and prioritize those for improvement (better prompts, tighter scope, or elimination). Our 80/20 split is a starting point — your ratios will differ. J

### 9.6 Adopt CPUO as Your Primary Cost Metric

**Impact: High | Effort: Low | Timeline: Ongoing**

Stop measuring cost per token. Start measuring cost per useful output. This reframes every cost discussion from "how do we use fewer tokens?" to "how do we produce more value per dollar?" The shift in framing changes which optimizations you prioritize. J

### 9.7 Use Hybrid Subscription + API

**Impact: Medium | Effort: Low | Timeline: 1 day**

Use subscription plans for interactive frontier-model work (where you'd exceed API costs). Use API for programmatic agent operations (where granular visibility and model routing matter). Don't default to one or the other — use each where it has an advantage. J

**Exhibit 7: Recommendation Priority Matrix**

| # | RECOMMENDATION | IMPACT | EFFORT | EXPECTED SAVINGS |
|---|---|---|---|---|
| 1 | Cron audit | High | 1 day | 30–90% |
| 2 | Multi-model routing | Very high | 1–2 weeks | 50–80% |
| 3 | Context measurement | High | 1 day | 25–75% |
| 4 | Log rotation | Medium | 2 hours | Incident prevention |
| 5 | Sub-agent tracking | Medium | 1 week | 10–20% |
| 6 | CPUO adoption | High | Ongoing | Decision quality |
| 7 | Hybrid sub + API | Medium | 1 day | Variable |

# 10. Methodology & Transparency

This report is based primarily on **first-party production data** from operating AI agents on the Ainary platform between October 2025 and February 2026. All cost figures, token counts, and optimization results are from our own systems unless otherwise attributed.

| | |
|---|---|
| **Data Source** | Ainary agent platform production logs, billing records, and operational metrics. Single-operator scale (not enterprise). |
| **Time Period** | October 2025 – February 2026 (5 months of continuous operation) |
| **Scale** | Solo operator / small team. Results may not generalize linearly to enterprise scale, where costs are higher but optimization leverage may differ. |
| **Model Pricing** | Based on published API pricing from Anthropic and OpenAI as of February 2026. Pricing changes frequently; specific numbers may be outdated by publication time. |
| **Industry Context** | Supplemented by web research on AI agent production costs, token pricing trends, and industry cost benchmarks (February 2026). |
| **Framework Originality** | The "Cost Per Useful Output" (CPUO) framework and "Cron Tax" concept are original to this report. Neither has been externally validated. |
| **Limitations** | Single-system data; results reflect our specific architecture and usage patterns. Hidden cost estimates include ranges to reflect uncertainty. Human oversight costs are particularly difficult to measure precisely. Sub-agent accuracy (80%) is an observed average — individual task types vary significantly. |
| **Claim Classification** | [E] Evidenced = backed by production data or external source. [I] Interpretation = reasoned inference. [J] Judgment |

= recommendation. [A] Assumption = stated but unproven.

| AI Involvement | This report was researched and drafted using AI agents — making us both the subject and the tool. We note this as both a feature (authentic first-party experience) and a limitation (potential bias toward systems we operate). |
| --- | --- |

## References

[1] ZDNET, "Why you'll pay more for AI in 2026, and 3 money-saving tips to try," January 2026. ByteDance research on quadratic token scaling in agent systems.

[2] ProductCrafters, "AI Agent Development Cost: $5K to $500K+ (2026 Pricing)," February 2026. Industry cost benchmarks including vector database hosting ($100–$2,000/month) and embedding generation costs.

[3] AgentFrameworkHub/MintSquare, "AI Agent Production Costs 2026: Real Data," January 2026. Budget recommendation of 5x expected token usage for agent workloads.

[4] OpenAI Pricing, api.openai.com/pricing, February 2026.

[5] Anthropic Pricing, anthropic.com/pricing, February 2026.

[6] PricePerToken.com, "LLM API Pricing 2026 — Compare 300+ AI Model Costs," February 2026.

**About the Author**

Florian Ziesche is the founder of Ainary Ventures, where he builds and operates AI agent infrastructure. This report reflects direct operational experience — every number in the case study came from systems he runs daily. He believes the most useful AI research comes from practitioners who ship, break, and fix things in production.

# Ainary

*Clarity in an age of agents.*

AI Agent Operations · Cost Optimization · Trust Infrastructure

Get in touch →