

BOOK THREE-A · THE AI ECONOMY MONETIZATION SERIES

The AI CFO Playbook

FinOps OS, revenue accounting, leakage, forecasting, and the CFO's mandate across SalesOps, RevOps, and FinOps

You cannot manage what you cannot measure. In AI, the measurement infrastructure must be built before the spending compounds.

The CFO who cannot see AI spending clearly cannot govern it. The CFO who cannot govern it cannot scale it. Build the visibility infrastructure before you need it.

Audience: CFOs, Controllers, FP&A leads, FinOps practitioners

PREFACE

The CFO's Mandate in an AI Business

Why AI changes the finance function's role — and what the new mandate requires.

There is a specific kind of discomfort that CFOs experience when AI deployment reaches a certain scale. It is not the discomfort of not understanding the technology — most finance leaders have made peace with not understanding transformers and diffusion models. It is the discomfort of not being able to see the economics clearly.

The AI spending is real. It appears in the cloud bills, the model API invoices, the infrastructure costs, the headcount. The question is whether it appears with enough precision to manage. Can the CFO attribute the AI spending to the business units

generating it? Can they forecast next quarter's AI costs with the same confidence they forecast marketing spend? Can they tell the board, with evidence, what the AI investment is returning? Can they identify, before the external auditor raises it, whether the variable consideration on outcome-based contracts is being estimated and constrained correctly?

For most CFOs managing significant AI deployments today, the honest answer to most of these questions is no. Not because the CFO is not capable of managing it — but because the financial infrastructure to manage it has not been built. The metering is not granular enough. The cost attribution is not precise enough. The revenue recognition policies have not been updated for AI-specific commercial models. The forecasting models are extrapolating SaaS patterns onto AI consumption that does not behave like SaaS.

This book is the CFO's guide to building that infrastructure. Not a guide to AI technology — this book does not explain how foundation models work. Not a guide to AI strategy — the strategic frameworks are in Books 0 and 1. This book is the guide to the financial infrastructure, governance mechanisms, accounting policies, and management reporting frameworks that allow a CFO to govern an AI-native business with the precision and confidence that the role requires.

The book is organized around the CFO's mandate in an AI business, which is broader than the traditional finance mandate. The CFO of an AI company is not just the guardian of the P&L and the keeper of accounting standards. They are the governing authority for AI spending across SalesOps, RevOps, and FinOps — the three operational functions whose activities directly affect whether AI investments are converted into commercial results. They are the architect of the token economy governance framework that prevents runaway AI spending. They are the accountable executive for the revenue recognition policies that determine whether AI revenue is recognized correctly and defensibly. And they are the board's primary source of informed judgment on whether the company's AI investments are generating returns commensurate with their cost.

This is a significant mandate. This book is the playbook for fulfilling it.

PART ONE

The CFO's Mandate

Why AI is the CFO's problem — and what governing it requires.

CHAPTER ONE

Why Monetization Is the CFO's Problem Now

AI spending is exploding. Revenue capture is inconsistent. The CFO as unifying force across SalesOps, RevOps, and FinOps.

The CFO of an AI company faces a problem that their predecessors in SaaS did not: the commercial mechanics of the business are more complex, more variable, and less visible than SaaS, at precisely the moment when the scale of the decisions being made — the AI infrastructure investments, the model API contracts, the enterprise AI deployments — demands the highest quality of financial governance.

Consider what the CFO needs to manage in an AI business that has no direct parallel in SaaS. Token consumption: millions of AI interactions per day, each consuming input and output tokens at a cost that varies with the model used, the prompt length, and the complexity of the task. The CFO's infrastructure bill is not a flat rate for servers running known workloads — it is a variable, model-dependent, consumption-driven cost that can spike dramatically when a new agent workflow is deployed or when a customer's usage pattern changes. Revenue variability: a SaaS business recognizes subscription revenue ratably and has good visibility into next quarter's revenue from the current ARR and historical churn rates. An AI business with significant consumption and outcome-based components has revenue that is genuinely variable — dependent on how much

customers use the AI, what outcomes it delivers, and how those outcomes are verified and recognized under ASC 606. The CFO's forecast model must accommodate this variability in a way that SaaS forecasting models do not. Agent governance: autonomous AI agents can consume millions of tokens, initiate API calls to external services, and incur costs at machine speed, without human approval of individual transactions. A single misconfigured agent workflow can generate a month's worth of token consumption in an afternoon. The CFO must establish the governance framework that prevents this — not after it happens, but before the agents are deployed.

These challenges converge on a single organizational reality: the CFO can no longer operate as a downstream recipient of commercial decisions made by sales, product, and engineering. The AI CFO must be upstream — involved in the pricing decisions that determine revenue structure, the product decisions that affect cost attribution, and the governance framework that prevents spending from outrunning visibility.

"The CFO of an AI company cannot operate as a downstream recipient of commercial decisions. They must be upstream — governing the pricing decisions that determine revenue structure and the deployment governance that prevents spending from outrunning visibility."

The Three-Function Model (Framework F11)

The three-function model (Framework F11) establishes the CFO's mandate in AI commercial operations with unusual precision. It defines three operational domains — SalesOps, RevOps, and FinOps — and assigns clear ownership for each.

SalesOps owns the offer: the commercial products, their pricing structures, and the deal terms through which they are sold. The SalesOps function creates the commercial objects that all downstream systems operate on. It is where the revenue trajectory of the business is fundamentally determined — not in the RevOps system that executes the sales motion, and not in the FinOps system that records the financial results, but in the offer design that defines what is sold and how.

RevOps owns the motion: the operational process of converting offers into executed contracts, activated entitlements, metered consumption, and issued invoices. RevOps is the commercial pipeline — the machinery that translates commercial intent into commercial reality. The quality of RevOps operations determines how much of the revenue potential defined by SalesOps is actually captured.

FinOps owns the result: the financial recording, governance, reporting, and forecasting of AI commercial activity. FinOps is where token consumption becomes an attributable cost, where consumption events become recognized revenue, and where the CFO's financial picture of the AI business is assembled.

The CFO's unique responsibility in this three-function model is integration: ensuring that the data definitions, the commercial objects, and the governance policies are consistent across all three functions. SalesOps creates a new pricing tier; RevOps must be able to bill it correctly; FinOps must be able to recognize the revenue and attribute the cost. When one function changes something without coordinating with the others, the downstream functions break in ways that are often invisible until quarter-close, when the reconciliation between what the billing system says and what the general ledger expects reveals the inconsistency.

The data-first principle — that the commercial data model must be defined before processes are designed and systems are selected — is the CFO's primary governance tool for maintaining three-function consistency. When the canonical object definitions (the thirteen monetization objects from Book 2a) are the shared reference point for all three functions, changes to pricing, billing, or accounting treatment can be evaluated for their

cross-function implications before they are implemented, not discovered after implementation when the inconsistency has already corrupted data.

The Three-Function Model — CFO's Role in Each Domain			
Function	What it owns	CFO's governance role	Failure mode without CFO governance
SalesOps	Commercial offers: products, pricing structures, deal terms	Review gate on new commercial structures — finance confirms rev rec treatment before launch	Products sold without defined rev rec policy; restatement risk at scale
RevOps	Commercial motion: Q2C pipeline, billing, entitlement activation	Data model consistency: canonical objects shared across RevOps and FinOps	RevOps and FinOps operate on different data models; reconciliation failures at period close
FinOps	Financial results: token governance, cost attribution, rev rec, reporting	Complete ownership: financial controls, governance infrastructure, reporting frameworks	No token governance; cost attribution gaps; revenue recognition not defensible

CFO GOVERNANCE GATE

The finance review gate prevents commercial structures from being sold without defined accounting treatment

Every new product tier, pricing model, or commercial structure must pass through a finance review gate before it is sold to any customer. The gate confirms three things: (1) the commercial terms can be precisely represented in the canonical data model; (2) the revenue recognition treatment has been determined and documented; (3) the cost attribution methodology is defined. Organizations that skip this gate discover, at audit time, that they have recognized revenue under policies they have not established — a finding that is expensive to remediate and potentially requires restatement.

Chapter One — The Essentials

- › The AI CFO must be upstream of commercial decisions — not a downstream recorder of their financial consequences.
- › The three-function model assigns governance responsibility: SalesOps (offer), RevOps (motion), FinOps (result).

- › The CFO's unique role is integration: ensuring data model consistency and policy alignment across all three functions.
- › The finance review gate prevents commercial structures from being sold before their accounting treatment is established.
- › Agent governance is a CFO responsibility: autonomous AI agents can incur costs at machine speed without human approval of individual transactions.

CHAPTER TWO

The Monetization Data Model: Objects the CFO Must Own

Which monetization objects finance must define, govern, and audit — and why each one affects financial statements.

The CFO must own the definition of specific monetization objects — not in the sense of managing them day-to-day, but in the sense of being the authority that ensures their definitions are precise enough to produce reliable financial statements, defensible revenue recognition, and accurate cost attribution.

The objects that finance must own are those whose definitions directly affect accounting treatment. The Product object's AI layer designation determines the revenue recognition treatment — a compute layer product is recognized differently from an outcome layer product. The Price object's variable consideration provisions determine the constraint analysis required under ASC 606. The Entitlement object's enforcement policy determines whether token budget exhaustion is a performance obligation delivery failure (potentially triggering a credit) or a consumption governance event (with no revenue recognition implication). The Allocation object is entirely a finance artifact — it exists to record the ASC 606 revenue allocation across performance obligations, and its configuration is a financial decision, not a commercial or technical one.

The CFO's practical task is to establish a formal review gate: before any new product, pricing tier, or commercial structure is launched, a finance review confirms that the commercial terms can be precisely represented in the data model, that the revenue recognition treatment has been determined, and that the cost attribution methodology is defined. This review gate prevents the most common financial infrastructure failure in AI companies: commercial structures that have been sold to customers before finance has established how they will be recognized.

Monetization Objects — Finance Ownership and Accounting Impact			
Object	Finance ownership requirement	Accounting treatment determined by	Audit risk if incorrectly defined
Product	AI layer designation must be confirmed by finance — determines rev rec treatment	Product type and AI layer determine whether recognition is ratable, consumption-based, or outcome-based	Wrong layer designation → wrong recognition pattern → potential restatement
Price	Variable consideration provisions require CFO sign-off on constraint methodology	Price type determines whether variable consideration analysis is required	Undisclosed or incorrectly constrained variable consideration → audit finding
Entitlement	Enforcement policy affects whether budget exhaustion is a rev rec event	Hard limit exhaustion may be a service failure with rev rec implications; soft limit is purely a governance event	Enforcement policy not documented → unclear accounting treatment for budget exhaustion events
Allocation	Finance owns entirely — this object exists solely for ASC 606 compliance	Allocation method and SSP determination determine when revenue from each component is recognized	Incorrect allocation → revenue recognized in wrong period → restatement risk
Contract	Performance obligations section must be defined with finance involvement	Performance obligation identification drives the entire ASC 606 analysis	Ambiguous performance obligations → inconsistent recognition → audit findings

Token Budget	Finance owns the budget hierarchy and enforcement policy design	Budget structure determines cost attribution methodology for management reporting	Budget hierarchy not aligned with cost attribution → management reports unreliable
--------------	---	---	--

Chapter Two — The Essentials

- › Six monetization objects require finance ownership: Product (AI layer), Price (variable consideration), Entitlement (enforcement policy), Allocation (entirely finance), Contract (performance obligations), Token Budget (budget hierarchy).
- › The Allocation object is entirely a finance artifact — it does not exist without deliberate finance design.
- › Product AI layer designation is a financial decision as much as a commercial one — it determines revenue recognition treatment.
- › The enforcement policy on Token Budget and Entitlement objects has accounting implications that must be evaluated before configuration.
- › Finance review of new commercial structures must happen before products are sold, not after recognition policies need to be established.

PART TWO

The AI FinOps OS

Token governance, cost attribution, and the financial visibility infrastructure.

CHAPTER THREE

The AI FinOps Imperative: Why Finance Is Flying Blind

The visibility gap. Token economics for finance. Building the infrastructure before you need it.

The FinOps imperative for AI is not a technology story. It is a financial governance story. And the problem it addresses is specific: finance leaders at most AI companies today are making spending decisions, budget approvals, and investment evaluations without visibility into the economics of their AI deployments at the granularity required for confident decision-making.

The visibility gap manifests in three recurring CFO experiences. The first is the surprise infrastructure bill: the AI team deploys a new agent workflow, the token consumption spikes, and the cloud bill arrives with a number that finance was not forecasting because the consumption data was not flowing to the financial planning function in real time. The second is the cost attribution question: the CFO is asked to allocate AI infrastructure costs to business units for divisional P&L reporting, and the answer is "we don't have that data" — because the event stream from the AI system does not carry team or cost center attribution. The third is the ROI question: the board asks whether the AI investment is paying off, and the CFO can describe the AI spending clearly but cannot describe the AI-generated revenue and value creation with the same precision.

All three experiences have the same root cause: the financial infrastructure for AI was not designed before the AI was deployed. Token consumption has been happening without attribution. Cost allocation rules have not been established. Value measurement has not been instrumented. The retrospective reconstruction of this financial picture — trying to build cost attribution after the fact, trying to establish ROI measurement for deployments that were not designed to be measured — is expensive, often incomplete, and sometimes impossible.

The principle that governs the FinOps approach in this book: build the financial visibility infrastructure before you deploy AI at scale, not after you need to explain the costs to the board.

⚠ The Cost of Deploying Without Financial Visibility Infrastructure

The most expensive FinOps implementation is the retroactive one. Organizations that deploy AI at scale without token attribution, cost allocation rules, and consumption governance must

eventually reconstruct the financial picture they did not capture in real time. This reconstruction is expensive (engineering time to query historical systems that were not designed for this analysis), incomplete (historical data may not carry the attribution fields needed for full cost allocation), and often produces financial statements that require restatement because the revenue recognition treatment was not established before recognition occurred. Build the financial visibility infrastructure before the first significant AI deployment, not after the first board question about AI ROI.

Token and Dollar Dashboards

Token and dollar dashboards — the real-time financial visibility tools that allow the CFO to see AI spending by team, agent, workflow, and model — are the foundational FinOps infrastructure. They are not analytics products built on top of the AI deployment. They are financial controls that must be in place before significant AI spending begins.

The token dashboard has four required views that correspond to the four financial management questions a CFO needs to answer about AI spending.

The organizational view shows total token consumption for the organization by day, week, and month, against the approved AI budget. This view answers the question every CFO needs to answer first: are we within budget? It shows the total token burn rate, the projected spend for the remainder of the period at current rates, and the comparison against the approved budget. Any organization without this view is operating its AI deployment on financial faith — hoping the costs will be within budget without any mechanism to detect when they are not.

The business unit view shows token consumption broken down by the teams or departments that are generating it. This view answers the cost attribution question: who is spending what? For organizations implementing AI cost chargeback, the business unit view is the financial basis for the internal billing. For organizations running showback, it is the visibility mechanism that drives behavioral accountability without internal financial charges. The business unit view requires that every token consumption event

carry a cost center or team identifier — an architectural requirement that must be built into the AI deployment, not retrofitted after the fact.

The model view shows token consumption broken down by the foundation models being used. This view answers the cost efficiency question: are we using expensive models for work that cheaper models could do? Different models have dramatically different cost profiles — a frontier model like GPT-4o or Claude 3.5 Sonnet might cost 10–20 times more per token than a smaller, faster model like GPT-4o-mini or Claude 3.5 Haiku. The model view identifies workflows that are consuming the most expensive models and surfaces the optimization opportunity of routing lower-complexity tasks to cheaper models.

The workflow view shows token consumption broken down by the specific AI workflows generating it. This view answers the ROI question at the granularity where it can actually be evaluated: is this specific workflow generating value proportional to its cost? A workflow that consumes 10 million tokens per month and has demonstrably improved the quality and speed of contract reviews is generating a measurable return. A workflow that consumes 5 million tokens per month and has produced no detectable improvement in the output it was designed to enhance is a candidate for redesign or elimination.

Token Dashboard — Four Required Views				
View	Financial question answered	Required data inputs	Update cadence	CFO action trigger
Organizational view	Are we within our approved AI budget?	Total token consumption · Approved AI budget · Current spend rate · Projection	Daily	Projected to exceed budget before period end
Business unit view	Who is spending what — cost attribution for P&L and chargeback	Token consumption by cost center / team_id · Cost rate table	Daily	Any BU above 90% of allocated budget · Chargeback disputes

Model view	Are we using expensive models for work cheaper models could do?	Token consumption by model_id · Cost per token per model	Weekly	High-cost model usage growing faster than value delivered
Workflow view	Is this specific workflow generating ROI proportional to its cost?	Token consumption by workflow_id · Outcome or value data for each workflow	Weekly	Workflow cost growing without corresponding value evidence

The Dollar Dashboard

The dollar dashboard translates the token-denominated consumption data into the financial terms that the CFO's budget management and board reporting require.

The key metric that the dollar dashboard must provide is real-time cost attribution: the dollar value of token consumption attributed to each business unit, workflow, and model, updated with a latency of no more than 24 hours. Real-time attribution requires two inputs: the token consumption data from the metering system (granular, timestamped, attributed to cost center) and the cost rate table from the FinOps system (the dollar cost per million tokens for each model, updated when model prices change).

The cost rate table is a financial governance document that must be maintained by the CFO's team. It records the cost per token for each model in use, the effective date of each rate (model prices change frequently), and the accounting treatment for model cost (infrastructure cost of goods sold versus technology and development expense, depending on the nature of the deployment). When model prices change — as they do frequently as the market evolves — the cost rate table must be updated, and the rate change must be applied to future consumption without retroactively recomputing historical cost attribution.

The dollar dashboard's most important derived metric is the chargeback versus showback comparison: the difference between what each business unit would be charged in a full chargeback model (their attributed token cost) versus what they are

actually being charged (which may be less if the organization is running a partial chargeback or showback model). This comparison surfaces the subsidy that some business units are receiving from others — information that the CFO needs to evaluate whether the current allocation model is producing the right behavioral incentives.

Dollar Dashboard — Cost Rate Table Structure			
Field	Description	Maintenance requirement	Accounting treatment note
model_id	Unique identifier for the foundation model being priced	One row per model, updated when new models are deployed	Model costs are COGS for production deployments; R&D expense for experimental
cost_per_million_input_tokens	Dollar cost per million input tokens at current model pricing	Updated within 24h of any model price change	Apply prospectively from effective date; do not retrocompute historical costs
cost_per_million_output_tokens	Dollar cost per million output tokens at current pricing	Updated within 24h of any model price change	Typically 2–5× input token cost — document the ratio for audit
effective_from	Date from which this rate applies	Historical rates must be preserved for period audit trails	Historical costs computed at the rate applicable at time of consumption
accounting_category	COGS · Technology expense · R&D · Allocated overhead	Finance determines at time of model deployment based on use case	Determines P&L line and cash flow treatment for AI infrastructure spend

FOR THE FINOPS LEAD

The chargeback vs showback decision is a behavioral governance choice, not just an accounting choice

Chargeback creates hard financial accountability: business units receive internal charges for their AI consumption and must manage against an allocated budget. This drives the strongest behavioral accountability but creates organizational friction and requires internal billing

infrastructure. Showback creates soft accountability: business units see their AI consumption costs without being charged. This drives awareness without the friction of internal billing. The right choice depends on the organization's maturity: showback first, then chargeback as governance disciplines are established. Neither is correct for all organizations at all stages.

Chapter Three — The Essentials

- › The visibility gap — finance managing significant AI spend without real-time attribution — is the starting point for every FinOps programme.
- › Four token dashboard views answer the four CFO questions: within budget? who is spending? using the right models? is this workflow generating ROI?
- › The cost rate table is a financial governance document maintained by the CFO's team — it must be updated within 24h of any model price change.
- › The chargeback vs showback decision is a behavioral governance choice — showback first, chargeback as governance maturity grows.
- › The most expensive FinOps implementation is the retroactive one — build visibility infrastructure before the first significant AI deployment.

CHAPTER FOUR

Agent Cost Attribution: Who Spent What and Why

Cost attribution for autonomous agents. Multi-agent workflow cost tracing. The attribution chain.

Agent cost attribution is the most technically challenging FinOps problem in AI finance — and the one that most organizations are solving least well. The challenge is structural: an AI agent executing an autonomous workflow may make dozens of model calls, consume hundreds of thousands of tokens, call external APIs, and trigger downstream processes, all within a single workflow execution. Attributing the cost of this activity to a specific customer, business unit, project, or cost center requires that the attribution chain be preserved across every step of the agent's activity.

The naive approach — attribute all agent costs to the cost center that deployed the agent — produces distorted financial reporting when a single agent serves multiple business functions. A research agent deployed by the strategy team but used by the product team, the marketing team, and the CEO's office is not purely a strategy cost. The correct attribution requires that each workflow execution carry the identity of the function that initiated it, which requires that agent invocations carry a cost center or project identifier in the request metadata.

The more sophisticated challenge is multi-agent workflow cost attribution. When an orchestrating agent delegates tasks to specialized sub-agents — a research orchestrator dispatching web search agents, summarization agents, and analysis agents — the total cost of the parent workflow includes the costs of all child agent invocations. If the child agents are shared infrastructure used by multiple workflows, their costs must be allocated to the parent workflows that invoked them in proportion to the resources each invocation consumed. This allocation chain can extend several levels deep in complex multi-agent deployments, producing a cost tree that must be traversed to arrive at the true cost of any specific workflow execution.

The practical FinOps implementation for multi-agent cost attribution uses the `workflow_id` and `session_id` fields from the Event object hierarchy. Every event generated by any agent in a workflow tree carries both the immediate agent's ID and the root workflow's ID. Aggregating costs by root `workflow_id` produces the total cost of each workflow execution, regardless of how many sub-agents were involved. This aggregation is the foundation for workflow-level ROI analysis: total workflow cost versus total value delivered by the workflow.

Token Budget Governance Architecture

Token budget governance is the CFO's primary financial control for AI spending. It is the mechanism through which the CFO translates an approved AI budget into enforceable spending limits that prevent the budget from being exceeded by autonomous agent activity.

The token budget hierarchy corresponds to the organizational structure of the AI deployment. At the top is the organizational budget: the total AI spending approved for the period, expressed in both token terms (the maximum tokens that can be consumed across all AI activities) and dollar terms (the financial equivalent at current model rates). Below the organizational budget are functional budgets for each major business unit — engineering, sales, customer success, finance, legal. Below each functional budget are project or workflow budgets for specific AI initiatives. The hierarchy enforces top-down spending controls: a team cannot exceed its functional budget regardless of individual workflow budgets, and the organization cannot exceed its total AI budget regardless of functional budgets.

The enforcement policy for each budget level must be explicitly configured, not defaulted. The organizational budget typically uses a soft limit with approval: consumption that approaches the limit triggers an alert and initiates a CFO review, but does not immediately halt AI activity (because halting all AI activity organizationally is typically too disruptive to be the correct response to a budget approaching its limit). Functional budgets for established AI deployments may use the same soft limit with approval approach. Functional budgets for experimental deployments may use a hard limit with a defined approval process for increases: the team cannot exceed their allocation without a formal budget amendment request.

The CFO override protocol — the mechanism through which the CFO can approve emergency budget increases in real time when a business-critical AI workflow is at risk of being throttled — is a governance mechanism that must be designed before it is needed. An override protocol that requires a committee meeting is not an override protocol — it is a process that will be bypassed when the situation is urgent. The effective CFO override is a configured authority mechanism in the Token Budget system: the CFO (and only the CFO) can approve an immediate budget increase up to a defined amount with a recorded rationale, with the approval logged in the governance audit trail and reviewed in the next scheduled FinOps review.

Token Budget Hierarchy — Governance Architecture				
Level	Budget entity	Enforcement policy	Alert threshold	Override authority
Level 1 (Organization)	Total organizational AI spend approved for period	Soft limit with CFO approval — halting all AI activity is too disruptive	85% consumed	CFO only — emergency increase up to 20% without board approval
Level 2 (Function)	Engineering · Sales · CS · Finance · Legal · Marketing budgets	Soft limit with VP approval for established teams; hard limit for experimental teams	80% consumed	Function VP up to 10% increase; CFO for larger increases
Level 3 (Project/Workflow)	Individual AI initiative or workflow budgets	Hard limit for experimental workflows; soft limit for production workflows	75% consumed	Project lead requests; function VP approves; no override without formal approval
Level 4 (Agent)	Per-agent consumption limit for autonomous workflows	Hard limit — no override without human approval — prevents runaway agent spend	70% consumed	Agent operator submits increase request; automated routing to project lead

Budget Governance Controls — By Level				
Control	Level 1 (Org)	Level 2 (Function)	Level 3 (Project)	Level 4 (Agent)
Enforcement policy	Soft limit with CFO approval	Soft limit (established) / Hard limit (experimental)	Hard limit (experimental) / Soft limit (production)	Hard limit — no exceptions without human approval
Alert threshold	85% consumed	80% consumed	75% consumed	70% consumed
Override authority	CFO only — documented in audit trail	Function VP + CFO notification	Project lead request +	Agent operator request + Project lead approval

			Function VP approval	
Review cadence	Weekly CFO FinOps review	Monthly function budget review	Weekly for active projects	Real-time alert; human review within 4 hours
Audit documentation	CFO override log + weekly FinOps report	Function budget utilization report	Project consumption report	Agent activity log with all invocations

CFO GOVERNANCE DESIGN

The CFO override protocol must be designed before the first budget alert fires

The CFO override protocol — the mechanism through which the CFO approves emergency budget increases in real time — must be designed, implemented, and tested before any budget alert is triggered. An override protocol that exists only in a policy document is useless when a business-critical AI workflow is about to be throttled at 3pm on the last day of the quarter. The effective override is a configured authority mechanism: the CFO can approve a specific budget increase amount with a single authenticated action, the approval is logged automatically, and the budget is increased in real time. Design it before you need it.

Chapter Four — The Essentials

- › Multi-agent workflow cost attribution requires root workflow_id aggregation across all child agent invocations.
- › The token budget hierarchy (org → function → project → agent) translates the approved AI budget into enforceable spending limits.
- › Every level of the budget hierarchy must have an explicit enforcement policy — do not default to whatever the billing platform provides.
- › The CFO override protocol must be designed and implemented before the first budget alert fires — not improvised under pressure.
- › Agent-level hard limits are non-negotiable: autonomous agents can incur costs at machine speed, and soft limits for agents create runaway spending risk.

PART THREE

Revenue Accounting

ASC 606 applied to AI services. Variable consideration. Multi-party allocation.

CHAPTER FIVE

ASC 606 and IFRS 15 Applied to AI Services

Performance obligations. Application to token, agent, and outcome billing. The five-step framework for AI commercial structures.

ASC 606 and IFRS 15 were designed to create consistent revenue recognition across industries by requiring that revenue be recognized when (and only when) a performance obligation is satisfied. Applied to AI services, this principle produces specific accounting treatments for each of the five AI economy layers — treatments that differ materially from each other and from the ratable recognition that governs most SaaS subscription revenue.

The five-step ASC 606 framework applied to AI services requires the CFO to answer five questions for every new AI commercial structure before revenue recognition policies are established.

Step one: identify the contract with the customer. For AI services, this is typically straightforward — a signed enterprise agreement, a cloud marketplace private offer, or a self-service subscription agreement. The complexity arises in agent-to-agent commerce, where the "contract" may be established programmatically without a signed document. The CFO must establish whether programmatic commercial agreements constitute binding contracts under applicable law, and whether the programmatic acceptance protocol generates the documentation required for an auditable contract record.

Step two: identify the performance obligations in the contract. This is where AI commercial complexity creates genuine accounting challenges. A hybrid AI contract that includes a platform subscription (access to the AI), a token allocation (a defined volume of AI consumption), and outcome delivery commitments (an agreed number of outcomes to be delivered) contains three distinct performance obligations — or potentially more, if individual product components are separately identifiable and distinct. The allocation of transaction price across these obligations directly affects revenue recognition timing, because each obligation may be satisfied at a different time.

Step three: determine the transaction price. For AI contracts with consumption-based or outcome-based components, the transaction price includes variable consideration — amounts that depend on future events (how much the customer consumes, how many outcomes the AI delivers). The CFO must estimate the variable consideration using either the expected value method (the probability-weighted average of all possible outcomes) or the most likely amount method (the most probable single amount). This estimate must be updated at each reporting period to reflect new information about actual consumption and outcome delivery rates.

Step four: allocate the transaction price. The transaction price must be allocated to each performance obligation based on the relative standalone selling price (SSP) of each obligation. For AI products where the standalone selling price of a specific component is not directly observable (because it is never sold separately), the CFO must estimate the SSP using one of the permitted estimation approaches: adjusted market assessment, expected cost plus margin, or residual approach.

Step five: recognize revenue when performance obligations are satisfied. For AI subscription access, performance obligations are typically satisfied ratably over the service period — revenue is recognized evenly over the subscription term. For AI token consumption, performance obligations are satisfied as tokens are consumed — revenue is recognized in proportion to consumption. For AI outcome delivery, performance obligations are satisfied when outcomes are verified — revenue is recognized at the point of verified outcome delivery. The practical implication: a hybrid AI contract requires

three recognition patterns applied simultaneously to three performance obligation components.

ASC 606 Five-Step Framework — Applied to AI Services			
Step	Standard requirement	AI commercial application	Key accounting judgment
Step 1 Identify contract	Binding commercial agreement exists	Enterprise contract, marketplace offer, self-serve subscription, or programmatic A2A agreement	Programmatic A2A agreements: do they constitute binding contracts under applicable law?
Step 2 Identify obligations	Distinct promises to deliver goods or services	Subscription access + token allocation + outcome delivery = up to 3 distinct obligations in a hybrid contract	Are components distinct? Is the customer's ability to benefit separately from each component?
Step 3 Determine price	Transaction price including variable consideration estimate	Base subscription + estimated consumption overage + estimated outcomes × outcome price	Variable consideration constraint: is it highly probable a significant reversal won't occur?
Step 4 Allocate price	Allocate to each obligation based on relative SSP	SSP for each component using adjusted market, cost-plus, or residual approach	SSP estimation when components not sold separately — document methodology carefully
Step 5 Recognize revenue	Recognize when/as obligation satisfied	Access: ratably · Consumption: as consumed · Outcomes: at verification	Three simultaneous recognition patterns for a single hybrid contract

"A hybrid AI contract with subscription, consumption, and outcome components requires three simultaneous revenue recognition patterns. This is not accounting complexity for its own sake — it is accounting precision for what the commercial structure actually is."

Variable Consideration

Variable consideration is the most significant revenue recognition complexity in AI commercial operations. It arises whenever the transaction price depends on future performance — specifically, whenever consumption-based or outcome-based pricing components are part of the commercial structure.

The constraint analysis for variable consideration requires the CFO to assess whether it is highly probable that including a variable consideration estimate in recognized revenue will not result in a significant revenue reversal when the uncertainty is subsequently resolved. This assessment must be documented and updated at each reporting date.

The constraint analysis for token consumption revenue is typically straightforward. Consumption-based token revenue has variability that is bounded by the contract's minimum commitment (the floor below which actual revenue cannot fall) and the customer's demonstrated consumption capacity (the ceiling above which actual consumption is unlikely to reach without a change in the customer's deployment). For deployments with stable consumption histories, the constraint is low and a high proportion of expected variable consideration can be recognized. For new deployments or volatile consumption patterns, the constraint is higher and a more conservative estimate is appropriate.

The constraint analysis for outcome-based revenue is more demanding. The variability in outcome-based revenue depends on: the AI's historical performance rate on the specific outcome type (more historical data reduces uncertainty), the measurement methodology (automated verification from customer systems is more reliable than manual verification, which may be contested), the attribution methodology (outcomes where causation by the AI is unambiguous are less variable than outcomes where attribution is contested), and the customer's operational stability (customers with stable workflows have more predictable outcome delivery rates than customers in organizational transition).

The practical implementation: the CFO must establish a variable consideration estimation policy for each product category that specifies the estimation method, the constraint analysis criteria, the evidence required to move between constraint levels, and the documentation required for audit. This policy must be reviewed at each major contract signing and updated at each quarterly close.

Variable Consideration Constraint Analysis — AI Revenue Categories						
Revenue category	Variability drivers	Low criteria (recognize more)	constraint (recognize less)	High criteria (recognize less)	constraint (recognize less)	Documentation required
Token consumption overage	Consumption rate variability; model changes; customer behavior	12+ months history; low variance; stable customer deployment; no significant model changes pending		New deployment; high consumption variance; customer in organizational transition		Consumption history; variance analysis; deployment stability assessment
Outcome-based components	AI performance rate; attribution methodology; measurement reliability	≥12 months outcome data; automated customer-system verification; unambiguous attribution; variance <10%		New deployment; manual verification; contested attribution; high variance		Performance history; verification methodology; attribution analysis; constraint calculation
Gain-share components	Business outcome measurement; attribution to AI vs other factors	Agreed measurement methodology; 3+ quarters of data; strong causal link to AI		Novel measurement approach; multiple contributing factors; measurement under development		Measurement methodology; causal analysis; gain calculation examples; baseline establishment
Marketplace royalties	Platform transaction volume; model usage by marketplace customers	Stable marketplace; predictable model usage; low customer concentration		New marketplace; high customer concentration; usage dependent on single large customer		Transaction volume history; customer concentration analysis; marketplace

				stability assessment
--	--	--	--	-------------------------

Multi-Party Revenue Allocation

Multi-party revenue allocation is required when AI revenue flows through arrangements involving more than one entity — model providers receiving royalties, channel partners receiving commissions, cloud marketplace operators taking a revenue share, or ecosystem participants sharing in AI-generated value.

The principal-versus-agent analysis governs whether the AI vendor recognizes revenue gross (as a principal controlling the service) or net (as an agent facilitating a transaction between the customer and the underlying service provider). For AI platform operators that aggregate model providers and application developers, the principal-versus-agent determination can be non-obvious and may vary by transaction type. The CFO must document the analysis and apply it consistently — switching between gross and net recognition for similar transactions without a substantive commercial change is an accounting quality issue that auditors will examine.

Commission accounting for channel partners under ASC 340-40 requires the CFO to determine whether channel partner commissions are incremental costs of obtaining a contract that should be capitalized and amortized over the contract term, or period costs that should be expensed as incurred. The determination depends on whether the commission would not have been incurred absent the customer contract — incremental commissions meet the standard for capitalization. The amortization period for capitalized commissions must be determined with reference to the expected customer relationship duration, including renewal periods that are reasonably certain to be exercised.

Model royalty accounting requires specific attention to the timing of royalty recognition in the context of the platform's revenue recognition. If the platform recognizes revenue on a consumption basis (as tokens are consumed by customers), and the model royalty is a percentage of that consumption revenue, the royalty expense should accrue in the

same period as the related revenue. If there is a timing difference — for example, if royalties are settled quarterly but revenue is recognized monthly — the CFO must ensure that a royalty accrual is recorded in the periods before settlement.

Multi-Party Revenue Arrangements — Accounting Treatment				
Arrangement type	Principal vs agent analysis	Revenue recognition	Commission/royalty treatment	
Direct AI vendor → Customer	Always principal — vendor controls the service and is responsible for delivery	Gross revenue recognized; COGS includes model inference and infrastructure costs	N/A — no multi-party arrangement	
Marketplace operator → Model providers	Principal if marketplace controls service experience; Agent if marketplace merely facilitates transaction	Principal: gross; Agent: net (take rate only)	Model royalties: COGS if principal; netted against revenue if agent	
Channel partner resale	Vendor is principal; partner is agent for the vendor	Vendor recognizes gross revenue; partner commission is Selling expense	Commission: capitalize if incremental cost of obtaining contract (ASC 340-40); expense if not	
Gain-share arrangement	Vendor is principal for AI service; customer and vendor share outcome upside	Vendor recognizes AI service revenue gross; gain-share is additional variable consideration	Customer's gain share: reduce COGS (if gain returned) or variable consideration (if recognized at higher rate)	

Chapter Five — The Essentials

- › The ASC 606 five-step analysis must be completed for every new AI commercial structure before revenue recognition policy is established.
- › Hybrid AI contracts require three simultaneous recognition patterns: ratable (subscription), consumption-based (tokens), point-in-time (outcomes).
- › Variable consideration constraint analysis must be documented at each reporting period — the estimate must be updated as actual performance data accumulates.

- › Principal-versus-agent determination for AI marketplace arrangements requires careful analysis — the same transaction structure can reach different conclusions based on control indicators.
- › Channel partner commissions under ASC 340-40: capitalize if incremental and recoverable; expense if not — document the determination and apply consistently.

PART FOUR

Leakage and Forecasting

Where AI companies lose money silently — and how to predict revenue accurately in a consumption business.

CHAPTER SIX

Revenue Leakage: Where AI Companies Lose Money Silently

The leakage map. Metering blind spots. Entitlement drift. Quantifying leakage. The audit methodology.

Revenue leakage in an AI business is the systematic failure to capture revenue for value that has been delivered. Unlike revenue that is invoiced and not collected (which appears in the AR aging report), revenue leakage never appears in the billing system because the revenue was never captured in the first place. It exists in the gap between what the AI delivered and what the billing system charged for.

The leakage map for an AI business identifies the specific points in the commercial pipeline where revenue is most likely to escape without detection. There are seven primary leakage points, each with a specific root cause and a specific remediation approach.

Metering blind spots are the most significant and the most invisible. Events that are generated by the AI system but not successfully written to the event store are consumed value that was never recorded. The most common cause is the event submission architecture — if the AI product submits events asynchronously with no retry logic, network failures create permanent event loss. A metering system that loses 2% of events creates 2% revenue leakage on all consumption billing — invisible on any single invoice, but compounding at scale.

Entitlement drift occurs when a customer's actual consumption entitlement — the commercial terms governing what they are allowed to consume — diverges from the entitlement configuration in the billing system. The most common cause is contract amendments that update the commercial terms without updating the billing configuration. A customer who negotiated a 20% token budget increase at renewal but whose billing system entitlement was not updated is consuming against their new entitlement without being billed for the additional consumption.

Attribution failures are events that reach the metering system but cannot be attributed to a billable commercial context — a customer, a product, an entitlement. Events without valid attribution are placed in the unattributed queue and may remain there indefinitely if the operations team does not actively manage the queue. Each event in the unattributed queue is potential revenue that has not been captured.

Outcome verification gaps occur in outcome-based billing when outcomes are delivered and verified but do not trigger a billing event. The most common cause is a broken integration between the outcome verification system and the billing engine — the verification event fires in the customer's system of record, but the webhook to the billing system fails and is not retried.

Underbilling on overages occurs when a customer exceeds their contracted consumption allocation but the overage pricing is not correctly applied. The most common cause is a billing configuration error — the overage pricing tier was not configured correctly, or the reset date for the allocation counter was set incorrectly, causing the overage calculation to apply the wrong pricing.

Model upgrade undercharging occurs when a customer's AI deployment is upgraded to a more expensive model without a corresponding update to the billing configuration. The customer receives a better AI product at the price they were paying for the previous, cheaper model.

The unrealized expansion gap is not technically a billing error — it is a commercial failure to capture value through expansion pricing. A customer who is consuming AI at a rate that significantly exceeds their contracted allocation but who has not been offered an expanded contract is receiving expansion value that the vendor is not capturing. The unrealized expansion gap is measured by comparing each customer's actual consumption to their contracted allocation and identifying the customers where actual consumption significantly exceeds the contracted amount.

Revenue Leakage Map — Seven Categories				
Leakage category	Root cause	Estimated prevalence	Detection method	Revenue impact formula
Metering blind spots	Event submission failures; no retry logic; network timeouts	High — affects most AI deployments without deliberate monitoring	Product activity log vs metering event log reconciliation	Lost events × applicable billing rate
Entitlement drift	Contract amendments not reflected in billing system	Very high — most common leakage category	Billing config vs contract terms audit per entitlement	(Actual entitlement – Billed entitlement) × consumption rate
Attribution failures	Missing or invalid commercial context on submitted events	Medium — peaks after org restructures and API key changes	Unattributed event queue volume and aging	Unattributed events × applicable billing rate

Outcome verification gaps	Broken webhooks integration between verification system and billing engine	Medium — affects all outcome-based products	Outcome verification events vs billing events reconciliation	Verified but unbilled outcomes × outcome price
Overage underbilling	Incorrect overage pricing configuration; wrong reset date	Medium — detected at billing cycle if BHI review is in place	Consumption vs entitlement vs invoice reconciliation	(Actual overage – Billed overage) × overage rate
Model upgrade undercharging	Billing config not updated after model upgrade	Low — but high per-instance impact	Deployed model vs billed model reconciliation	(New model rate – Old model rate) × consumption since upgrade
Unrealized expansion gap	Pricing model cannot capture consumption above contracted amount	Always present — magnitude varies by customer	Actual consumption vs contracted consumption by customer	(Actual consumption – Contracted allocation) × overage rate (opportunity cost)

Leakage Quantification Methodology

Quantifying revenue leakage requires a structured audit methodology that systematically examines each leakage category and estimates the revenue impact.

The metering leakage audit compares the product system's activity records to the metering system's event records. For each product, the product system should have a record of every significant activity that generates a billable event — every API call served, every agent workflow initiated, every outcome verified. The metering system should have a corresponding event for each. A systematic comparison of the two records identifies the events that were generated in the product system but not recorded in the metering system — the metering blind spots. Multiplying the count of missing events by the applicable billing rate produces the metering leakage estimate for the audit period.

The attribution audit examines the unattributed event queue. For each event in the queue, the audit determines whether the event can be attributed retroactively (by tracing the customer and product from the event's metadata), or whether the attribution information is irrecoverably missing. For events that can be attributed retroactively, the audit produces a correction event that attributes the original event to its correct billing context and adds the missed revenue to the current period. For events that cannot be attributed, the audit documents the reason and the revenue impact for the period close package.

The entitlement drift audit compares the entitlement configurations in the billing system to the entitlement terms in the signed contracts. For each entitlement, the audit verifies that the token budget matches the contracted amount, the reset period matches the contracted reset cadence, the SLA configuration matches the contracted SLA, and the price_id references the correct price for this customer and contract vintage. Discrepancies are flagged for correction and the revenue impact of each discrepancy is estimated.

The total leakage estimate from the audit — the sum of metering leakage, attribution leakage, entitlement drift leakage, and expansion gap — is typically expressed as a percentage of total billings for the period. Industry experience suggests that AI companies conducting their first formal leakage audit typically find leakage in the range of 3–8% of billings. At \$50M ARR, 5% leakage is \$2.5M of unrecognized revenue per year. At \$200M ARR, it is \$10M. The audit investment — typically 2–4 weeks of RevOps and engineering time — is justified on a straightforward ROI calculation for any AI company above \$20M in consumption-based revenue.

Revenue Leakage Audit — Methodology and Expected Findings				
Audit component	Methodology	Typical finding	Remediation	Time investment
Metering completeness	Compare product activity log to metering event count by customer/period	Event loss rate of 0.5–3% for organizations	Fix event submission retry logic; implement continuous	3–5 days engineering + 1 day RevOps

		without explicit event monitoring	metering reconciliation	
Attribution quality	Review unattributed event queue: volume, aging, categories	0.1–2% of events unattributed; peaks after org restructures	Attribution correction events for recoverable items; API key governance fix	2–3 days RevOps + 1 day engineering
Entitlement accuracy	Billing config vs contract terms: token budget, reset period, price_id, SLA	5–15% of entitlements have at least one parameter discrepancy	Systematic entitlement audit and correction; contract intelligence implementation	1 week RevOps
Outcome billing completeness	Outcome verification events vs billing events: are all verified outcomes billed?	1–5% of verified outcomes may not trigger billing events if integration is fragile	Webhook reliability improvements; outcome-to-billing event tracing	2–3 days engineering
Expansion gap quantification	Actual consumption vs contracted allocation by customer	10–30% of customers consuming significantly above contracted allocation	Expansion outreach with consumption data as evidence; proactive upsell programme	1 day RevOps analytics

FOR THE CONTROLLER**The leakage audit is a revenue assurance investment with a defined payback period**

A rigorous revenue leakage audit for an AI company with \$50M in consumption-based revenue typically costs \$150–250K in internal time (RevOps, engineering, finance) and produces leakage findings of \$1.5–4M annually (3–8% of revenue). Payback period: approximately 3–6 weeks of recovered revenue. The audit should be conducted annually for any AI company above \$20M in consumption billing, and quarterly for companies above \$100M where the leakage dollar amounts justify the continuous investment.

Chapter Six — The Essentials

- › Revenue leakage has seven categories — metering blind spots and entitlement drift are the most prevalent and highest-value to fix.

- › The leakage estimate for a first-time audit is typically 3–8% of consumption billing — significant at any scale above \$20M ARR.
- › The metering completeness audit compares product activity logs to metering event counts — the gap is the metering leakage.
- › The entitlement accuracy audit compares billing system configuration to contract terms — 5–15% discrepancy rate is typical before systematic governance is in place.
- › Annual leakage audits are justified on pure ROI for any AI company above \$20M in consumption billing.

CHAPTER SEVEN

Consumption Cohort Modelling and Agent Utilisation Forecasting

Usage cohort design. Token burn curves. Agent spend forecasting. The rolling 12-month view.

Consumption cohort modelling is the forecasting methodology that replaces SaaS ARR-based forecasting for AI businesses with significant consumption-based revenue. The fundamental difference: SaaS ARR-based forecasting assumes that a contract's revenue is determined at signing and changes only through renewal actions (churn, contraction, expansion). AI consumption forecasting recognizes that the revenue from any given contract depends on how much the customer uses the AI — a quantity that changes continuously as the deployment matures, as new use cases are deployed, and as the customer's business grows or contracts.

The cohort model groups customers by the month they were contracted and tracks their consumption trajectory over time. The key insight of cohort analysis for AI products is the consumption maturation curve: newly deployed AI customers typically have low initial consumption (they are still deploying and learning) that grows significantly over the first 6–12 months as the deployment matures. If the CFO forecasts consumption based on current rates without accounting for the maturation curve, they systematically

underforecast revenue for recently acquired customers and overforecast revenue for mature customers who may be approaching their consumption ceiling.

Building a reliable cohort model requires 12–18 months of consumption history. With that history, the CFO can model the consumption trajectory for the average customer at each stage of deployment maturity: what percentage of maximum consumption does the average customer reach in month 1, month 3, month 6, month 12, and month 24 after contract signing? This maturation profile, combined with the current cohort distribution of the customer base (how many customers are at each stage of deployment maturity), produces a consumption forecast that accounts for the natural growth in consumption as deployments mature.

The practical implementation uses the event data from the metering system to produce per-customer consumption time series, which are aggregated into cohort consumption curves. The cohort curves are then used to project future consumption for each customer based on their current deployment stage and the historical trajectory of similar customers. This approach produces forecasts that are more accurate than SaaS ARR-based extrapolation for AI businesses with significant consumption variability.

Consumption Maturation Curve — Typical AI Deployment Profile			
Deployment month	% of mature consumption rate	Revenue forecast implication	FP&A action
Month 1	15–25%	Significantly below run-rate — do not use as forecast basis	Apply maturation curve to project forward; do not annualize month 1 consumption
Month 3	40–55%	Growing but still below mature rate — cohort model critical for accuracy	Use cohort model; supplement with deployment pipeline data for expansion projection
Month 6	65–80%	Approaching mature rate — cohort model reliable for planning horizon	Cohort model + burn curve analytics provides reliable 6-month forward projection
Month 12	85–95%	Near mature rate — consumption stabilizing — minor growth from optimization	Rolling forecast with cohort model; identify expansion opportunities from usage patterns

Month 18+	90–100% (with periodic expansions)	Mature deployment — growth from expansion use cases, not maturation	Customer P&L analysis; identify expansion opportunities; flag customers below maturation benchmark
-----------	------------------------------------	---	--

Token Burn Curves

Token burn curves are the granular view of consumption that the CFO needs for both short-term budget management and longer-term financial planning. A token burn curve plots cumulative token consumption against time for a specific entity — a customer, a team, a workflow, or the entire organization — and extrapolates the trajectory to estimate when the current budget will be exhausted.

The burn curve's most important derived metric is the projected budget exhaustion date — the date at which cumulative consumption will equal the budget, assuming the current consumption rate continues. This date drives three financial management actions. If the projected exhaustion date is before the end of the budget period, the CFO must either approve a budget increase, implement consumption governance controls to reduce the burn rate, or accept that some AI activity will be throttled when the budget is exhausted. If the projected exhaustion date is well beyond the end of the budget period, the CFO may have over-allocated budget to this entity — budget that could be redeployed to higher-return activities.

The burn curve also reveals consumption pattern anomalies that are not visible in aggregate metrics. A customer whose token burn is normally smooth but spikes dramatically in week three of the month is probably running a batch processing workflow — a pattern that has implications for both billing (the spike may trigger consumption tier transitions) and financial planning (the batch workflow may represent an expansion opportunity). A team whose token burn is decelerating may be encountering model quality issues, reduced workload, or the beginning of a churn risk.

The burn curve analytics infrastructure requires that the metering system maintain per-entity consumption time series at daily granularity, with a retention period sufficient for the cohort analysis. The CFO's FinOps dashboard should surface the burn curves for the

top 20 customers by consumption volume and all entities approaching their budget limits, updated daily.

Token Burn Curve Analytics — KPIs and Actions				
Metric	Definition	Normal range	Action trigger	Management response
Daily burn rate	Average tokens consumed per day in the rolling 7-day window	Within $\pm 15\%$ of prior-week rate	Spike $> 50\%$ above prior week	Investigate workflow changes; check for misconfigured agents; verify billing config
Projected exhaustion date	Date when cumulative consumption reaches the budget ceiling at current rate	≥ 30 days before period end	< 14 days before period end	Budget increase approval workflow; or consumption optimization programme
Budget utilization rate	Cumulative consumption / Total budget at each day in the period	Roughly linear — no dramatic acceleration	Non-linear acceleration in the last 25% of period	Identifies batch processing patterns; may indicate overage pricing triggers approaching
Cost per outcome	Total token cost for a workflow / Number of outcomes produced by the workflow	Stable or declining as AI improves	Rising cost per outcome	Investigate model degradation, prompt drift, or task complexity increase

Agent Utilization Forecasting

Agent utilization forecasting is the most technically demanding AI forecasting problem because it requires predicting not just the consumption rate but the deployment breadth

— how many agent workflows will be running, across how many teams, on what volume of work.

The challenge is that agent deployments tend to expand non-linearly. A legal team that deploys an AI contract review agent for routine NDAs will, if the agent performs well, expand the deployment to more complex agreements, then to due diligence review, then to regulatory compliance monitoring. Each expansion increases token consumption, typically at a higher rate than the initial deployment because more complex workflows consume more tokens per document. Forecasting this expansion requires understanding not just current consumption but the pipeline of future deployments — which requires engagement with the business units managing AI deployments, not just analysis of historical data.

The agent utilization forecast is built in three components. The run-rate component projects current agent consumption forward using the cohort model: given the current deployment's maturation stage and the historical trajectory of similar deployments, what is the expected consumption in each future period? The pipeline component adds the expected consumption from agent deployments that are currently in progress or planned: new use cases being developed, new teams beginning AI adoption, new workflows being piloted. The upside component estimates the potential consumption from expansion opportunities that have been identified but not yet committed: customers whose current consumption patterns suggest they would benefit from additional AI deployment.

The total agent utilization forecast is the sum of these three components. The run-rate component provides high forecast confidence. The pipeline component provides medium confidence, because deployment timelines frequently slip. The upside component provides low confidence, but is important for the long-range financial planning that the board requires. The CFO should present the agent forecast with confidence intervals that reflect these three confidence levels — not as a single point estimate.

Agent Utilization Forecast — Three Components				
Component	What it covers	Confidence level	Data source	Update cadence
Run-rate	Current agent deployments projected forward using cohort maturation model	High — based on historical consumption data	Metering system + cohort model	Monthly (updated with latest consumption data)
Pipeline	Expected consumption from deployments in progress or planned	Medium — timeline slippage common	Business unit deployment plans + engineering roadmap	Monthly (reviewed with business unit leads)
Upside	Potential consumption from identified but uncommitted expansion opportunities	Low — many will not convert in the forecast horizon	CS expansion opportunity pipeline + customer health data	Quarterly

Chapter Seven — The Essentials

- › Consumption cohort modelling replaces ARR-based forecasting for AI businesses with significant consumption revenue.
- › The consumption maturation curve is essential: new customers consume at 15–25% of mature rate in month 1, reaching 85–95% by month 12.
- › Forecasting without the maturation curve systematically underforecasts revenue for new customers and overforecasts for mature ones.
- › Token burn curves are proactive governance tools — projected exhaustion dates drive budget management decisions, not surprises.
- › Agent utilization forecasts have three confidence tiers: run-rate (high), pipeline (medium), upside (low) — present to board with explicit confidence levels.

PART FIVE

AI-Powered FP&A

Agents in the finance function. The rolling forecast. What the CFO delegates — and what they do not.

CHAPTER EIGHT

AI-Powered FP&A: Agents in the Finance Function

The FP&A co-pilot. Rolling forecast architecture. What the CFO delegates to agents.

The finance function is among the highest-ROI applications of AI in an enterprise — and paradoxically among the most underdeployed. The reason for the paradox is cultural: finance teams are conservative about technology adoption (with good reason — errors in financial reporting have serious consequences), and the early AI tools for finance were general-purpose assistants that could not be trusted with the precision that financial analysis requires.

The next generation of AI for finance is different: specialized agents designed for specific financial tasks, with access to the company's actual financial data, governed by the same financial controls that govern human analysts, and auditable in their decision-making. These agents do not replace the judgment of the CFO or the Controller — they amplify it by doing the data assembly, the calculation, and the initial analysis that currently consumes most of the finance team's time.

The FP&A co-pilot is the primary AI application for finance team productivity. It is an AI agent with access to the financial data warehouse, the billing system, the budgeting system, and the prior periods' financial statements. Its core capability is translating a financial question — "what drove the variance in token costs this quarter versus plan?" — into a structured analysis by retrieving the relevant data, performing the calculation, and generating a narrative explanation with supporting evidence.

The FP&A co-pilot's most valuable application in an AI business is the weekly consumption variance report: a structured analysis of the difference between planned and actual token consumption by business unit, workflow, and model, with a narrative explanation of the drivers and a projection of the impact on full-period financial results. Preparing this analysis manually takes a finance analyst 4–6 hours. The FP&A co-pilot prepares it in minutes, freeing the analyst to focus on the interpretation and the management conversation.

The Rolling Forecast Architecture

The rolling forecast is the quarterly FP&A deliverable that has become the primary financial planning tool in AI businesses, replacing the annual budget as the primary planning document for most AI-specific spending.

The replacement of the annual budget with a rolling forecast reflects a specific characteristic of AI businesses: the spending curve is too uncertain and too dynamic for a single annual budget to be meaningful. A company that approves a \$5M annual AI infrastructure budget in January may find by March that a new model capability has changed their deployment plans dramatically — either expanding the opportunity (requiring more budget) or reducing the cost (freeing up budget for other investments). The annual budget process, which is typically finalized 2–3 months before the fiscal year begins, is too infrequent to capture this dynamism.

The rolling forecast solves this by maintaining a 12-month forward financial projection that is updated monthly. Each monthly update incorporates: actual results for the most recently completed month, updated consumption model projections for the remaining periods based on current consumption data, pipeline updates from the commercial team, and any changes to model pricing or infrastructure costs. The rolling forecast gives the CFO and the board a current-information view of AI financial trajectory that an annual budget cannot provide.

The AI-specific inputs to the rolling forecast that differ from SaaS forecasting: the consumption model projection (token burn curves and cohort trajectories for each customer segment and deployment), the model cost rate table (current prices per million tokens for each model in use, updated when prices change), the agent deployment pipeline (planned new agent deployments and their expected consumption impact), and the variable consideration estimate for outcome-based revenue (updated based on current AI performance data).

What the CFO delegates to agents in the rolling forecast process: the data assembly (pulling current period actuals from each source system, updating consumption projections from the metering system, refreshing cost rate inputs), the calculation layer (applying the cohort model to project consumption, calculating variable consideration estimates, running sensitivity analyses), and the first-draft narrative (generating the management commentary that explains the key variances and risks). What the CFO retains: the judgment on the variable consideration constraint (the CFO must assess and approve the estimate, not just accept the agent's calculation), the communication of financial results to the board and investors, and the interpretation of what the financial results mean for the business strategy.

Rolling Forecast — Agent vs Human Responsibilities			
Forecast component	Agent responsibility	Human (CFO/FP&A) responsibility	Why the distinction matters
Data assembly	Pull current period actuals from each source system; update consumption projections from metering; refresh cost rate inputs	Review data quality flags; resolve missing data sources	Agents handle volume; humans handle judgment on data quality questions
Calculation layer	Apply cohort model to project consumption; calculate variable consideration estimates using configured methodology; run sensitivity analyses	Review and approve the variable consideration constraint estimate; approve key assumptions	CFO must personally assess and approve the variable consideration constraint — not delegate to agents

Sensitivity analysis	Model financial outcomes under base, bear, and bull consumption scenarios; calculate impact on key metrics (NRR, gross margin, DSO)	Determine the scenarios to model; interpret the results in strategic context	Agent precision + human strategic judgment = better decisions
Variance analysis	Calculate actual vs planned variance by category; generate narrative explanation of top 5 variance drivers	Interpret variances in context; determine management actions	The narrative explanation draft is agent work; the management response is human work
Board commentary	First draft of management commentary: key variances, risks, opportunities	Edit, approve, and own the board communication	CFO owns the board communication; agent does the first-draft data work
Variable consideration update	Calculate updated variable consideration estimates based on current performance data	Assess constraint: is it highly probable that the updated estimate will not result in a significant reversal?	The constraint assessment is a judgment call that must be made by an accountable human – not automated

What the CFO Delegates to Agents — and What They Do Not

Task	Delegate to agents?	Rationale
Pulling period actuals from source systems	Yes	Volume work with no judgment required — agents do it faster and more completely
Consumption forecast calculation using established cohort model	Yes	Mechanical application of documented methodology — agents do it accurately
Variance calculation (actual vs plan)	Yes	Arithmetic — agents do it without error
First-draft variance narrative	Yes	Structured writing from structured data — agents produce a useful first draft
Variable consideration constraint assessment	No	Requires judgment about probability — 'highly probable' is a legal standard that a human must assess and sign off
Revenue recognition policy determination	No	Accounting standards require professional judgment; policy must be approved by an accountable person

Board communication	No	The CFO owns this communication; agents assist, they do not own
Audit response	No	Auditors are engaging with the organization, not with its agents — humans must respond
Material judgment calls (impairment, contingencies, concern) going	No	Materiality judgments require experience, context, and accountability that agents cannot provide

FOR THE CFO**The variable consideration constraint is the one FP&A judgment that must never be delegated to an AI agent**

Under ASC 606, the variable consideration constraint requires the CFO (or their accountable designee) to assess whether it is 'highly probable' that including a variable consideration estimate will not result in a significant revenue reversal. This is a legal standard with restatement implications if incorrectly applied. An AI agent can calculate the estimate and model the scenarios, but the constraint assessment — the judgment about whether the conditions for high probability are met — must be made by a person who is accountable for the financial statements. This judgment cannot be automated, and organizations that allow agents to make it are assuming accounting risk they may not fully appreciate.

Chapter Eight — The Essentials

- › The FP&A co-pilot handles data assembly, calculation, and first-draft narrative — freeing the finance team for interpretation and judgment.
- › The rolling 12-month forecast updated monthly is the right planning tool for AI businesses — annual budgets cannot accommodate AI's pace of change.
- › The variable consideration constraint assessment must never be delegated to an AI agent — it is an accountable human judgment under ASC 606.
- › The table of 'what the CFO delegates and does not' is the governance boundary — document it and enforce it consistently.
- › AI in the finance function compounds over time: each quarter of clean data makes the FP&A agent more accurate and more useful.

CLOSING

Build the Visibility Infrastructure — Before You Need It

The CFO who can see AI economics clearly can govern them. The CFO who cannot is navigating without instruments.

The CFO of an AI business is governing something genuinely new: a commercial operation where costs are variable at machine speed, revenue has a significant component that depends on AI performance, and the spending governance mechanisms that exist for human-paced procurement do not apply to autonomous agent workflows.

The infrastructure described in this book — the token budget hierarchy, the consumption dashboards, the cost attribution methodology, the variable consideration policies, the cohort forecasting model, the leakage audit framework, the FP&A agent overlay — is not the infrastructure of a mature AI-native finance function. It is the infrastructure of a finance function that is beginning the transition to AI-native operations, built deliberately, with the data precision and governance rigor that the CFO role demands.

The maturity curve for AI-native finance operations runs three to five years for most organizations. It begins with visibility: the token dashboard and cost attribution infrastructure that makes the financial picture of AI spending legible. It advances through governance: the budget hierarchy, the approval workflows, and the entitlement enforcement that translates financial visibility into financial control. It matures through forecasting: the cohort models, the burn curve analytics, and the rolling forecast infrastructure that makes AI financial planning as rigorous as any other business planning function. And it reaches full maturity when the AI agents deployed in the finance function are themselves governing the AI agents deployed in the business — a recursive application of AI governance that is the organizational expression of what it means to be AI-native.

The CFO who builds this infrastructure is not just making the finance function more efficient. They are making the entire AI commercial operation more governable — which makes it more investable, more scalable, and more likely to generate the returns that justified the original AI investment.

Build the visibility infrastructure. Build it before you need it. And build it with the conviction that the CFO who can see AI economics clearly is the CFO who can make the AI investment case to the board with the confidence that role requires.

"The CFO who builds this infrastructure is not just making the finance function more efficient. They are making the entire AI commercial operation more governable — which makes it more investable, more scalable, and more likely to generate the returns that justified the original AI investment."

The AI Economy Monetization Series continues in Book Three-B:

Revenue Integrity and A/R Governance