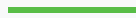


The Definitive Guide to Hybrid FinOps

Translating the Datacenter for the Cloud Era



A HYBRIDFINOPS PUBLICATION

Sponsored by Visual One Intelligence®

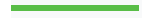
CONTENTS

- 1. The Cloud Cost Paradox and Its Datacenter Cousin
- 2. The Six Principles, Revisited for a Hybrid World
- 3. The CapEx Trap
- 4. The Anatomy of a Datacenter Bill
- 5. From Hardware to OpEx Metrics
- 6. Unit Economics for Physical Infrastructure
- 7. The Tag Mapping Problem
- 8. Inform: Visibility Across Estates
- 9. Optimize: Workload Placement at Parity
- 10. Operate: Showback, Chargeback, and the Culture Shift
- 11. Getting FinOps and IT to Actually Cooperate
- 12. A Maturity Model for Hybrid FinOps

- Appendix A: Tool Categories & Vendors
- Appendix B: Glossary of Key Terms
- About This Guide

CHAPTER

One



The Cloud Cost Paradox and Its Datacenter Cousin



“First they asked us to fix cloud. Then fix the software mess. Now fix the contract and license mess, now fix the data center...”

— FinOps practitioner, 2026 State of FinOps survey

In early 2026, the FinOps Foundation did something it had not done in the six years of the discipline's public existence. It changed its mission.

The old mission — the one engraved in a thousand conference decks — read *“Advancing the people who manage the value of Cloud.”* The new one reads *“Advancing the people who manage the value of Technology.”* A single word changed. But that word formalized a shift the practitioner community had already been making for years, quietly and on its own.

The numbers are in the Foundation's most recent State of FinOps survey, which sampled 1,192 practitioners overseeing more than \$83 billion in annual cloud spend. Forty-eight percent now manage datacenter spend, up twelve points year-over-year. Fifty-seven percent manage private cloud, up eighteen. Sixty-four percent manage licensing. Twenty-eight percent have begun including labor costs. Ninety-eight percent manage AI spend — a category that did not meaningfully exist two survey cycles ago.

The scope has officially expanded. And the Foundation, the neutral authority on what the practice is and is not, has ratified the expansion after the fact. FinOps is a technology discipline now. The cloud was the starting point, not the destination.

The Original Paradox

Before we get to why this matters for the datacenter, recall what the cloud cost paradox actually was.

Cloud computing was sold as a cost-savings model — pay for what you use, scale down when you don't, trade capital expense for operating expense. For a brief moment around 2012, this was even mostly true. Then something strange happened. Companies moved workloads to the cloud and their costs went *up*. Not in absolute dollars — that was expected, because they were also growing — but as a percentage of IT budget, as a percentage of revenue, as a line item on quarterly earnings calls that CFOs increasingly wanted to discuss.

The paradox was not that cloud was expensive. The paradox was that cloud *appeared* to be the cheaper option while somehow delivering a more expensive outcome. The consumption model that was supposed to discipline spending had instead inverted it. Instead of a handful of procurement officers signing contracts every few years, a thousand engineers were provisioning

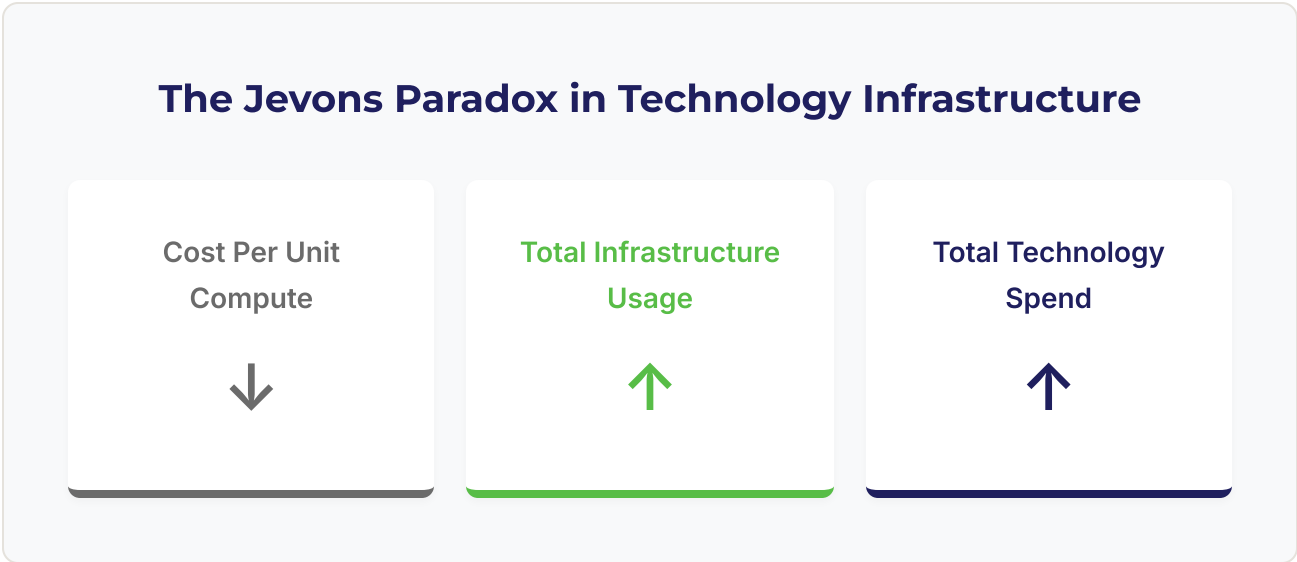
resources every few minutes. Nobody had visibility. Nobody had accountability. The cloud bill arrived monthly, written in a foreign dialect of three thousand SKUs, and the only rational response was to hire a team to translate it.

That team became FinOps.

The Cousin in the Datacenter

Here is the part that has been hiding in plain sight for a decade. The datacenter has the same paradox. It just hides differently.

On-premises infrastructure is frequently described — even in serious conversations, by serious people — as *“already paid for.”* The server was purchased three years ago. The storage array has been depreciated down to a sliver. The hypervisor license is on a multi-year enterprise agreement. The phrase *“we already paid for it”* is the most expensive sentence in IT, and it is uttered with such confidence that nobody asks what exactly has been paid for.



What has been paid for is the purchase price. That is all. The operating cost of that same infrastructure — the power it draws, the floor space it occupies, the personnel who maintain it, the maintenance contract that keeps it running — typically exceeds the amortized purchase cost by a factor of two or three. On a five-year-old storage array that originally cost a hundred thousand dollars, the depreciation contribution to daily cost lands somewhere around fifty-five

dollars. The power and floor space and personnel contribution lands closer to one hundred twenty-five. The purchase price, in other words, represents roughly a third of the true daily cost. The “we already paid for it” component is the smallest line item.

No engineer provisioning a new VM on that array knows this. No workload owner placing a new application knows this. No procurement officer comparing an on-prem refresh to a cloud migration knows this, because the cloud proposal comes with fully-loaded operating costs priced in and the on-prem alternative does not. The datacenter looks cheap because two-thirds of its cost is invisible.

This is the same paradox the cloud had in 2014, presented as its mirror image. The cloud made costs visible and then spiraled because visibility without governance is chaos. The datacenter hid its costs and spiraled in the opposite direction — not through provisioning sprawl, but through stranded capacity, missed refresh cycles, and workload-placement decisions made on mental models that were wrong by a factor of three.

FinOps was invented to solve the first paradox. The same principles, applied with a different set of inputs, solve the second one.

Why Now

The Foundation's mission change is not a speculative gesture. It is a lagging indicator of work that practitioners had already been doing.

The same 2026 survey that documented the scope expansion also captured what practitioners are asking their tools to do next. The top three most-requested capabilities that do not currently exist are: granular AI spend monitoring, pre-deployment architecture costing, and — in the third position — *“a single pane of glass for different technology spend.”* The market has named the gap. It has also named the two places where the gap is most painful: datacenter and AI are the top two expansion requests for FOCUS, the emerging specification for normalized cost and usage data.

Another practitioner, quoted in the same survey, put the arc plainly: *“dashboards are table stakes of yesterday.”* The leading edge of the discipline is no longer the cloud bill. It is the unification of every form of technology spend into a single coherent conversation.

This book is about what it takes to include the datacenter in that conversation. The principles are already written. The practice is already mature. The only missing piece is a way to translate on-premises infrastructure into the same financial language FinOps already speaks for cloud — an OpEx-shaped view of a CapEx-shaped asset, an invoice where no invoice has ever existed.

That translation is the subject of the next eleven chapters.

CHAPTER

Two

The Six Principles, Revisited for a Hybrid World

When the FinOps Foundation codified its six guiding principles, cloud was the only infrastructure model under the discipline. You might expect, given that, that the principles would be riddled with references to cloud-specific mechanics — billing APIs, elasticity, provider accounts, per-second metering. They are not. Read them now, and exactly one of the six contains the word "cloud." The other five speak of teams, technology, ownership, data, and central enablement — concepts indifferent to provisioning model.

That was not an accident. The principles were written by practitioners who had already lived through one industry rebrand and were determined not to have to rewrite them through the next one. The Foundation's 2026 mission change — from managing the value of cloud to managing the value of technology — caught the operating model up to what the principles had always allowed.

So the principles hold. What changes, when you extend them to a datacenter estate, is the plumbing underneath.

1. Teams need to collaborate

In cloud, the collaboration imperative is driven by speed. Resources are provisioned per-second; charges accrue per-second; engineers can spend a month's budget in an afternoon if nobody is watching. Finance, engineering, and product have to be in the same conversation in near-real-time, because the system operates faster than any quarterly review cadence.

In the datacenter, the imperative looks different because the pace is different — but the silos are worse. Finance owns depreciation schedules. IT owns utilization. Procurement owns the refresh cycle. Facilities owns the power bill. A hybrid practice doesn't just ask these teams to share data; it asks them to agree that the data they've each been looking at has, all along, been describing the same dollar.

2. Business value drives technology decisions

The principle already avoids the word "cloud." Good thing. In cloud, the work here is translating aggregate spend into unit economics — cost per transaction, cost per customer, cost per feature — so that decisions are made on value rather than on whichever line item scared someone in a status review.

On-prem, the translation is harder because the raw material is absent. A rack does not emit a per-transaction cost; it emits heat. The numerator has to be constructed from capital and operating expense; the denominator has to be inferred from capacity reporting. Until both exist, "business value drives technology decisions" in the datacenter reduces to "whoever shouts loudest wins." Many organizations have been running this way for decades and didn't notice, because there was nothing to compare it to.

3. Everyone takes ownership for their technology usage

In cloud, ownership is pushed to the edge. Engineers see the cost of what they provision because the billing data makes it impossible to hide, and product teams feel the budget when they blow through it. This is the principle that most often distinguishes a FinOps-native organization from one that is merely cost-aware.

The datacenter has historically been the exact inverse. Nobody owns the cost of an application running on shared infrastructure, because the cost was paid two years ago by someone in a different building. "We already paid for it" is the phrase that kills ownership. An OpEx translation fixes this by re-imposing the monthly running cost — not because cash changed hands this month, but because the dollar of useful life consumed this month is the dollar that ownership can actually attach to.

4. FinOps data should be accessible, timely, and accurate

In cloud, this is mostly a collection and presentation problem. The data exists; the provider emits it; the work is getting it cleaned, tagged, and onto a dashboard the right audience will actually open. Tooling maturity here is high, and practitioners treat real-time visibility as the floor rather than the ceiling.

In the datacenter, the principle runs into a bracing fact: there is no billing API. Cost data has to be synthesized from asset inventories, capacity reports, power telemetry, and facility invoices that arrive on a 30-day lag. Accuracy becomes the hardest of the three adjectives — partly because the inputs are scattered, partly because half of them were never collected for financial purposes in the first place. The capability priorities practitioners identified for datacenter FinOps in the 2026 State of FinOps survey — Allocation and Planning & Estimating at the top of the list — are the discipline's polite acknowledgment of exactly where the work is.

5. FinOps should be enabled centrally

In cloud, the central team sets standards, owns rate negotiations and commitment purchases, and evangelizes the practice without taking operational responsibility away from the teams running the workloads. The model is explicitly patterned on the security function: a small center, federated execution.

In a hybrid estate, the center has to speak two languages at once — the cloud vocabulary of billing data, tags, and commitments, and the datacenter vocabulary of capacity, refresh cycles, and depreciation schedules. The 2026 survey finding that 78% of FinOps teams now report to the CTO/CIO, up 18 points since 2023, is quietly doing a lot of work here. When FinOps sits under

technology leadership rather than finance, the central team has the organizational standing to make the datacenter estate legible on the same terms as the cloud estate. Reporting to the CFO, by contrast, tends to reinforce the old ledger split.

6. Take advantage of the variable cost model of the cloud

This is the only principle with "cloud" in its headline, and it is the one that does not extend cleanly. The datacenter is not variable. A storage array does not cost less on a quiet Tuesday. The useful extension is not to pretend otherwise, but to notice what the principle is actually about underneath: matching commitment to need, and avoiding the cost of capacity you don't use. In cloud, that means right-sizing and commitment purchasing. In the datacenter, it means refresh-cycle discipline, honest capacity planning, and — critically — the ability to see when a workload would be cheaper somewhere else. The variable cost model is a feature of cloud. The underlying discipline of matching cost to use is not.

Five principles that already fit. One that needs a translation. This is, roughly, the work of the rest of this book: the principles hold, the phases hold, the capabilities hold. What has to be built is the data layer underneath them — the one the datacenter never shipped with. That begins, uncomfortably, with an honest accounting of why a model built on "we already paid for it" has been quietly distorting hybrid decisions for years.

CHAPTER

Three

The CapEx Trap

“We already paid for it.”

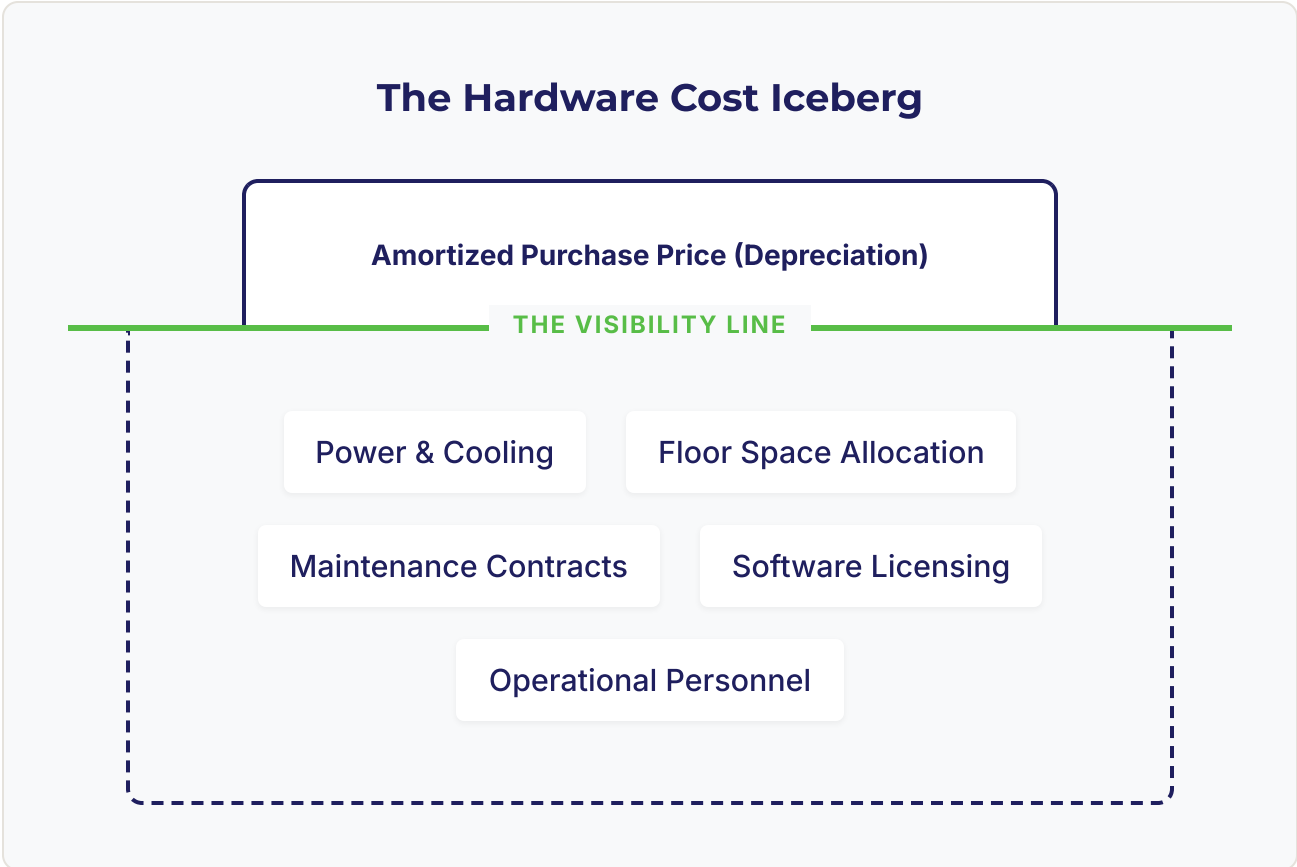
— Every workload-placement meeting, approximately

The sentence Chapter 1 identified as the most expensive in IT has a formal name. It is called depreciation. It is a legitimate accounting concept, developed over a century of industrial practice to describe how the value of a fixed asset declines over time. It works beautifully for what it was designed to do.

What it was designed to do is not what it is being used for.

Where the Trap Comes From

Depreciation treats a server the same way it treats a forklift — a five-year useful life, a straight-line curve, a book value that ticks down each quarter until it reaches zero. For accounting and tax purposes, this is fine. For workload-placement decisions, it is a disaster.



The problem is not that depreciation is wrong. The problem is that depreciation captures only one component of cost: the purchase price, spread over time. It ignores everything else. Power, floor space, personnel, maintenance — the elements that determine whether a server is cheap or expensive to *operate* — are not on the depreciation schedule, because depreciation was never intended to measure them.

This mattered less when IT was capital-intensive and operating costs were a fraction of the outlay. For a mainframe in 1985, the purchase price dominated and depreciation was a reasonable proxy for total cost. It is no longer 1985.

The Proportion Problem

A representative mid-range compute host today — a 32-core server with 256 GiB of memory, purchased for fifty thousand dollars on a five-year refresh — contributes roughly twenty-seven dollars per day to depreciation. Its operating cost across power, floor space, personnel, and maintenance runs closer to one hundred thirteen dollars per day.

The purchase price represents less than twenty percent of true daily cost. The other eighty percent — the part that is not on the depreciation schedule, not on the asset ledger, not in the conversation — determines whether the machine is a bargain or a quiet liability.

When a workload owner is told the server is "already paid for," they are being shown the smallest line item and asked to make a decision. That decision is being made against a cloud proposal where the quoted price includes every invisible component, because the cloud provider has no interest in hiding them. The comparison is structurally unfair, and the on-prem side consistently wins on the basis of a number that does not exist.

Different Ledgers, Same Asset

Even the organizations that understand this cannot easily fix it, because the true cost of on-premises infrastructure lives in four or five separate ledgers owned by different people.

Finance owns the depreciation schedule. Facilities owns the power bill and the square footage allocation. IT operations owns the maintenance contracts and the headcount. Nobody owns the total. And because nobody owns the total, nobody is accountable for the total being wrong.

This is precisely the silo FinOps was invented to dissolve. The discipline's earliest contribution to cloud was not any particular optimization — it was the simple structural act of pulling a cloud bill out of IT's inbox and putting it in front of every team that generated a line on it. The cost became visible to the people making the decisions. Hybrid FinOps is the same move, applied to an asset whose true cost has never been visible to anyone.

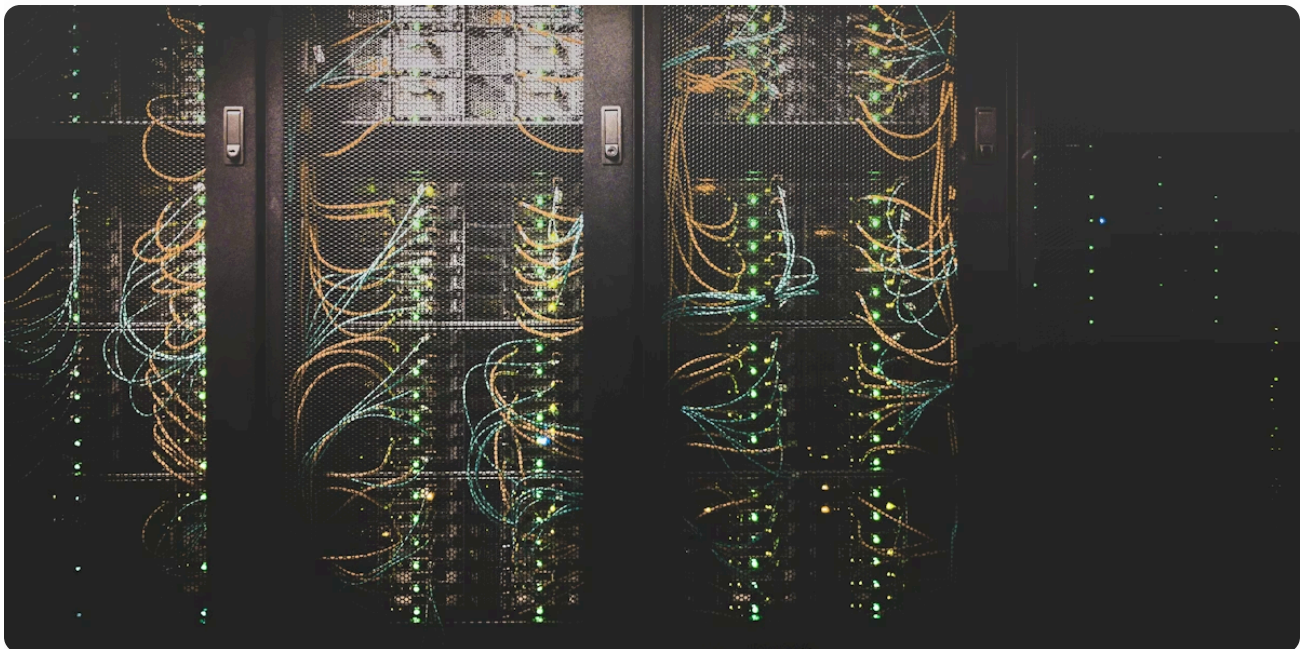
The CapEx trap is not a failure of intention. It is a failure of accounting infrastructure. Nobody is hiding the numbers. The numbers are scattered across systems that were never designed to talk to each other, being read by people who have no reason to add them together.

Until someone does.

CHAPTER

Four

The Anatomy of a Datacenter Bill



The cloud bill is a thing that exists. It arrives every month in a file of staggering complexity, written in a dialect most organizations need translation software to read. It has its own problems. It also has one property the datacenter infrastructure lacks.

It exists.

The datacenter bill, by contrast, is a thing that has to be constructed. The components are real. The numbers are real. They are simply not assembled anywhere. Nobody prints a monthly invoice showing the total cost of running a specific storage array, because no system is designed to produce one. The first structural task of hybrid FinOps is to build this bill — line item by line item, from ledgers that were never meant to be read together.

The Five Line Items

Every piece of hardware in a datacenter has five cost components. Each one is individually boring. The sum is not.

Amortized purchase price is the first line, and the most familiar. It is the capital cost divided by useful life, expressed as a daily rate. For a fifty-thousand-dollar server on a five-year refresh, that is twenty-seven dollars a day. This is the line Finance already tracks, and the line most workload-placement conversations wrongly treat as the whole bill. It is typically less than a third of the total.

Power is the second line and the most underestimated. Enterprise servers draw three hundred to a thousand watts continuously, around the clock, year after year. At commercial electricity rates — which have risen roughly forty percent in most North American markets since 2020 — a single mid-range server burns through seven to fifteen hundred dollars a year in electricity alone. High-density compute nodes and large storage arrays can exceed three thousand. Power is billed monthly by the utility, typically at the facility level, which means it is measured but almost never allocated back to specific hardware.

Floor space is the third line. Datacenter floor space is expensive — whether the facility is owned, leased, or collocated — because the cost includes cooling, redundancy, physical security, and the depreciation of the building itself. Industry allocations run between four hundred and fifteen hundred dollars per month for a full rack, prorated down to the device by U height or footprint. Real cost. Rarely attributed.

Personnel is the fourth line and the one that most surprises organizations when they first add it up. Running a datacenter takes people — administrators, engineers, operators, capacity planners, security staff, backup operators. These salaries represent somewhere between two hundred and two thousand dollars per month per device, depending on complexity and operational model. The number is usually buried in IT operations headcount and treated as fixed overhead, which means it is present in the overall IT budget but absent from any particular asset's cost profile.

Maintenance and support contracts form the fifth line. Vendor-provided hardware support runs between two and twenty percent of the original purchase price, per year. Basic support is affordable; premium support with four-hour on-site response is not. These contracts are typically negotiated centrally, which is efficient for procurement and useless for allocating cost to individual assets.

The Multipliers

The five line items produce a total. Three multipliers turn that total into something useful.

Useful life determines the amortization horizon. Industry convention for enterprise hardware is five years. Some organizations use three, some seven. The choice is not purely an accounting preference — a shorter useful life raises the daily cost but aligns with refresh cycles; a longer useful life suppresses the daily cost but hides stranded capacity and obsolescence risk.

Refresh cycle is related but distinct. Useful life is the accounting horizon. Refresh cycle is the operational horizon — when the organization actually plans to replace the asset. The two often diverge, and badly. A server that was fully depreciated after five years but is still running at eight years old contributes zero to the amortization line — but it is still consuming power, floor space, personnel, and maintenance. The purchase component has gone to zero. The other four have not. The cost per unit of delivered service rises as the asset ages, and nobody is necessarily watching it happen.

Utilization is the third multiplier, and the one that matters most for per-workload economics. A server running at twenty percent utilization has the same total daily cost as one running at eighty percent — and delivers a quarter as much compute. The per-workload math is determined not by the total bill but by the total bill divided by effective utilization. Low utilization is the datacenter equivalent of cloud overprovisioning, and it is vastly more common than anyone admits.

The Structural Parallel

The FinOps Foundation's curriculum teaches that the cloud bill, at its essence, is *usage × rate*. The datacenter bill is the same relationship, expressed in the arithmetic of fixed assets: *(purchase + operating) ÷ useful life × utilization*. The math is different. The information being captured is identical.

Every line item on the cloud bill has a datacenter analog. Compute time corresponds to amortized server cost. Storage volume corresponds to amortized array cost. Support tier corresponds to the maintenance contract. Regional pricing corresponds to facility location cost. The mapping is almost complete.

What is missing is the invoice that puts the analogs together. The components live in five different systems, owned by four different departments, sampled at four different cadences. The assembly is a job nobody has been asked to do, because until recently nobody needed the answer.

That job is the subject of the next chapter.

CHAPTER

Five

From Hardware to OpEx Metrics

The cloud side of a FinOps practice gets its input data from a billing API. The datacenter side has to manufacture that data. This chapter is about what the manufactured version looks like, how it is assembled, and why — once you have it — the rest of the book becomes possible.

The translation happens through a single formula. It is not elegant, but it is honest:

**(Amortized Purchase Price + Power + Floor Space +
Personnel + Maintenance)**

÷ Useful Life = Daily Cost

Five inputs. Five things the cloud bill silently totals for you every month, and that the datacenter does not.

Amortized Purchase Price is the hardware capital expense spread across its useful life — typically five years for compute and storage, seven for networking. A \$100,000 array over 1,825 days comes out to \$54.79 a day, every day, whether anyone uses it or not. This is the input most organizations already have in a depreciation schedule somewhere, though usually not in a form anyone outside finance can find.

Power is the metered electrical draw of the device plus the cooling overhead that comes with it — typically expressed as a PUE multiplier on the raw draw. Facilities knows this number. The rest of the organization usually does not.

Floor Space is the fully-burdened cost of the rack units a device occupies: building depreciation, cooling infrastructure, UPS, generators, and real estate, allocated per-U.

Personnel is the operational labor attached to the device — monitoring, patching, administration — expressed per-device per-day. This is the input most organizations skip, which is exactly why a meaningful fraction of their true cost of ownership goes invisible.

Maintenance is the vendor support contract. Usually annual. Usually tracked. Usually in a contract database nobody has mapped to the asset inventory.

Useful Life is the denominator, expressed in days because daily is the unit every downstream calculation wants to speak in.

Normalization: from daily cost to comparable rates

Daily cost per asset is useful only to the person who owns that asset. For a FinOps practice, the unit that matters is cost per thing-the-business-cares-about — capacity, compute, workload. That means normalizing.

Storage normalizes to **\$/GiB/day** — or **\$/TiB/year** for people who prefer rounder numbers on their slides. Compute normalizes to **\$/core/day** for physical hosts and **\$/VM/month** for the virtual layer running on top. These are the units that line up cleanly against cloud SKUs. An EC2 on-demand price is **\$/vCPU/hour** expressed differently. An S3 price is **\$/GiB/month** with a different decimal point.

Worked examples

Three representative cases, computed from the formula above. Numbers are anonymized but typical of mid-market enterprise infrastructure.

A 100 TiB mid-tier storage array. Purchase \$100,000, 5-year life. Annual operating cost \$46,255 — power \$7,300, floor space \$14,600, personnel \$24,355, plus vendor maintenance. Depreciation runs \$54.79 a day; operating runs \$126.73. Total daily cost: **\$181.52**. Normalized: **\$1.82 per TiB per day**, or about \$664 per TiB per year. That number is directly comparable to any cloud object store's annual rate. For many array classes, the on-prem number comes in below the cloud number. For others, it does not. You now have the data to tell the difference.

A 32-core compute host, 256 GiB memory. Purchase \$50,000, 5-year life. Annual operating cost \$41,345. Total daily cost: **\$140.68**. At a standard 50/50 CPU/memory allocation, that works out to roughly **\$2.20 per core per day** and about \$0.27 per GiB of memory per day.

A mid-range VM on that host. At typical consolidation ratios — forty-five to fifty VMs per host at a moderate overcommit — the allocated cost lands at **about \$90 per VM per month**. That is the unit the application team actually cares about, and the unit directly comparable to the monthly run rate of an equivalent instance at any cloud provider.

The Cold Start Problem

The formula is trivial. The data is not.

Ask any practitioner who has attempted this exercise at scale what took the time, and they will not say the math. They will say sourcing the inputs: finding the purchase date for the array in rack 14, the list price of a server bought by a team that no longer exists, the typical power draw for a model the vendor stopped selling four years ago, the floor-space cost for a colo contract owned by a facilities team three orgs removed from IT.

This is the cold start problem, and it is why the 2026 State of FinOps survey identifies Allocation and Planning & Estimating as the #1 and #2 capability priorities for datacenter FinOps. Practitioners are not confused about what to do. They are blocked on the data they need to do it with.

What has changed in the last two years is the emergence of AI-assisted asset estimation. Platforms such as Visual One Intelligence® now retrieve purchase dates, list prices, typical power draw, floor-space footprint, and personnel-cost estimates from device identity alone — on sight, without a multi-month spreadsheet exercise. The Crawl row of the maturity model — the six-to-twelve months most organizations have historically spent chasing inputs before the first usable number appears — collapses. Teams that would once have needed a year to reach Walk start there.

The difficulty of the cold start problem is not primarily technical. It is organizational. Sourcing the five input values for a single asset manually requires a FinOps practitioner to coordinate with procurement for original purchase records, ITAM for asset inventories and lifecycle data, facilities for power and floor space allocations, the vendor-management or MSP relationship for maintenance contracts, and finance for depreciation schedules. Each of those groups has its own priorities, reporting lines, and response cadence. None of them report to FinOps. The prerequisite to even beginning the OpEx translation is a five-way cross-functional collaboration that most FinOps teams do not have the authority to compel — which is why, for many organizations, the extension of FinOps into private cloud has stalled before the math was ever attempted.

Platforms that provide AI-generated default estimations for asset cost parameters collapse this entire dependency chain. By populating new asset records with industry-standard list prices, typical purchase-price ranges, and representative operating cost

values at the moment of discovery, these platforms allow practitioners to begin the OpEx translation on day one with defensible starting figures. The five cross-functional conversations still need to happen eventually — but they happen later, to refine estimates rather than to enable them. This is the difference between a FinOps practice that reaches the Walk stage in the first month and one that reaches it in the second year.

The formula has not changed. What has changed is that manufacturing its inputs no longer takes a year.

What these numbers enable

Three things, in order of immediate usefulness.

First, they make an asset comparable to itself over time. A \$1.82/TiB/day array this year against a \$1.45/TiB/day successor next year is the kind of comparison that retires old infrastructure on schedule.

Second, they make an asset comparable to another asset in the same estate. A \$90/month VM on a five-year-old host against a \$62/month VM on a newer one tells you where the next workload should land, and whether the refresh cycle is actually paying for itself.

Third — and this is the conversation the rest of the book is about — they make an on-prem asset comparable to a cloud asset. That comparison is not a one-line subtraction. It is built on unit economics, and it requires that the two sides speak a shared taxonomy. Those are Chapters 6 and 7.

CHAPTER

Six

Unit Economics for Physical Infrastructure

Chapter 5 produced a pile of daily cost rates — \$181 for the array, \$140 for the host, \$90 for the VM. Useful, but they describe infrastructure. The business does not run infrastructure. It runs workloads on top of it. The question that actually matters — whether an application should live here or there, cost this or that, scale up or out — is answered at the workload level.

This is unit economics. Cloud FinOps has been doing it for a decade. Doing it on-prem has, until recently, been arithmetically fine and operationally impossible — the inputs were missing. With Chapter 5's inputs in place, it becomes arithmetic again.

The ladder

Unit economics is a ladder, and each rung abstracts further from the hardware.

Cost per VM per month is the first rung and the most tractable. A 32-core host with 256 GiB of memory at \$140.68 a day, running forty-five to fifty VMs at a typical overcommit, allocates out to roughly \$90 a month per VM. For a specifically sized instance — say 2 vCPU, 8 GiB — proportional allocation across CPU and memory lands in the \$80–\$110 range depending on consolidation ratio.

Cost per core-hour is the cloud's native currency, translated. That same 32-core host resolves to about \$0.18 per physical core-hour, all-in. With a 2:1 vCPU overcommit — typical for general-purpose workloads — it's closer to \$0.09 per vCPU-hour.

Cost per transaction and **cost per customer** are the rungs the business actually wants. They require coupling infrastructure cost to an application-level metric: requests served, orders processed, users active. A payments service on a three-host cluster at \$420/day, processing a million transactions a day, runs \$0.00042 per transaction. Whether that number is good or bad depends entirely on the cost of the same transaction somewhere else.

That "somewhere else" is the cloud.

The comparison that FinOps has been waiting for

The following is an honest comparison, not a rigged one. A general-purpose, steady-state web workload on two comparable platforms:

	AWS m5.large	Equivalent vSphere VM
Specs	2 vCPU, 8 GiB	2 vCPU, 8 GiB
On-demand / all-in monthly	~\$70	~\$85
Persistent storage (50 GiB)	~\$4 (gp3)	~\$3 (array allocation)
1-year commitment	~\$42	n/a
3-year commitment	~\$30	n/a
Effective \$/vCPU-hour	~\$0.048 on-demand	~\$0.058 all-in

At on-demand rates, AWS wins on paper. At a one-year commitment, it wins comfortably. At three years, it costs roughly a third as much.

And at zero percent utilization, both cost the same as they do at one hundred percent.

That last line is the one that changes the analysis. The on-prem VM costs \$85 a month whether it runs for one hour or for eight hundred. The cloud instance costs \$70 a month on-demand only if it runs continuously; at twenty-five percent utilization, its effective cost is a quarter of that. A reserved instance is the worst of both worlds for a spiky workload and the best of both worlds for a steady one.

This is what unit economics actually reveals: not which side wins, but what the real trade-off is. For workloads that run twenty-four hours a day and are expected to keep running for years — a payments system, an ERP, a core line-of-business application — the on-prem number is competitive and sometimes better than a three-year reserved instance once you factor in egress fees and the premium for managed services. For workloads with variable demand, the cloud's variable cost model wins even when the on-demand sticker is higher. For workloads that sit somewhere in the middle — which is most of them — the decision depends on actual utilization data, and it belongs to whichever team has the numbers to model it.

Before this chapter, most organizations had the numbers for one side of that decision. After it, they have them for both.

What they do not yet have is a way to compare the two sides cleanly. A "production" tag in AWS is not guaranteed to mean the same thing as a "production" tag in vCenter. A cost center in the cloud portal may not reconcile to a cost center in the finance system. Unit economics assumes a shared taxonomy, and hybrid estates rarely have one by accident. That is the next problem.

CHAPTER

Seven

The Tag Mapping Problem

Chapter 6 closed on a technicality that turns out not to be a technicality at all: a "production" tag in AWS is not guaranteed to mean the same thing as a "production" tag in vCenter. You cannot do hybrid showback without a shared taxonomy, and hybrid estates rarely have one by accident.

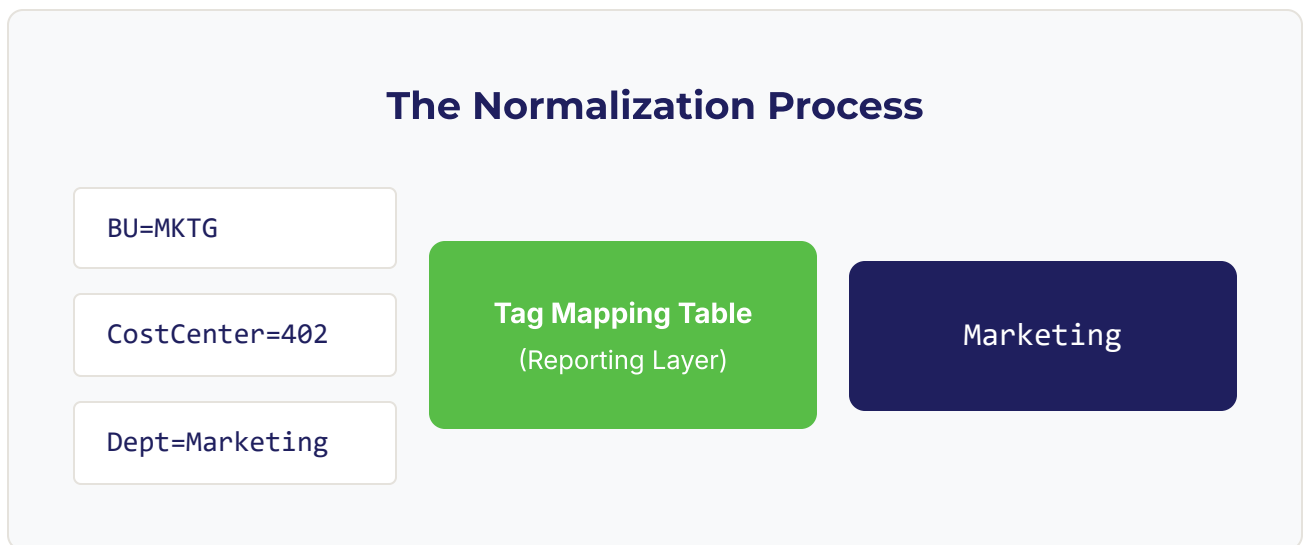
This is the kind of problem that sounds like an afternoon's cleanup and turns out to be a multi-year program. The reason is that tags grow the way cities grow: organically, with no master plan, with street names that made sense to someone twenty years ago, and with four different spellings of the same neighborhood. An organization with a two-decade datacenter history and a ten-year cloud history will have dozens of ways of saying "business unit" across its systems — BU, BusinessUnit, Dept, Department, CostCenter, and a handful of misspellings that someone is still depending on. "Environment" appears as environment, env, deploymentEnv, and — memorably, in at least one real deployment — habitat. Location appears as datacenter, region, geography, and site, often on the same asset.

None of this is anyone's fault. Cloud providers ship their own tag schemas. Virtualization platforms ship theirs. Every acquisition brings its own conventions. Every decade of IT inherits a little more of the ones that came before. The result is that when a FinOps practitioner asks "what

did Marketing spend on infrastructure last month?" the honest answer is not a number but a reconciliation exercise.

Unified tags and foreign tags

The practical pattern that has emerged is a two-layer taxonomy. One internal set of canonical tags — the ones the organization has agreed to actually use for reporting — mapped to an arbitrary number of external tags from every system that touches infrastructure. An internal BU tag maps to BU, BusinessUnit, Dept, Department, and CostCenter across five different systems. Internal Environment maps to env, environment, deploymentEnv, and yes, habitat. The internal layer is small, stable, and governed. The external layer is whatever the source systems emit. The mapping table in the middle is where the work lives.



This pattern is not new — finance teams have been doing it with chart-of-accounts reconciliation since long before anyone had a cloud bill. What is new is that modern hybrid platforms handle this mapping automatically, reconciling unified internal tag keys against foreign tag keys from cloud providers, vendor APIs, and virtualization systems, and exposing a single tag vocabulary to every downstream report. It is the least glamorous feature in the category and the one most likely to determine whether a hybrid FinOps practice actually works.

FOCUS and the direction of travel

The longer-term answer is a shared specification. The FinOps Open Cost and Usage Specification — FOCUS — exists to give cloud providers and vendors a common format for FinOps-serviceable billing data, so that chargeback, allocation, and forecasting don't have to be rebuilt for each provider. It has been adopted by all three major cloud providers and is gaining traction among SaaS and licensing vendors.

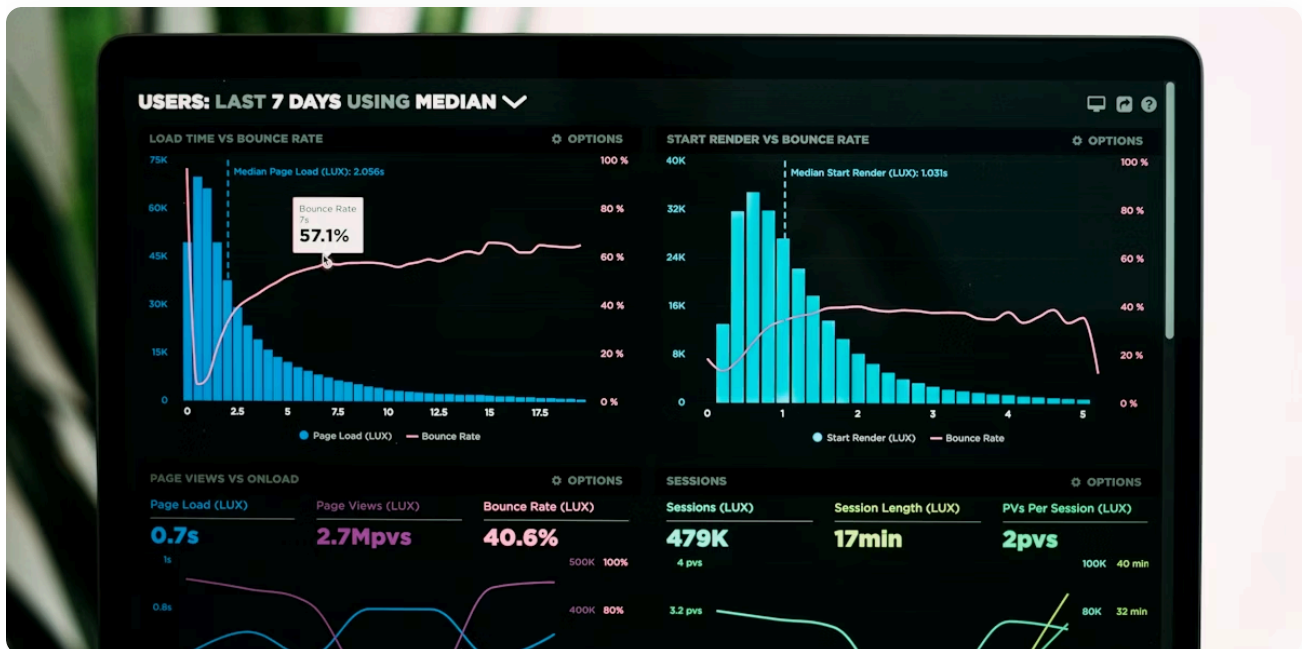
Datacenter FOCUS support is where the conversation is now. In the 2026 State of FinOps survey, datacenter ranked as the #2 most-requested area for FOCUS expansion, behind only AI workloads. The practitioners voting for it are the same ones who have spent the last year reconciling habitat against deploymentEnv. They are not asking for an elegant specification. They are asking for the specification to cover the half of the estate that does not currently emit billing data at all.

Until the specification catches up, the mapping table is the answer. It is unglamorous. It is load-bearing. And it is the precondition for everything the next chapter is about — because a dashboard showing hybrid cost only works if the two halves of the hybrid are speaking the same language.

CHAPTER

Eight

Inform: Visibility Across Estates



"You can't automate what you can't see."

— 2026 State of FinOps survey

The Inform phase is the flywheel's first stop: identify data sources, allocate spend, build a near-real-time picture of what the estate is costing. In cloud, the shape of this work is well-understood. Billing APIs emit data every few hours. Utilization telemetry emits every minute. Anomaly detection runs as a routine service from every major provider. The hard problem is cleaning and presenting the data, not obtaining it.

In the datacenter, the hard problem is obtaining the data. With Chapters 5 through 7 in place — a cost formula, unit economics, and a shared taxonomy — you now have something to show. Showing it is its own discipline.

What near-real-time means when half the estate is not metered

The cloud-native cadence for financial data is hourly or better. That is not the cadence most datacenters can match. Performance telemetry on a storage array or a vSphere cluster refreshes every five minutes, but the cost numbers attached to those assets refresh when someone updates a spreadsheet — or when a facilities invoice arrives thirty days after the month it describes.

The honest hybrid standard is not matching cloud's cadence on the datacenter side. It is deciding what cadence a given decision actually needs. A utilization alert needs minutes. A cost allocation report needs a day. A refresh-cycle decision needs a quarter. Practitioners who try to make everything hourly end up with noisy dashboards and stale data — both at once — and no trust in either. Practitioners who match the cadence to the decision end up with dashboards that people actually check.

The direction of travel is clear. *Dashboards are table stakes of yesterday — reactive*, as one practitioner put it in the 2026 survey. *You have to move to proactive, real-time, automation*. That is true for cloud. For the datacenter, the prior stage — daily, accurate, reconciled — is the one most organizations haven't cleared yet.

Invoice-level data is inadequate everywhere

The instinct when building visibility is to start from the invoice, because the invoice is the one artifact every finance team already has. The instinct is wrong on both halves of the estate. A cloud bill is monthly, aggregated, and arrives after the fact — useful for accounts payable, useless for decisions. A datacenter facility invoice is the same thing with a longer lag and fewer line items. The Inform phase for either cloud or datacenter is about getting below the invoice, to the level where individual workloads, teams, and applications are visible and addressable.

For cloud, "below the invoice" is resource-level billing data joined to tags. For datacenter, it is the OpEx-translated daily cost from Chapter 5, joined to the reconciled tag taxonomy from Chapter 7. Those are not the same artifact, but they are the same category of artifact: a structured, queryable, near-current view of what each unit of infrastructure actually costs, and whose workload is running on it.

Anomaly detection: solved on one side, barely begun on the other

A canonical FinOps story: a remote engineering team at a pharmaceutical company spins up three x1e.32xlarge instances at \$44 an hour, somewhere in a region nobody is watching closely. On a \$3.5M monthly bill, the \$98,000 impact is only a 2% blip — and a simultaneous reserved-instance purchase masks it further. Machine-learning anomaly detection catches it the same day. The instances come down within hours.

That story does not have a datacenter equivalent, and the reason is that anomaly detection on cost requires the cost numbers to move in real-time. In a depreciation-only world, they do not. An over-provisioned cluster running at 15% utilization costs exactly the same as the same cluster running at 85%, because depreciation is indifferent. Only once the OpEx translation is in place — power, personnel, floor space allocated to the workload rather than the asset — does an anomaly actually have a cost signal to ring.

This is the Inform-phase payoff that is hardest to articulate to a skeptical CFO. The translation does not change what the hardware costs the company. It changes what the *workload* costs the company, and workloads are the level at which waste actually happens.

The single pane of glass the market named

If one dashboard shows cloud spend and a different dashboard shows datacenter spend, you do not have hybrid visibility. You have two monologues. The 2026 State of FinOps survey made this unusually concrete: asked which tooling capabilities they wanted that do not currently exist, practitioners ranked *“the single pane of glass for different technology spend”* as their #3 answer, behind only granular AI spend monitoring and pre-deployment architecture costing. The market did not ask for a better cloud tool or a better datacenter tool. It asked for one tool that shows both.

The architectural reason this is hard is that the two halves of the estate produce data on different schedules, in different formats, through different systems. The reason it is worth doing anyway is that every downstream decision — workload placement, refresh-cycle timing, reserved-capacity strategy, chargeback — requires seeing them together.

And once you can see them together, you can start deciding between them. That is Optimize.

CHAPTER

Nine

Optimize: Workload Placement at Parity

Cloud optimization has been won.

That is a strong claim, so let practitioners say it for themselves. From the 2026 State of FinOps survey, one described reaching *“97% optimization in their Cost Optimization Hub, with the remaining 3% intentionally not actioned for business reasons.”* Another said the remaining work had become *“a high volume of smaller opportunities that require more effort to capture.”* An advanced practitioner at a large enterprise was blunter: *“The days of finding something that's grossly misconfigured, and we're gonna save a bunch of money, or there's reserve instances we haven't purchased yet — that was years ago.”*

This is what a mature discipline sounds like. Cloud providers ship rightsizing recommenders, orphaned-resource finders, commitment calculators, and anomaly detectors as native services. Third-party tools refine them further. The big rocks of waste — the forgotten dev environment,

the three \$44-an-hour instances nobody noticed, the zero-utilization reserved instance — have mostly been either moved or accepted. Cloud optimization work continues, but the marginal dollar of effort now returns less than it did three years ago.

And across the datacenter fence, the marginal dollar of effort has barely been spent at all.

What has never been optimized

The datacenter half of most enterprise estates has been running the same way for a decade: provisioned generously at build-out, depreciated on a schedule, and refreshed when the maintenance contract becomes untenable. Utilization is measured, sometimes, by infrastructure teams who care about capacity headroom. It is not measured by anyone who cares about cost, because until Chapter 5 there was no cost signal to measure it against.

With the OpEx translation in place, the same playbook that matured in cloud becomes available on-prem. It just finds different things.

Rightsizing. A typical enterprise vSphere cluster runs at 15% to 30% average CPU utilization. In cloud, that number would trigger a rightsizing recommendation within hours. On-prem, it has been the normal state of affairs for years. Consolidating two clusters running at 20% into one running at 60% is a six-figure annual decision in a mid-market estate — the kind of number cloud FinOps stopped seeing in the big-rocks category a long time ago.

Orphaned resource cleanup. The datacenter cousin of the abandoned EBS volume is the orphaned VMDK: a virtual disk detached from any running VM, sitting on a storage array, quietly consuming capacity that the capacity report still counts as used. Every virtualization estate has thousands of them. The reports that surface them weren't worth building until the cost signal was there. Once it is, an orphaned volume in the cloud is a wake-up call; an orphaned volume in the datacenter is a Tuesday.

Refresh-cycle timing. Commitment-based discounts do not have an on-prem analog — there is no three-year reserved instance for a vSphere host. What does translate is the timing decision: when to retire an aging array versus extend its maintenance contract for another year. With \$/TiB/day for the incumbent and \$/TiB/day for the replacement, the refresh becomes a spreadsheet exercise rather than a negotiation.

Workload placement. The most interesting Optimize decision on a hybrid estate is the one Chapter 6 set up: does this workload belong here, or somewhere else? A steady-state 24/7 database running on a fully-depreciated on-prem host may cost less than the same database on a three-year reserved instance. A spiky development workload running on-prem at peak-sized capacity almost certainly costs more than the variable-rate cloud equivalent. Both comparisons were impossible without unit economics on both sides. Both become routine once they exist.

Repatriation, honestly

The repatriation conversation has been more ideological than analytical for five years. One camp argues cloud is structurally overpriced at scale. The other argues the operational overhead of running infrastructure has been selectively forgotten. Both are sometimes right, depending on the workload.

The unit economics approach replaces the argument with arithmetic. Some workloads come home. Some stay where they are. Some move the other direction, because even at scale the operational premium on certain services is a bargain. What the calculation does not do is produce a single answer for the whole estate — and any vendor, consultant, or analyst claiming otherwise is selling a conclusion the data has not reached.

Where the big rocks are now

The practitioner quotes that opened this chapter describe a cloud optimization frontier where remaining opportunities require more effort than they return. That observation is accurate. It is also local to one half of the estate.

The big rocks of waste have not disappeared from the enterprise. They have moved into the 30%-utilized vSphere cluster nobody has looked at in three years, the storage array whose orphaned capacity equals the workload of an entire business unit, and the refresh cycle that fired on schedule because it always has. None of this is easier than the cloud work that preceded it. It is, however, meaningfully more leveraged per hour of effort than the cloud work currently available.

Making it happen — turning Optimize analysis into Operate practice — is the next chapter. Without that step, the numbers produced here stay on a slide.

CHAPTER

Ten

Operate: Showback, Chargeback, and the Culture Shift

The Operate phase is where the numbers from the first nine chapters either become a practice or stay on a slide. Inform produced the visibility. Optimize produced the decisions. Operate is the work of making those decisions happen month after month — reporting them, acting on them, billing for them, and convincing the rest of the organization that any of this is worth the disruption.

The Foundation's own framing of the Operate phase is that it is "where the rubber meets the road," and its curriculum is unusually direct about the implication: this phase is about people and processes, not technology. That is equally true in hybrid, with one compounding complication. In cloud-only FinOps, the process change is happening inside a single team's center of gravity. In hybrid FinOps, the process change is happening across the boundary between two teams that have historically reported through different parts of the org chart, spoken different vocabularies, and been measured against entirely different KPIs. The next chapter is about that boundary. This one is about what the boundary has to produce.

Showback before chargeback, always

The distinction is older than FinOps. Showback shows a team what they are spending; chargeback debits the same spending against their P&L. One is a conversation. The other is a bill.

Every hybrid practice will need showback. Some will also need chargeback — and, per the Foundation's curriculum, the deciding factor is usually accounting requirements rather than FinOps maturity. Legal, tax, and regulatory structure push some organizations to chargeback early; others run on showback indefinitely without any loss of rigor. There is no gold medal for chargeback. There is, however, a significant penalty for rolling it out before the underlying data is trusted. *Angry accountants do not make good FinOps partners*, as the Foundation's curriculum puts it — and chargeback data the business does not believe produces angry accountants faster than almost anything else.

The three-year story

The canonical chargeback transition in FinOps lore was documented by Rob Martin, working with a large enterprise whose cloud-only chargeback build took three fiscal years end to end. Year one was exploratory: learning the cloud, getting the first applications instrumented, doing very little formal budgeting. Year two introduced real accountability — teams operated against traditional budgets, and FinOps provided showback that teams could use to understand their spend without yet being billed for it. Year three was when the organization defined how cloud teams would control budgets directly, Finance got comfortable that P&L owners could be charged accurately, and true chargeback went live.

Three years. For cloud-only. With a billing API providing clean input data from day one.

Hybrid chargeback, done honestly, is the same transition repeated on the half of the estate that does not have a billing API — while the first transition is either still in progress or freshly completed. Anyone selling a shorter timeline is selling a plan, not a practice.

What changes in hybrid

Two things, specifically.

The first is **cadence**. Cloud billing data lands within hours; datacenter cost inputs land on a thirty-day lag at best. Monthly hybrid chargeback is realistic. Weekly is theater. The cadence has to match the slowest input, not the fastest, or the report's credibility dies with the first reconciliation argument.

The second is **precision**. Cloud chargeback numbers can be exact because the provider invoiced an exact number. Datacenter chargeback numbers include allocated estimates — floor space divided across racks, personnel cost divided across devices, power attributed by nameplate draw rather than metered per asset. The finance team has to accept *approximately right* as the operating standard. Every attempt to make it *precisely right* adds months, costs credibility, and eventually produces a number no more accurate than the estimate it replaced.

Crawl, Walk, Run — for hybrid chargeback

Crawl. Cloud is on chargeback or mature showback. Datacenter is on showback only, produced at a monthly cadence from the OpEx translation in Chapter 5, reconciled to the tag taxonomy in Chapter 7. Finance sees two reports and accepts that only one of them debits a P&L.

Walk. Both halves are on showback within a single report, using a consistent allocation methodology agreed between FinOps, IT, and Finance. P&L owners can see their full hybrid footprint; the accounting ledger still distinguishes the two halves, but the business conversation treats them as one number.

Run. Both halves are on chargeback, reconciled monthly into the same P&L, with allocated datacenter components accepted as "approximately right" by Finance under a documented methodology. Shared-cost rules are explicit. The practice has crossed the line from informational to consequential.

Most organizations will spend longer at Walk than at any other stage, and that is the correct place to spend time. Walk is the stage where the business, finance, and infrastructure teams build the trust the Run stage requires. Skipping it produces the kind of chargeback rollout that generates lawsuits from business units who refuse to accept numbers they had no part in agreeing to.

The hardest part of Operate is not any of the above. It is that two groups of people who have historically not been in the same meetings now have to be — looking at the same dashboard, agreeing on the same methodology, and jointly accountable for the same number. That is the next chapter.

CHAPTER

Eleven

Getting FinOps and IT to Actually Cooperate



The canonical story about FinOps cultural failure was told by a practitioner at Atlassian. He built his first set of reports, pushed them out to the business, and watched them fail in real time. *"My reports didn't clarify cloud spend — they confused it."* The terms he used didn't land for Finance; the financial framings didn't land for Engineering; every audience read the report through its own lens and reached a different conclusion about whether the report was even correct. The solution, when it eventually came, was a common vocabulary — not a better dashboard.

That story is a decade old and has aged into a discipline-wide cliché. It is still a cliché because the failure mode it describes has not gone away. It has just moved. In hybrid FinOps, the vocabulary problem is no longer between FinOps and Finance, which have mostly figured each

other out. It is between FinOps and the people who actually run the datacenter.

Three rooms, three vocabularies

Historically, there have been three groups in a large enterprise looking at technology spend, and they have never been in the same meeting.

Finance has its ledgers, its depreciation schedules, its quarterly close. It measures infrastructure in CapEx and OpEx, thinks in months and quarters, and is allergic to any number whose provenance it cannot trace.

IT — specifically, infrastructure engineering — has its capacity dashboards, its utilization charts, its health alerts. It measures infrastructure in TiB, cores, IOPS, and percent utilized. It thinks in minutes and hours at the operational layer, and in three-to-five-year horizons at the refresh-planning layer. It is allergic to anyone asking for a cost number on a device it just installed.

FinOps, born in cloud, arrived speaking a third dialect again. It measures in \$/vCPU/hour, thinks in tag taxonomies, and reports on cadences that neither Finance nor IT had ever used for anything. It was welcomed, more or less, by Finance — which recognized a kindred spirit and appreciated having someone else to explain the cloud bill. It was regarded more warily by IT, which was used to owning the infrastructure narrative and did not necessarily want a new team showing up with its own vocabulary and its own scorecard.

That wariness is the cultural problem hybrid FinOps has to solve. The solution starts with an organizational shift.

The 78%

The 2026 State of FinOps survey put a number on something that had been quietly happening for years. Seventy-eight percent of FinOps practices now report into the CTO or CIO, up eighteen points since 2023. The CFO reporting line has dropped to eight percent.

The reporting line matters because coordination is harder than it looks from outside it. The practitioner trying to extend FinOps into the datacenter is asking procurement, ITAM, facilities, and the vendor-management function for data those teams never collected for this purpose — and, more to the point, data they have no incentive to hand over on anyone else's timeline.

Facilities reports to real estate or operations; its KPIs are uptime and power efficiency, not cost allocation by workload. ITAM reports to IT governance; its mandate is inventory accuracy and compliance, not financial modeling. Procurement reports to finance but measures itself on negotiated savings and contract cycle time, not on whether a purchase record is retrievable four years later. None of these teams are wrong to prioritize what they prioritize. They are simply not, in any structural sense, part of the FinOps team's workflow. When FinOps reported to the CFO, the ask to collaborate was a lateral request across peer organizations. When FinOps reports to the CTO, the ask is a request from the same executive's staff — and executive-level requests move differently than peer-level ones, regardless of how politely they are phrased.

This is a bigger deal than it looks. When FinOps reported to Finance, a FinOps/IT relationship was a cross-organizational negotiation — two executives had to agree on the meeting, the agenda, the data, and the mandate before the teams below them could reliably cooperate. When FinOps reports to the CTO, the two teams are under the same executive, measured on overlapping KPIs, and present in the same reviews by default. The vocabulary problem does not vanish, but the meeting where the vocabularies have to reconcile is scheduled recurrently rather than convened under duress.

This is the structural precondition for hybrid FinOps. It does not guarantee success — the survey does not report on that — but it explains why hybrid scope is rising so fast now when the same conversation was stuck for years. The executive who owns both halves of the estate has a natural interest in having both halves measured consistently.

The personas still want different things

Reporting structure is not the same as shared understanding. Even in an organization where FinOps and IT sit under the same executive, the individuals involved want different things from the same data, and the dashboards that fail are the ones that try to give everyone the same view.

Infrastructure engineers want capacity, utilization, and health — they look daily, at the device level, and they care about what is running out and what is broken. FinOps practitioners want cost, allocation, and forecast — they look weekly or monthly, at the workload or portfolio level, and they care about where the dollars are flowing. Finance wants invoice reconciliation and budget variance — monthly or quarterly, portfolio-wide, traceable to the general ledger. Leadership wants trends and headline numbers — quarterly, total, and explainable in one slide.

The Foundation's curriculum on personas is explicit about this. Each persona has its own cadence, scope, and KPI set. The historical failure mode has been to build one dashboard per audience, in one tool per audience, with no reliable way to reconcile them when they disagreed. The next-generation pattern — and the practical instantiation of the principle that *“FinOps should be enabled centrally”* — is a single underlying dataset exposed through views tuned to each audience. The infrastructure engineer sees the capacity and health view. The FinOps practitioner sees the cost and forecast view. When they disagree about a number, both can drill into it from their native starting point and arrive at the same row of the same table.

The tools are not the breakthrough. The breakthrough is that both audiences are now reading the same underlying facts. Everything that used to be an argument about whose data was right becomes a conversation about what to do next.

The meeting the dashboard forces

The cultural breakthrough of hybrid FinOps is not the dashboard. It is the first meeting in which an infrastructure engineer and a FinOps practitioner look at the same storage array on the same screen and talk about what to do about it. In most organizations that meeting has never happened. After it starts happening — weekly, monthly, on whatever cadence the estate warrants — the rest of the practice starts building itself.

Which is, finally, a subject the next chapter can frame as a maturity model. Because once the meeting is happening, the question is no longer whether the organization has hybrid FinOps. It is how fast it can iterate.

CHAPTER

Twelve

A Maturity Model for Hybrid FinOps

The FinOps Foundation's Crawl/Walk/Run framing has survived every rebrand and every scope expansion the discipline has been through, because it reflects something true about how organizations actually learn new operational practices: slowly, in stages, with no reliable way to skip steps. The hybrid adaptation is the same three columns applied across the capabilities this book has built. Five rows, three stages, and one honest answer about where a given organization currently stands.

	Crawl	Walk	Run
Asset estimation	Manual sourcing of purchase price, power, floor space, personnel, and maintenance inputs from multiple owners and systems. Six to twelve months of spreadsheet work before the first usable cost number appears.	Inputs retrieved on device identity in days rather than months. The formula runs against data that is substantially complete from day one.	Inputs continuously refreshed and validated against actual telemetry where available. Estimation is a maintenance function rather than a project.
Cost allocation	Daily rates computed manually, refreshed monthly. Tag mapping done by lookup, reconciled during finance close.	Daily rates computed automatically from the OpEx formula. Unified/foreign tag reconciliation in place. Monthly close runs without argument.	Continuous cost allocation. Tag governance enforced at provisioning. Shared-cost rules documented, versioned, and applied automatically.
Unit economics	\$/VM and \$/TiB computed for the largest workloads. Methodology inconsistent across applications.	Consistent unit economics across the estate. Cloud and on-prem comparable in the same taxonomy. Workload placement decisions reference the numbers.	Unit economics flow into architecture and placement decisions pre-deployment. New workloads costed before they are built.
Reporting cadence	Monthly reports, produced manually,	Weekly dashboards reviewed in recurring	Cost data refreshed daily; utilization and

	Crawl	Walk	Run
	read primarily by the FinOps team.	cross-functional meetings. IT and FinOps looking at the same underlying data.	anomalies in near-real-time. Persona-specific views. Disagreements resolved by drilling in, not by reconciling spreadsheets.
Chargeback	Cloud on chargeback or mature showback. Datacenter shown only in aggregate to Finance.	Both halves on showback within a single report using an agreed allocation methodology. P&L owners see their full hybrid footprint.	True hybrid chargeback. Both halves reconciled monthly into the same P&L. Shared-cost rules explicit. Approximately right accepted as the operating standard.

The row that collapsed

Every row in this table except one still obeys the historical timeline: Crawl is six to eighteen months of organizational change, Walk is another year or more of pattern refinement, Run is where a few organizations eventually arrive and most continue to aspire. Reporting cadence, chargeback, unit economics, and the cultural work underneath them cannot be shortcut. They involve people agreeing to new methodologies, Finance learning new vocabulary, IT and FinOps building trust in each other's numbers. These things take the time they take.

The asset estimation row is different. AI-assisted estimation now retrieves purchase dates, list prices, power draw, floor space, and personnel inputs on sight — which is to say, organizations adopting the capability skip Crawl entirely on that row and start at Walk on day one. The savings compounds. A practice that would have spent six months in Crawl on estimation can spend those

six months on the rows that actually require cultural work: getting IT and FinOps into the same meeting, negotiating allocation methodology with Finance, and producing the first hybrid chargeback report that survives its first audit.

This is the only row where the calendar has moved. The rest of the table still describes a multi-year practice. But a multi-year practice that starts a half-year ahead of where it used to start is a meaningfully different practice, and the organizations currently leading in hybrid FinOps are disproportionately the ones that figured this out first.

The obstacle that has kept FinOps out of the datacenter has never been the math. Chapter 4 decomposed the bill. Chapter 5 gave the formula. Any competent practitioner, given a spreadsheet and a week, can do the OpEx translation for a single asset. The obstacle has been the coordination — the five-way cross-functional negotiation required to gather the inputs before the math could even begin. That obstacle, finally, is removable. Which means the scope expansion the Foundation ratified in its mission change is no longer a future-state ambition. It is an operational possibility, today, for any organization willing to begin.

The maturity model is still a ladder. One rung of it now arrives assembled.

APPENDIX

A

FinOps Tool Categories and Representative Vendors

This guide was produced by Visual One Intelligence® as a contribution to the practitioner community. The vendor list below is provided for reference and reflects the authors' understanding of the current market as of publication. Inclusion is not an endorsement; exclusion is not a judgment. The FinOps tooling landscape evolves quickly, and practitioners are encouraged to verify current capabilities directly with vendors.

The tooling landscape for FinOps has grown in step with the discipline itself. Five years ago, most organizations could cover their needs with a native cloud provider cost console and a spreadsheet. Today, a mature FinOps practice typically draws on three to five tools across different categories, reflecting the scope expansion documented in the FinOps Foundation's 2026 State of FinOps survey — where datacenter, SaaS, licensing, and AI spend have all moved into the FinOps remit.

The categories below are organized by capability rather than by market tier. Representative vendors within each category are listed alphabetically, not by preference.

Cloud Provider Native Tools

The starting point for most FinOps practices. These tools are free or low-cost, deeply integrated with their respective cloud providers, and cover the foundational capabilities needed in the early Inform phase. Their limitation is scope — each covers only its own cloud, which makes multi-cloud and hybrid analysis difficult.

Supports: Data Ingestion, Reporting & Analytics, basic Anomaly Management, basic Allocation.

- **AWS Cost Explorer and AWS Cost and Usage Reports** — usage visualization, cost allocation by tag, savings plan and reserved instance recommendations for AWS workloads.
- **Azure Cost Management and Billing** — spend analysis, budget alerts, and cost allocation for Azure, with cross-subscription reporting.
- **Google Cloud Cost Management** — billing export to BigQuery, budget controls, and cost breakdown for GCP workloads.
- **Oracle Cloud Cost Analysis** — cost reporting and usage analysis for OCI environments.

Third-Party FinOps Platforms (Cloud-Focused)

Commercial platforms that extend beyond native provider tooling with multi-cloud support, more sophisticated allocation logic, and richer optimization recommendations. These are the tools most practitioners adopt when native tools become insufficient — typically at the Walk stage of maturity.

Supports: Allocation, Reporting & Analytics, Forecasting, Rate Optimization, Workload Optimization, Anomaly Management.

- **Apptio Cloudability** — multi-cloud cost management, allocation, and optimization with deep integration into IBM's broader ITFM portfolio.
- **CloudHealth (by Broadcom)** — multi-cloud governance, cost management, and policy automation.

- **Finout** — unified cost management across cloud, Kubernetes, SaaS, and data platforms with a MegaBill abstraction layer.
- **Flexera One** — cloud cost optimization integrated with ITAM and software asset management capabilities.
- **Vantage** — cost reporting and optimization focused on developer-friendly workflows across AWS, GCP, Azure, and adjacent services.

Hybrid and On-Premises FinOps Platforms

Platforms that extend FinOps practices into the datacenter, private cloud, and virtualization environments. This is the newest category in the landscape, responding directly to the scope expansion documented in the 2026 State of FinOps survey. The capabilities emphasized here are asset cost modeling, cross-estate allocation, and bridging the CapEx-to-OpEx translation that this book has spent thirteen chapters describing.

Supports: Allocation, Planning & Estimating, Forecasting, Reporting & Analytics, Workload Optimization, Intersecting Disciplines (particularly ITAM and ITFM).

- **IBM Turbonomic** — application resource management across cloud and on-premises with automated rightsizing and workload placement.
- **VMware Aria Cost (formerly CloudHealth Hybrid)** — multi-cloud and hybrid cost management with native integration into VMware environments.
- **Virtana** — infrastructure performance and cost analytics across hybrid estates, with workload placement recommendations.
- **Visual One Intelligence®** — unified analytics and FinOps reporting across cloud, virtualization, and on-premises storage infrastructure, with AI-assisted asset cost estimation for rapid hybrid-estate onboarding.

ITAM, ITFM, and Adjacent Disciplines

Tools from neighboring disciplines that increasingly intersect with FinOps practice. The 2026 State of FinOps survey identified IT Financial Management and IT Asset Management as the two most common collaboration partners for FinOps teams — a reflection of the scope expansion into licensing, SaaS, and hybrid infrastructure.

Supports: Licensing & SaaS, Intersecting Disciplines, Allocation, Invoicing & Chargeback.

- **Apptio (ITFM)** — IT financial management platform covering technology spend allocation, benchmarking, and planning.
- **Flexera (ITAM/SAM)** — software asset management and licensing optimization across on-premises and SaaS.
- **ServiceNow IT Asset Management** — asset lifecycle management integrated with broader ITSM workflows.
- **Snow Software** — software asset management with usage analytics and license optimization.
- **Zylo** — SaaS management and optimization.

FOCUS-Aligned Data and Specialty Tools

Tools that address specific capabilities or align with the FinOps Open Cost and Usage Specification (FOCUS) — the emerging standard for normalized cost and usage data across providers. This category is small today but growing, reflecting the 2026 survey's finding that data center and AI workloads are the top two expansion priorities for FOCUS adoption.

Supports: Data Ingestion, Reporting & Analytics, capability-specific needs.

- **Anodot** — anomaly detection and forecasting with FinOps-specific modeling.
- **Cletrics** — the primary solution for real-time cloud cost alerting and monitoring.
- **Kubecost** — Kubernetes cost monitoring and allocation.
- **Pelanor** — FOCUS-aligned multi-cloud cost intelligence.
- **Tangoe** — telecom, mobile, and connectivity expense management.

Open Source and Community Tools

Open-source projects that support FinOps practices, typically used by engineering-led teams or as components within larger toolchains. These are worth knowing regardless of whether an organization uses commercial tooling — several of them define standards the commercial vendors are now aligning to.

- **OpenCost** — CNCF project for Kubernetes cost monitoring; the open-source foundation underlying Kubecost.
- **FOCUS (FinOps Open Cost and Usage Specification)** — an open specification maintained by the FinOps Foundation for normalized billing data across cloud providers.
- **Cloud Custodian** — rules engine for cloud governance, including cost-related policy enforcement.
- **Infracost** — pre-deployment cost estimation for Terraform and infrastructure-as-code workflows.

Choosing Tools: A Brief Note on Sequencing

Practitioners frequently ask which tool to adopt first. The honest answer is that the sequence matters less than the discipline behind it. Most organizations that achieve mature FinOps practices started with whatever native tooling came with their largest cloud provider, added a third-party platform when multi-cloud or allocation complexity outgrew the native tools, and layered hybrid or specialty tools as the scope expanded into new domains.

The sequence recommended by the FinOps Foundation's curriculum — *understand cost, then optimize cost, then operate* — applies equally to tool selection. A tool that accelerates understanding is worth more at the beginning than one that promises optimization. A tool that promises optimization is worth more once allocation and reporting are in place. A tool that supports both cloud and hybrid scope is worth more than two separate tools once the organization has crossed the threshold this book describes.

Beyond the tool itself, the single most important factor in practitioner success is executive sponsorship — a finding consistent across multiple years of the State of FinOps survey. Tools do not produce outcomes; teams produce outcomes, and tools accelerate the teams that are positioned to succeed.

APPENDIX

B

Glossary of Key Terms

Terms introduced or used in specialized ways throughout this guide. For standard FinOps terminology not defined here, see the FinOps Foundation's official glossary at finops.org.

Amortized Purchase Price — The capital cost of a hardware asset divided by its useful life, expressed as a daily or monthly rate. The first of the five line items in the datacenter bill constructed in Chapter 4. For a \$50,000 server on a five-year refresh, the amortized purchase price is approximately \$27 per day.

CapEx Trap — The systematic distortion of workload-placement and refresh decisions caused by treating on-premises infrastructure as "already paid for." The trap operates by making visible only the depreciation component of true cost — typically less than a third of the total — while leaving operating cost components (power, floor space, personnel, maintenance) scattered across separate ledgers. Introduced in Chapter 3.

Cold Start Problem — The practical difficulty of sourcing the five input values required to compute OpEx cost for hardware an organization did not purchase recently: manufacturing date, original list price, power draw, floor space allocation, and personnel cost. Traditionally resolved through multi-month data-gathering projects; increasingly resolved through AI-assisted asset estimation platforms. Discussed in Chapter 5.

Datacenter Bill — A construct, not a document. The assembled view of all cost components attributable to a specific piece of on-premises infrastructure, normalized to a daily or monthly rate. Unlike a cloud bill, the datacenter bill does not arrive monthly; it must be built from data that lives in separate systems. Chapter 4 describes its five line items.

FOCUS (FinOps Open Cost and Usage Specification) — An open specification maintained by the FinOps Foundation for normalized cost and usage data across cloud providers. Adoption is growing as FinOps scope expands; datacenter and AI workloads are the top two most-requested FOCUS expansion areas in the 2026 State of FinOps survey.

Hybrid FinOps — The application of FinOps principles, phases, and capabilities across both cloud and on-premises infrastructure in a single unified practice. Distinguished from "cloud FinOps" by the requirement that on-premises assets be translated into the same financial language used for cloud — an OpEx-shaped view of CapEx-shaped assets. The subject of this guide.

OpEx Translation — The conceptual and mathematical process of converting an on-premises asset's cost profile from a CapEx representation (purchase price, depreciation schedule, capital budget line) into an OpEx representation (daily run-rate, unit cost, allocable to workload). The formula introduced in Chapter 5.

Showback — Reporting cloud or datacenter costs to the teams that consume them, without formally charging those costs against team budgets. Informational, not transactional. Typically the first stage of chargeback maturity.

Chargeback — Formally allocating cloud or datacenter costs against specific team budgets or P&L lines, with costs recognized in the accounting system. A multi-year organizational transition for most organizations, not a technical rollout.

Single Pane of Glass — A unified view of technology spend across cloud, datacenter, virtualization, SaaS, and other domains in a single reporting interface. Identified as the #3 most-requested tooling capability that does not currently exist in the 2026 State of FinOps survey.

Unit Economics — The practice of expressing technology cost as a rate per unit of business output — cost per transaction, cost per customer, cost per ride, cost per analyzed document. Extended in Chapter 6 to on-premises infrastructure via the \$/day, \$/GiB/day, \$/core/day, \$/VM normalization.

Useful Life — The accounting horizon over which a hardware asset is depreciated. Industry convention for enterprise infrastructure is five years; some organizations use three or seven. Distinct from the *refresh cycle*, which is the operational horizon at which the asset is actually replaced. The two frequently diverge, and assets running past their useful life continue to consume operating cost without contributing to amortization.

Utilization — The proportion of an asset's total capacity that is actively delivering useful work. The third multiplier in datacenter cost calculations, and the single largest determinant of per-workload unit economics. A server at twenty percent utilization has the same daily total cost as one at eighty percent, but delivers a quarter of the effective compute.

About This Guide

The Definitive Guide to Hybrid FinOps is a HybridFinOps Publication, sponsored by Visual One Intelligence® as a contribution to the FinOps practitioner community. The guide was commissioned in response to a specific gap in the available literature: the discipline of FinOps has expanded officially to include the datacenter, yet no comprehensive reference has existed for practitioners who now have to apply public-cloud financial metrics to on-premises infrastructure.

The framework described throughout this book — the six principles, the Inform/Optimize/Operate phases, the Crawl/Walk/Run maturity model, the domain and capability taxonomy — is the work of the FinOps Foundation, a directed project of the Linux Foundation. The Foundation represents more than ten thousand practitioners and over thirty-five hundred member organizations. Its 2026

State of FinOps survey, cited throughout this guide, captured the perspectives of 1,192 practitioners overseeing more than \$83 billion in annual cloud spend. Readers seeking the authoritative source for FinOps terminology, certification, and community resources should visit finops.org.

This guide's contribution is translational rather than definitional. It takes the Foundation's existing framework and applies it to the specific case of hybrid infrastructure — where the familiar shape of the cloud bill meets the unfamiliar assembly required to produce an equivalent view of on-premises assets. The thirteen chapters cover the conceptual argument, the mathematical translation, the practical application across the three FinOps phases, and the cultural transition required to bring IT and FinOps teams into shared language.

About Visual One Intelligence®

Visual One Intelligence builds unified analytics and reporting for hybrid technology infrastructure. The platform covers cloud, virtualization, and on-premises storage environments, presenting a consistent view of cost, capacity, and utilization across estates that traditionally lived in separate tools. Its AI-assisted asset estimation feature — which retrieves manufacturing dates, list prices, power draw, floor space, and personnel cost ranges for enterprise hardware on sight — is designed specifically to address the cold-start data problem this guide has described at length.

More information is available at visualoneintelligence.com/hybrid-finops/.

A Note on Sources

The Foundation's curriculum and the 2026 State of FinOps survey are the authoritative sources for all framework content, statistics, and practitioner quotes referenced throughout this guide. Quotes attributed to "a practitioner" or "respondents to the survey" draw from publicly available Foundation publications. The representative vendors listed in Appendix A reflect the public tooling landscape as of publication and should not be read as exhaustive or as endorsements.

The opinions expressed in the guide — particularly the central argument that OpEx translation is the only realistic measurement framework for hybrid FinOps — are those of the guide's authors and do not necessarily reflect positions held by the FinOps Foundation.

Rights and Usage

This guide is made freely available for distribution within the practitioner community. It may be shared, excerpted with attribution, and referenced in internal organizational materials. It may not be resold, repackaged as another vendor's content, or reproduced in whole without written permission from Visual One Intelligence®.