

# The Hybrid Equilibrium: A Practical Guide to Optimizing Costs and Performance with Bare Metal and Cloud.

Farhan Javed<sup>1</sup>

farhan.jaaved@outlook.com

<sup>1</sup>Kalyani University

7 July, 2025

## Abstract

For over a decade, a "cloud-first" doctrine has dominated IT strategy, promising unparalleled scalability and ease of use. While the cloud remains a revolutionary force for innovation, a growing number of organizations are awakening to the significant, often unsustainable, costs associated with a purely cloud-based infrastructure. This paper argues not against the cloud, but for a more mature, financially-aware approach to infrastructure—a practice we call being "cloud-smart." We will explore a hybrid equilibrium where workloads are placed on the most logical platform. This involves leveraging bare metal for its raw performance and superior cost-efficiency for stable, predictable workloads, while strategically utilizing the cloud for its elasticity and high-value specialized services. Through a detailed analysis of real-world case studies, including the dramatic migrations of Dukaan, Basecamp, and OneUptime, and an examination of modern management tools that have democratized bare metal operations, this paper serves as a practical guide for CTOs, DevOps engineers, and FinOps professionals. Our goal is to empower companies to move beyond default choices and implement superior infrastructure strategies that drive both technical excellence and long-term financial sustainability.

## 1 Introduction

The advent of hyperscale cloud providers like Amazon Web Services (AWS), Google Cloud Platform (GCP), and Microsoft Azure was nothing short of a paradigm shift. They democratized access to enterprise-grade infrastructure, allowing startups to compete with incumbents by trading large, upfront capital expenditures (CapEx) for a flexible, operational expenditure (OpEx) model. For businesses navigating the uncertainty of product-market fit and unknown scaling trajectories, the cloud was—and still is—an indispensable catalyst for innovation.

However, as businesses mature and their workloads stabilize, the very OpEx model that once provided flexibility can become a financial liability. The layers of abstraction that make the cloud easy to use also obscure its true cost, leading to a phenomenon of unchecked spending. This isn't a niche problem. The research firm IDC estimates that a staggering 20-30% of all cloud spending is wasted (IDC, 2023)<sup>1</sup>. This waste stems from over-provisioning, idle resources, suboptimal service selection, and punitive data egress fees—a slow, persistent drain on capital that could be reinvested into product development, talent, or growth.

The response to this challenge is FinOps—a cultural and operational practice that brings financial accountability to the variable spend model of the cloud. This paper extends that philosophy, arguing that a truly

---

<sup>1</sup>Source: IDC Future Enterprise Planning Guides, Control Cloud Costs and Expand Transparency with FinOps, IDC #US50654223, May 19, 2023

mature FinOps strategy must look beyond just optimizing cloud spend and ask a more fundamental question: Is the cloud the right home for every workload?

## 2 The Bare Metal Advantage: Deconstructing Performance

To grasp the "why" behind a bare metal strategy, one must first deconstruct performance and understand the architectural differences between a cloud VM and a dedicated bare metal server.

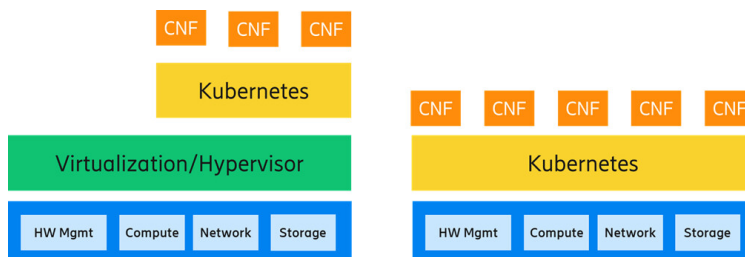


Figure 1: Kubernetes over virtualized infrastructure versus bare metal infrastructure

The core difference is the hypervisor. In a cloud environment, the hypervisor is a software layer that sits between the physical hardware and the operating system. Its job is to partition the physical resources—CPU cores, RAM, storage, and network—into smaller virtual slices (VMs) to be sold to multiple tenants. While a marvel of engineering, this layer introduces two significant performance penalties:

1. **The Hypervisor Tax (Virtualization Overhead):** The hypervisor itself consumes computational resources to do its job. A portion of the CPU cycles and RAM of the physical machine is reserved for managing the various VMs. This "tax" means that even if you are paying for 4 vCPUs, you are not getting the full, unadulterated power of four physical CPU cores. This overhead impacts every single operation your application performs, creating a constant drag on performance.
2. **The "Noisy Neighbor" Effect:** In a multi-tenant environment, you are sharing the underlying physical infrastructure—disk controllers, network cards, CPU cache—with other customers. If one of these "neighbors" runs a highly intensive workload (like a large data processing job), it can monopolize these shared resources, causing your application's performance to degrade unpredictably. You can experience sudden increases in latency or drops in throughput for reasons that are completely outside of your control and invisible to your monitoring tools.

Bare metal eliminates both of these problems. By giving you direct, uncontested access to the entire physical machine, you get:

- **Raw I/O Performance:** For databases and applications that depend on fast storage, this is a game-changer. Direct access to high-speed NVMe SSDs without a hypervisor bottleneck results in dramatically higher Input/Output Operations Per Second (IOPS) and lower latency. Database queries run faster, APIs respond quicker, and user experience improves.

A detailed performance comparison conducted by The New Stack confirms this, showing that for CPU, RAM, storage, and network-intensive tasks, Kubernetes clusters running on bare metal consistently and measurably outperform identical clusters running on virtual machines (The New Stack, n.d.).

While some studies suggest that for certain High-Performance Computing (HPC) workflows, cloud platforms can deliver comparable performance when hardware is identical and virtualization overhead is low, this often comes at a premium price and doesn't negate the inherent overhead for the vast majority of common business applications

### 3 Case Studies in Cloud Repatriation: The Overwhelming Financial Argument

The theory of cost savings and performance gains is powerful, but real-world data from companies that have made the switch provides undeniable proof.

#### Case Study 1: Dukaan's Radical 50x Cost Reduction

- **The Problem:** Dukaan, a fast-growing e-commerce platform, was a model cloud user. They leveraged the scalability of their cloud provider to grow rapidly. However, this growth came with an astronomical price tag. Their cloud bills were soaring to unsustainable levels with their peak at \$90,000 per month, eating directly into their margins. Performance, particularly for their databases, was also a concern due to the inherent limitations of virtualized disk I/O.
- **The Solution:** Dukaan made the bold decision to repatriate their core application infrastructure to bare metal. They didn't abandon the cloud entirely; instead, they adopted a surgically precise hybrid model. Their primary application and databases now run on high-performance bare metal servers, taking full advantage of the raw disk speed. For disaster recovery—a perfect use case for the cloud's strengths—they continue to use a cloud provider for Point-in-Time Recovery (PITR) backups with real-time replication and archival.
- **The Results:** The outcome was transformative. Dukaan slashed its infrastructure spending by **nearly 50x from its peak**, a figure that fundamentally changed their business economics. Now their monthly costs are around \$1500 per month. Furthermore, they saw significant performance improvements across their platform. They even engineered a clever solution for peak traffic: by distributing their bare metal servers globally, they can use monitoring to detect which region is experiencing night-time (low traffic) and use a proxy to route peak traffic from another region to that underutilized hardware, effectively load-balancing globally without spinning up costly new instances.

#### Case Study 2: Basecamp's Principled and Profitable Cloud Exit

- **The Problem:** 37signals, the company behind the popular project management tool Basecamp, is known for its principled approach to business. After years on the cloud, they publicly detailed their concerns, which were less about a single technical issue and more about the long-term financial model. They felt they were "renting computers at an absurd markup" and that the complexity of cloud service offerings often obscured the true cost.
- **The Solution:** They embarked on a well-documented journey to move their entire suite of applications off the cloud and onto their own purchased hardware, co-located in data centers. They invested in modern Dell servers with top-tier specifications (e.g., dual 32-core AMD EPYC CPUs, 1.5TB of RAM) for a fraction of what comparable cloud instances would cost over a few years.
- **The Results:** The financial impact was staggering. Basecamp is on track to save approximately \$10 million over five years. This wasn't a one-time saving; it's a permanent reduction in their operational overhead. Their story became a landmark case, proving that even for a highly successful, modern software company, the economics of owning (or renting bare metal) can be overwhelmingly superior to renting virtual machines in the long run.

### Case Study 3: OneUptime’s Lean 90% Savings

- **The Problem:** As an observability platform, OneUptime processes a significant amount of data. Running their stack on AWS was functional but expensive, directly impacting their ability to offer competitive pricing for their open-source tool.
- **The Solution:** The team migrated their entire stack from AWS to a BMaaS (Bare Metal as a Service) provider. This allowed them to get the performance and cost benefits of dedicated hardware without the large upfront capital expenditure of buying servers.
- **The Results:** As documented in their blog, the move resulted in a **90% reduction in their monthly infrastructure costs** (OneUptime, 2023). This not only improved their profitability but also empowered their team with deeper knowledge of their own infrastructure, leading to better optimization and a more resilient product.

## 4 Debunking the Myths: Modern Management and True Reliability

### Myth 1: Bare Metal is a Management Nightmare

The evolution of the cloud-native ecosystem has provided the tools to manage bare metal with the same agility as the cloud.

- **Kubernetes as the Great Equalizer:** The most significant development is Kubernetes. It provides a universal abstraction layer for deploying and managing containerized applications. Running Kubernetes on bare metal provides a unified, cloud-native operational experience (Ericsson, 2022). Your DevOps team doesn’t need to learn a new “bare metal” workflow as Kubernetes is cloud-agnostic. They use the same kubectl commands, the same YAML manifests, and the same container images. They gain the performance and cost benefits of the underlying hardware without sacrificing the agility of container orchestration.
- **A Rich Ecosystem of Modern Tooling:** The complexity of manual configuration is a thing of the past.
  - **CI/CD Pipelines:** Modern CI/CD tools (like GitLab CI, Jenkins, GitHub Actions) integrate seamlessly with Kubernetes on bare metal, automating everything from code commit to testing and deployment.
  - **Management UIs:** Tools like **Lens** and **Portainer** provide an intuitive graphical user interface for observing, managing, and troubleshooting Kubernetes clusters, making them accessible even to team members who aren’t command-line wizards.
  - **Proven Migration Paths:** For teams considering the move, detailed, step-by-step guides from engineers who have already made the journey are widely available, providing a clear roadmap for a successful migration.

### Myth 2: The Cloud is Inherently More Reliable

The idea of the cloud as an infallible utility is a dangerous misconception. Murphy’s Law—“anything that can go wrong will go wrong”—applies equally to all complex systems, including those run by hyperscalers.

The Google Cloud incident of June 2024 is a perfect illustration. During a major outage affecting multiple services in a region, Google’s official status page offered a clear workaround: “Customers can failover to

zones in other regions” (Google Cloud Status, 2024).<sup>2</sup>

While technically correct, this advice reveals a deep disconnect from the reality of most cloud customers. Architecting, implementing, testing, and maintaining a true multi-region, active-active failover strategy is a monumental engineering challenge. It is incredibly complex and, more importantly, prohibitively expensive. It often doubles your infrastructure costs at a minimum. The vast majority of startups and even established mid-sized companies are simply not prepared for this. They operate under the assumption of regional reliability.

This incident teaches a crucial lesson: **true reliability is engineered, not purchased.** Learning to operate on bare metal forces a team to confront the realities of redundancy, networking, and failure modes head-on. This hands-on experience often leads to a more genuinely resilient system than a single-region cloud deployment that operates on a foundation of untested assumptions.

## 5 A Pragmatic Path Forward: The Hybrid Model and FinOps Maturity

This paper does not advocate for a full-scale, dogmatic exodus from the cloud. It advocates for intelligence, intention, and financial accountability.

### The Right Tool for the Right Job

A “cloud-smart” strategy involves a nuanced assessment of each workload:

- **Cloud Strengths:** The cloud remains the undisputed champion for: Early-Stage Development: When you are pre-product-market fit and need maximum speed and flexibility. Cost-Effective Testing: Using spot instances for CI/CD runners and other stateless, fault-tolerant tasks is an brilliant, low-cost use of the cloud’s scale.
- **Bare Metal Strengths:** For the core of your business—the applications that handle the majority of your traffic and are relatively predictable—bare metal offers an unbeatable combination of price-to-performance.

### Operationalizing FinOps: Beyond Simple Cost-Cutting

Mature FinOps is about continuous, proactive optimization. This can happen within the cloud, too. The Spot.io by NetApp white paper on operationalizing FinOps provides a framework for this, and their case studies show its power:

- **Finova saved 70% on their Azure spend** by intelligently leveraging spot instances for their workloads.<sup>3</sup>
- **Siemens, guided by CloudCheckr, saved 30-40% on their RDS costs** simply by identifying that a cheaper edition of the service met their technical requirements.<sup>4</sup>

These examples show that vigilance pays off. However, they are still optimizations within the cloud’s pricing model. The biggest optimization of all can be moving the workload entirely.

A critical, often overlooked aspect of cloud cost is data egress fees. Cloud providers charge a significant premium to move your data out of their network. For applications that serve large files, video, or simply have a lot of data transfer, these fees can become a massive and unpredictable line item on your bill. In

---

<sup>2</sup>Source: Incident Report - Google Cloud Service Health — <https://status.cloud.google.com/incidents/dS9ps52MUnxQfyDGPfkY>

<sup>3</sup>Source: <https://spot.io/case-studies/how-finova-saved-70-percent-azure-cloud-spend/>

<sup>4</sup>Source: <https://spot.io/case-studies/cloudcheckr-is-the-engine-driving-cloud-governance-at-siemens/>

a direct challenge to this model, providers like Cloudflare have launched their R2 storage service with zero egress fees, leading to a competitive charge that highlights just how punitive these fees can be. Most bare metal providers offer far more generous and predictable bandwidth pricing, which can lead to enormous savings.

### A Prime Example: Blockchain Validators and the Egress Cost Trap

Nowhere is this egress cost issue more critical than in the world of blockchain infrastructure. Consider the task of running a Solana validator. These nodes are the backbone of the network, responsible for processing transactions and maintaining consensus. They are characterized by extremely high transaction throughput and constant communication with other nodes.

For a Solana validator, almost every processed transaction that is gossiped to the rest of the network constitutes data egress. With the volume of transactions on the network, this results in a constant, high-volume stream of outbound data, potentially amounting to many terabytes per month.

On a major cloud provider, this business model becomes instantly unviable. The per-gigabyte egress fees would systematically erode, and likely exceed, any staking rewards earned by the validator. The operational cost would be dominated not by compute, but by data transfer. In this context, choosing a hyperscaler is not just a suboptimal choice; it is a financially ruinous one.

This is why the vast majority of serious Solana validators, and indeed validators for many other high-throughput blockchains, run on bare metal. Bare metal providers typically include a massive—often tens or hundreds of terabytes—or even completely unmetered bandwidth allowance in their predictable monthly fee. For this specific use case, the decision to use bare metal is not a performance optimization; **it is an absolute business necessity dictated by the crippling reality of egress costs.** This powerful example demonstrates that for certain applications, understanding the cloud pricing model is just as important as understanding the technology itself.

## 6 Conclusion: From Cloud-First to Cloud-Smart

The evolution from a "cloud-first" to a "cloud-smart" philosophy is a sign of a maturing industry. It marks a transition away from chasing trends and towards making deliberate, engineering-led, and financially sound decisions. It is the very essence of a robust FinOps and DevOps culture.

By understanding the raw power and unparalleled economics of bare metal for core, stable workloads, and surgically combining it with the unique, flexible services of the cloud, organizations can architect a truly superior infrastructure. This hybrid model delivers:

1. **Superior Performance:** Faster applications and a better user experience.
2. **Predictable and Lower Costs:** Breaking free from the punitive elements of cloud pricing and achieving long-term financial sustainability.
3. **Deeper Operational Knowledge:** Empowering your team with a true understanding of how their applications run.
4. **Engineered Reliability:** Building resilience based on sound principles, not just expensive service tiers.

The question for tech leaders, architects, and engineers to ask is no longer "Should we be in the cloud?" but rather, "Are we using every tool at our disposal—cloud, bare metal, and everything in between—in the smartest way possible to achieve our business goals?"

## 7 Frequently Asked Questions (FAQ)

### 1. Is the cloud bad?

**Answer:** Absolutely not. The cloud is a powerful and essential tool that has enabled incredible innovation. It excels at providing speed, agility, and scalability, especially for new companies or for workloads that are highly variable and unpredictable. The argument is not to abandon the cloud, but to use it strategically for its strengths, rather than as a costly default for all infrastructure needs. A smart strategy uses both.

### 2. Isn't moving to bare metal a huge upfront cost (CapEx)?

**Answer:** It doesn't have to be. While you can purchase your own hardware, the rise of Bare Metal as a Service (BMaaS) from providers like Equinix Metal, Packet, or OVHcloud allows you to rent dedicated physical servers on a monthly or even hourly basis. This gives you the OpEx flexibility of the cloud with the performance and cost benefits of dedicated hardware, eliminating the need for large upfront investments.

### 3. We are a small team. Is managing bare metal too complex for us today?

**Answer:** It is far more accessible than you might think. With modern tools like Kubernetes for orchestration, your team manages applications and containers, not individual physical machines. The day-to-day workflow for a DevOps engineer using `kubectl` is nearly identical whether the underlying nodes are cloud VMs or physical servers. The benefit is that your team gains invaluable, deep knowledge of the full stack—from the network to the application—which leads to better long-term engineering.

### 4. What are egress fees and why do they matter so much?

**Answer:** Egress fees are charges cloud providers apply when data leaves their network to reach the public internet. These costs are often overlooked at first but can become a major and unpredictable expense for high-bandwidth applications like video streaming, file delivery, or blockchain validators. For example, a Solana validator may push 60–100 TB of outbound data monthly. On AWS, that could cost \$4,000–\$7,000/month in egress alone, whereas bare metal providers often include generous or unmetered bandwidth for a flat fee—e.g., \$38–\$360/month. This makes bare metal not only predictable but dramatically cheaper, which is crucial when calculating total cost of ownership (TCO) at scale.

## References

- Spot by NetApp White Paper. How to operationalize FinOps to drive cost and cloud efficiency. <https://spot.io/white-paper/how-to-operationalize-finops-to-drive-cost-and-cloud-efficiency-white-paper/>
- IBM. What is a bare metal server? <https://www.ibm.com/think/topics/bare-metal-dedicated-servers>
- Rodrigo Mompo Redoli, Amjad Ullah (15 Apr 2025). Kubernetes in the Cloud vs. Bare Metal: A Comparative Study of Network Costs. arXiv:2504.11007 [cs.DC]
- Nishanth Reddy Pinnapareddy (06 May 2025). Cloud Cost Optimization and Sustainability in Kubernetes. <https://doi.org/10.52783/jisem.v10i45s.8895>
- Gideon Juve, Ewa Deelman, Karan Vahi, Gaurang Mehta, Bruce Berriman, Benjamin P. Berman, Phil Maechling. Scientific Workflow Applications on Amazon EC2. arXiv:1005.2718 [astro-ph.IM]
- 8grams. Migrate from Public Cloud: Building Kubernetes Bare-Metal Infrastructure. <https://8grams.medium.com/migrate-from-public-cloud-building-kubernetes-bare-metal-infrastructure-136b74888f34>
- Arpit Bhayani, Subhash Choudhary. How Dukaan moved out of Cloud and on to Bare Metal. <https://youtu.be/vFxQyZX84Ro?si=8ELZFwBLrYGzKOMY>
- Bare Metal. Why Bare Metal Beats the Cloud: Massive Cost Savings Without the Lock-In. <https://bare-metal.io/why-bare-metal-beats-the-cloud-massive-cost-savings-without-the-lock-in>
- Ericsson. How Kubernetes over bare metal infrastructure improves TCO. <https://www.ericsson.com/en/blog/2022/6/how-kubernetes-over-bare-metal-infrastructure-improves-tco>
- OneUptime. How moving from AWS to Bare-Metal saved us \$230,000 /yr. <https://oneuptime.com/blog/post/2023-10-30-moving-from-aws-to-bare-metal/view>
- Lyrid. Why Managed Hosting Providers are Clamoring Over Bare Metal Kubernetes. <https://www.lyrid.io/post/bare-metal-an-alternative-for-kubernetes-applications>
- Ericsson. Why Kubernetes over bare metal infrastructure is optimal for cloud native applications. <https://www.ericsson.com/en/blog/2022/5/kubernetes-over-bare-metal-cloud-infrastructure-why-its-important-and-what-you-need-to-know>
- Oleg Zinovyev. Does Kubernetes Really Perform Better on Bare Metal vs. VMs? <https://thenewstack.io/does-kubernetes-really-perform-better-on-bare-metal-vs-vm>