# HIVELOCITY

# Rethinking Bare Metal Infrastructure in 2026

The wait tax, the virtualization trap,
and the bare metal advantage.

## DEDICATED SERVERS | BARE METAL

Why dedicated bare metal servers deliver the performance,
predictability, and savings your business deserves.

# What's Inside

# Introduction: The wait tax

For nearly two decades, the economic law of the data center was simple: wait six months, and it will get cheaper. We lived in a deflationary era where Moore's Law manifested directly in pricing — computing became denser, faster, and more accessible with every passing quarter.

## Welcome to 2026. The rules have changed.

The old playbook of infinite cloud scaling is no longer a safety net.

We have entered an era defined by the Wait Tax — a market condition where deferring infrastructure decisions results in compounding cost increases due to scarcity and inflationary pressure.

It's a structural shift that demands a strategic pivot. The solution is a return to the fundamentals of performance, isolation, and control found in dedicated bare metal servers.

Here is the forensic analysis of why infrastructure costs are shifting in 2026, and how smart organizations are using dedicated servers to secure their future.

**Total Cost Over Time: Cloud vs. Bare Metal**

— Cloud (variable + egress + licensing)

- - - Bare metal (predictable TCO)

<- The Wait Tax

Month 1    Month 3    Month 6    Month 9    Month 12

Cloud costs compound month over month. Bare metal TCO stays flat and predictable.

# The silicon squeeze: The physics of scarcity

To understand why a standard server costs more in 2026 than it did in 2024, you have to look upstream to the fabrication plants. Global silicon production capacity is finite, and allocation is ruthless. The semiconductor industry has made a decisive pivot toward high-margin AI components, creating a phenomenon known as wafer displacement.

## The high bandwidth memory effect

The production of high-bandwidth memory (HBM) — essential for the GPUs driving the AI boom — is incredibly resource-intensive. Producing 1 GB of AI-ready HBM consumes roughly three times the wafer capacity of producing 1 GB of standard DDR5 server memory. Major manufacturers like Samsung and SK Hynix have reallocated production lines to chase this demand, creating a structural shortage of the commodity components that form the backbone of general-purpose hosting.

| | | |
|---|---|---|
| **3×**<br>Wafer capacity per GB of HBM vs standard DDR5 | **55–60%**<br>Server DRAM price increase forecast Q1 2026 | **10%+**<br>Price jump possible while waiting for quote approval |

## The return of the memory supercycle

The result is a hyper-bull cycle for memory pricing. Contract prices for server DRAM are forecast to rise by 55–60% in Q1 2026 alone. This hits virtualization hosts particularly hard. The transition to DDR5 is mandatory for modern high-core-count CPUs, but the high-density modules required for dense virtualization are seeing the steepest price hikes.

## The Wait Tax in action

This supply constraint creates the Wait Tax. In the past, you could procure hardware just-in-time. Today, hardware costs can appreciate week-over-week during peak scarcity. If you wait for a quote to be approved, the inventory might be gone — or the price might have jumped 10%.

# The power premium:
# When electricity becomes a luxury good

If silicon is the first constraint, power is the second. In 2026, power availability has replaced rack space as the primary limit on data center capacity. The cost of electricity has transitioned from a stable utility expense to a volatile commodity, subject to extreme inflationary pressures in major hubs.

## The PJM shock

The epicenter of this crisis is the PJM Interconnection region, which covers Northern Virginia — the world's largest data center market. In the capacity auction for the 2025/2026 delivery year, clearing prices for power capacity surged dramatically. This cost is being passed directly through to tenants in the form of higher base rack rates and aggressive power surcharges.

The era of cheap, flat-rate power in major interconnectivity hubs is effectively over.

## The density challenge

| Rack Metric | Previous Decade | 2026 |
|---|---|---|
| Standard kW/rack | 5–10 kW | 50–100 kW |
| AI workload density | Rare | Standard |
| Power pricing model | Flat rate | Variable surcharge |
| Power as TCO % | Minor | Major |

Power density has increased 10× over the previous decade, dramatically raising operating costs in traditional data center hubs.

# The licensing crisis:
# The virtualization tax

Perhaps the most disruptive force in 2026 is the decoupling of software costs from hardware value. We are witnessing a licensing cliff where aggressive monetization strategies by incumbent software vendors are destroying the economics of traditional virtualization.

## 01
**TheSilicon Squeeze**
AI chip demand crowds out standard server memory, pushing DRAM prices +55-60%

## 02
**ThePower Premium**
Power costs now rival hardware in major data center hubs. Flat-rate power is gone.

## 03
**TheVirtualization Tax**
VMware/Broadcom licensing shocks are destroying TCO. Open-source is the exit.

Three compounding forces are reshaping infrastructure economics in 2026.

## The Broadcom-VMware fallout

The acquisition of VMware by Broadcom has resulted in a pricing shock that continues to reverberate. The shift to bundled subscriptions (VMware Cloud Foundation) and the elimination of perpetual licenses has led to renewal cost increases that are reshaping procurement decisions across the industry.

> "The Virtualization Tax penalizes modern hardware. Deploy a high-density AMD EPYC Turin processor to save on hardware, and the core-based licensing model eats your savings."
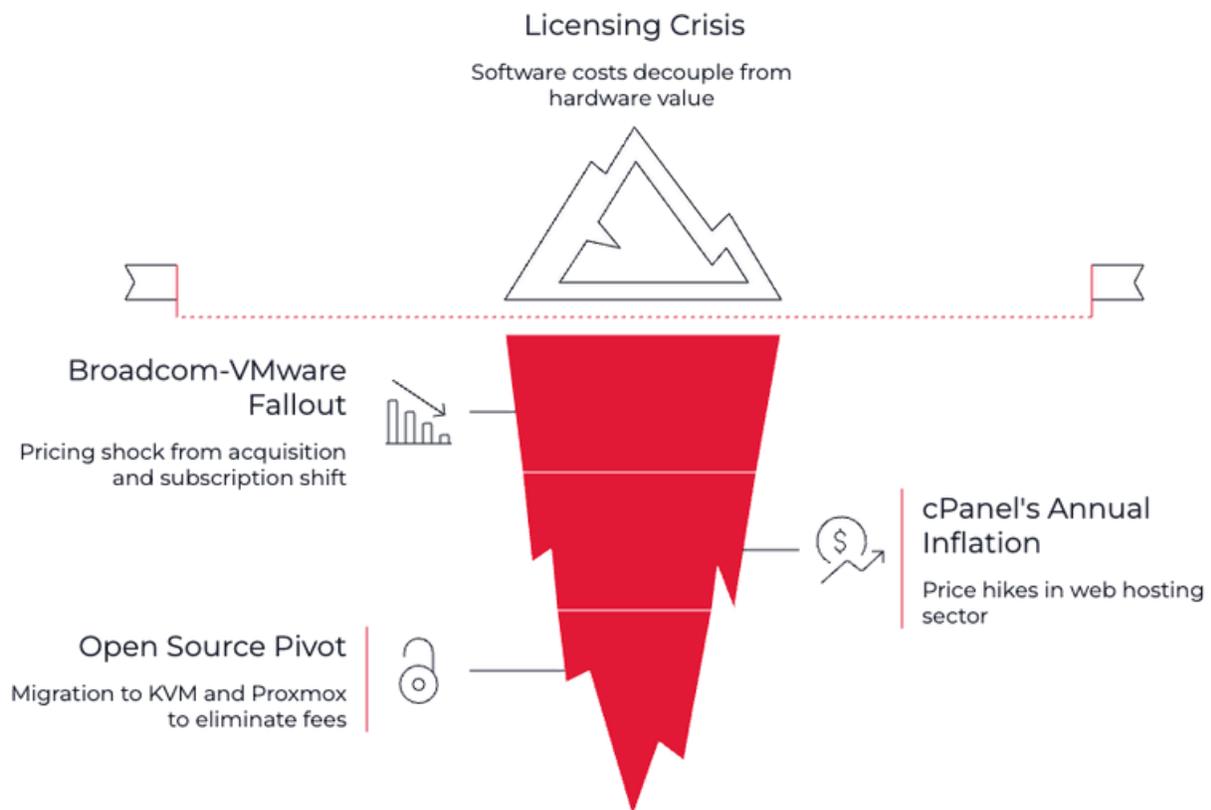
## cPanel's annual inflation

The web hosting sector faces a similar pressure. As of January 2026, cPanel implemented another round of price hikes, with the Premier tier now sitting at roughly $69.99/month plus $0.49 per additional account. For a dedicated server hosting 1,000 accounts, the software license now rivals the cost of the hardware itself.

## The open source pivot

These taxes are driving a massive migration toward KVM, Proxmox, and other open-source hypervisors. By moving to a dedicated server running a Linux-based open-source stack, businesses can eliminate the software layer's exorbitant fees entirely. It is the single most effective lever for reducing total cost of ownership in 2026.

**Software Licensing Crisis: The Virtualization Tax**

### Licensing Crisis

Software costs decouple from hardware value

### Broadcom-VMware Fallout

Pricing shock from acquisition and subscription shift

### cPanel's Annual Inflation

Price hikes in web hosting sector

### Open Source Pivot

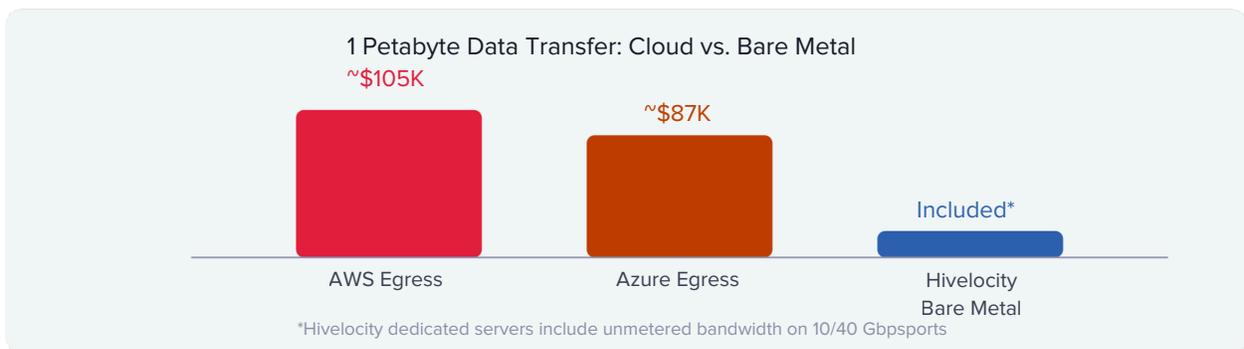Migration to KVM and Proxmox to eliminate fees

# The great repatriation: Cloud vs. bare metal

The convergence of these factors — expensive hardware, volatile power, and predatory licensing — has catalyzed The Great Repatriation. Organizations are moving steady-state workloads away from hyperscale public clouds back to dedicated infrastructure.

Why? Because the public cloud business model relies on data gravity and egress fees that have become indefensible for data-intensive applications.

## The egress trap

Consider the math of moving data out of AWS. At approximately $0.09 per GB, transferring just 1 petabyte of data out of the cloud can cost roughly $90,000 to $120,000.



1 Petabyte Data Transfer: Cloud vs. Bare Metal

~$105K — AWS Egress

~$87K — Azure Egress

Included* — Hivelocity Bare Metal

*Hivelocity dedicated servers include unmetered bandwidth on 10/40 Gbpsports

For streaming services, AI training datasets, or large-scale backups, egress fees are a massive variable risk. Hivelocity dedicated servers include unmetered bandwidth on 10 Gbps and 40 Gbps ports.

## The 37signals benchmark

"By repatriating their workloads to owned hardware in colocation, 37signals (makers of Basecamp) projected over $10 million in savings over five years."

# The strategic playbook for 2026

Navigating this landscape requires a shift in procurement strategy. Passive buying is a liability. Here is how smart infrastructure leaders are adapting.

### 01 — Leverage instant provisioning as a hedge

In a supply-constrained market, availability is a feature. Hivelocity's instant dedicated servers are pre-racked, pre-wired, and ready to deploy. This is secured inventory. It lets you bypass OEM lead times and spot-market price volatility. You get the agility of cloud — deploying in minutes — without the variable billing.

### 02 — Consolidate with high-density silicon

Use next-generation silicon like AMD EPYC Turin to consolidate workloads. A single Turin-based server with 192 cores can replace 3 to 4 legacy servers, drastically reducing your rack footprint and simplifying management overhead.

### 03 — Right-size your AI hardware

Don't use a sledgehammer to crack a nut. While training foundation models requires H100s, running them (inference) does not. The NVIDIA L40S or high-density CPUs are significantly more cost-effective for inference workloads, avoiding the massive premiums attached to training-grade GPUs.

### 04 — Hybrid architecture: base on metal, burst to cloud

Host your predictable, steady-state workload (the 70%) on dedicated servers to minimize total cost of ownership and ensure isolation. Reserve the public cloud only for the unpredictable 30% of traffic spikes — minimizing egress fees while maintaining flexibility.

# Conclusion: Return to fundamentals

The dedicated server market in 2026 is defined by a return to fundamentals. The allure of infinite cloud elasticity has faded in the face of rising costs and complexity. In its place, a pragmatic, performance-first mindset has emerged.

While hardware and power costs are creating inflationary headwinds, the strategic value of owning your resources and avoiding the hidden taxes of the hyperscale cloud has never been higher.

Hivelocity delivers the dedicated servers, managed hosting, and bare metal expertise to help you navigate 2026 with certainty.

**KEY TAKEAWAYS**

- **The Wait Tax is real.**
  Procurement delays now cost money. Inventory is constrained and prices rise week over week. Lock in your hardware now.

- **Cloud egress is a trap.**
  Moving 1 PB of data out of AWS costs $90K–$120K. Unmetered bare metal bandwidth eliminates that risk entirely.

- **Open source wins in 2026.**
  Eliminating VMware and cPanel licensing is the single most effective TCO lever available today.

- **Hybrid beats all-in.**
  Base 70% of workloads on bare metal. Burst the rest. This is the optimal cost architecture for SMBs in 2026.

## Let's find your perfect bare metal solution.

Stop paying for capacity you don't own.
Get a custom bare metal quote built around your exact workload.

[Build your server](#)