

The great unlock: How platform engineering creates AI-native enterprises





Breaking *the* AI implementation plateau

AI has seen the fastest rate of adoption of any technology in history. Almost 90% of organizations report their usage of it. The potential gains are massive, and so the hype and optimism are massive too. Most enterprises however, are not seeing these gains. They remain stuck at what we call the “AI implementation plateau”, the stage where organizations stop seeing rapid gains from AI adoption, and face slower ROI, integration hurdles, and the need for deeper cultural or process change.

They've provisioned GPUs, experimented with models, and deployed Copilot across development teams. They use AI to generate code and documentation. But, they are far from the sustainable business transformation that drive the true returns of AI implementation. Revenue optimization systems remain theoretical. Operational AI that drives profit per customer stays on whiteboards.

One key reason this plateau exists is because enterprises treat AI infrastructure as a special case requiring new organizational patterns. They spin up siloed GPU clusters managed by data science teams. They bypass governance frameworks that took years to establish for cloud-native workloads. They fragment ownership across business units, each building custom pipelines without shared standards.

At the same time, they may lack the confidence in skillset, and especially security and governance to expand their AI workflows from simple code generation prompts to full fledged Agentic workflows.

The result? Cost proliferation, compliance gaps, failed projects that never reach production, and a confusing disappointing dive into AI enablement.



This is not a set fate for enterprises however. Achieving the aspired results of AI is possible. Platform engineering is the unlock. The same platform teams that built golden paths, and composable infrastructure templates for cloud-native applications must now do the same for AI. They must abstract GPU orchestration complexity, provide pre-built application stacks, embed governance before workloads go to production and enable downstream developers to focus on business outcomes with AI rather than infrastructure tuning.

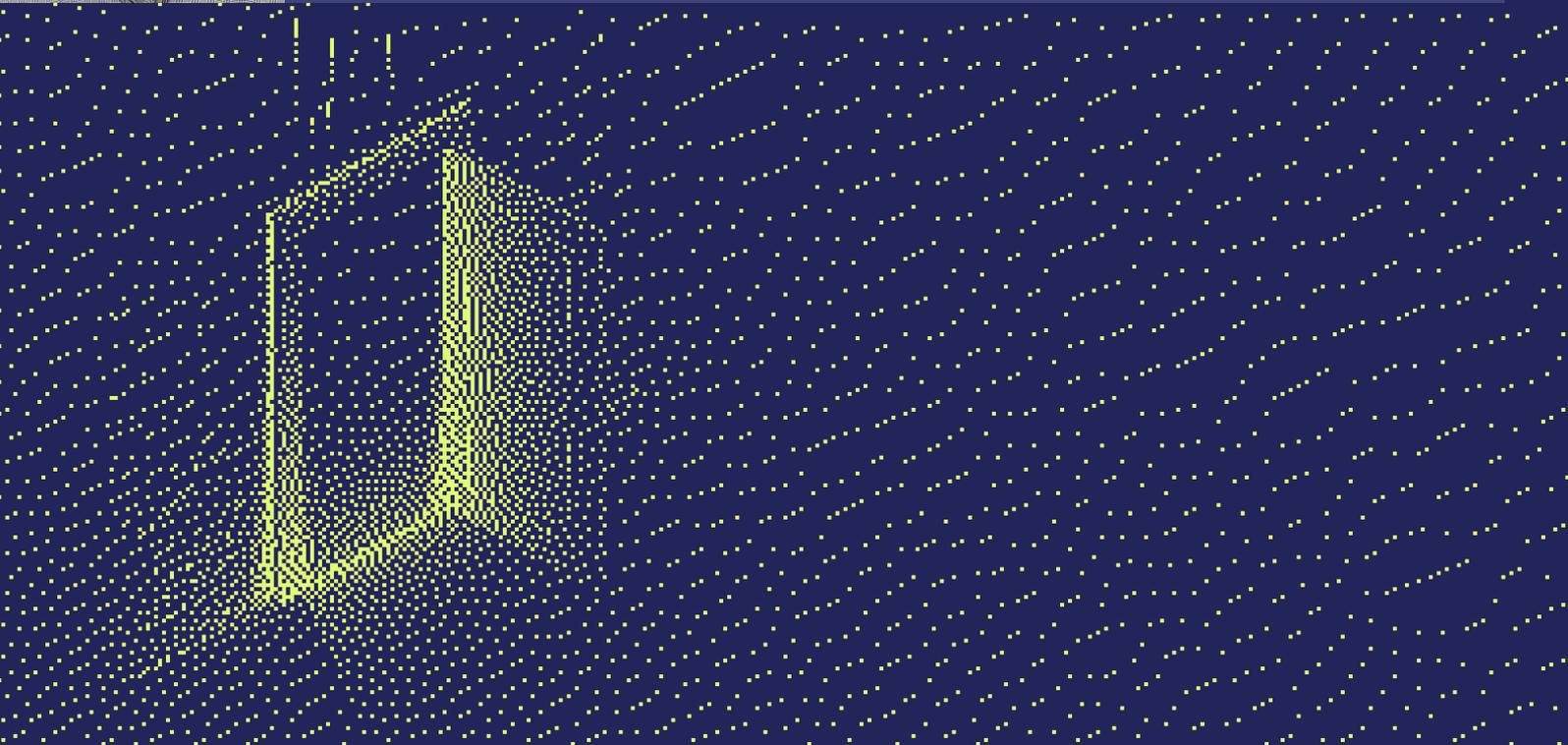
This whitepaper delivers the architectural patterns, ownership boundaries, and implementation roadmap required to make this transformation. You'll see the three-layer reference architecture that separates concerns and enables scale. You'll understand exactly how enterprises can execute on this module through a case study of how a hospitality organization can double room revenue through AI-native guest experiences deployed in weeks. You'll learn what AI-native really means, and how platform engineering enables you to attain it.



It's like a star collapse before going supernova, right? So there needs to be a collapsing of the star. Down into platform engineering as a choke point for governance and compliance. Then you can go supernova.

Kevin Cochrane

CMO, VULTR





What is an AI-native enterprise?

An AI-native enterprise embeds artificial intelligence into every business process and customer experience as a first-class operational capability, not an experimental add-on. This distinction matters. Every organization today is AI-enabled now. They all have access to AI tools, whether it's Copilot, ChatGPT or their own internal model (Copilot with lipstick). But AI-native is far more than this. It is the organizations that treat models, datasets, and inference endpoints as versioned, containerized platform services with the same rigor they apply to microservices and APIs. It's the organizations that build with AI, not that merely use AI.

The difference shows up in architecture. AI-enabled organizations merely use AI tools within existing workflows. Developers might rely on Copilot for coding assistance, teams might use ChatGPT-style interfaces, and applications may occasionally call external models through APIs. AI-native organizations take a different approach.

They design products, services, and internal workflows with AI as a core capability. Integrating models, data

pipelines, and inference services directly into their application architecture and operational systems. They version control datasets alongside code. They maintain audit trails showing which data trained which model for which application. They deploy models through CI/CD pipelines with build, test, and production environments. Simply, AI-enabled teams use AI within their existing paradigm. AI-native teams build with AI to explore new paradigms.

Platform engineering is a key enabler of this difference. Just as platform teams built golden paths for cloud-native applications, they must now provide the same for AI workloads. This means centralized model hubs containing pre-approved LLMs and custom models that downstream teams can consume. It means composable infrastructure templates that spin up GPU clusters, storage, and networking with a single click. It means pre-built application architectures defining entire RAG pipelines, data flows, and inference endpoints - all blessed, governed, and known to work.



These golden paths for AI infrastructure abstract GPU complexity while maintaining flexibility. They are pre-built, pre-tested, secure templates that ensure downstream development success.

For example, a golden path might define the entire stack for a customer service agent including GPU clusters for inference, vector databases for RAG, data pipelines for real-time updates, and API endpoints for application integration.

Developers would consume the template, customize business logic, and deploy without touching Kubernetes configurations or GPU memory management.

This puts the power of advanced AI workflows in the hands of your developers at rapid fire speed while still maintaining all necessary security, and governance protocols.



This is what allows enterprises to execute on safe, repeatable, and massively valuable AI workflows. More than just blasting through the AI implementation gap, these are the enterprises that will own the future of their industries.



Why platform engineering creates AI-native organizations

A big part of why the AI implementation plateau exists is because enterprises lack the operational playbook to industrialize AI delivery. In the State of AI in Platform Engineering report, 89% of platform engineers report using AI daily, yet most usage remains tactical, focused on code generation and documentation rather than production systems that drive measurable business outcomes. At the same time, on the infrastructure layer, organizations experiment heavily, hiring data scientists, launching pilot projects, and accessing GPU resources through cloud providers or internal infrastructure.

But experimentation alone does not translate into enterprise impact.

The gap between experimentation and business impact widens because infrastructure complexity, fragmented ownership, and governance concerns increasingly block progress.

Platform engineering solves this by applying proven cloud-native principles and product management to AI infrastructure. The same discipline that enabled enterprises to deploy thousands of microservices now enables them to deploy hundreds of AI models. The mechanism is identical. You centralize governance, abstract complexity, provide golden paths, enable distributed innovation and treat your users like they are your customers. The technology stack differs, but the organizational pattern remains constant.



Consider what happens without platform engineering. Each business unit spins up its own GPU cluster. Data science teams build custom pipelines without shared standards. Developers integrate AI through third-party APIs, sending proprietary data to external services. Security teams discover compliance gaps after deployment or worse, and out of fear simply block all experimentation. Finance sees GPU costs proliferate without clear ROI, and race to shut it all down.



Projects fail time and time again because teams spend months on infrastructure instead of business logic.

This is the mainframe anti-pattern applied to AI. Enterprises centralize GPU resources in a single location, force global teams to schedule time on shared clusters, and overprovision for peak capacity. When demand spikes in Europe at midnight, underutilized GPUs sit idle in North America. When a new model requires retraining, teams wait weeks for cluster availability. The economics don't work, the developer experience suffers, and AI remains experimental.

Platform engineering breaks this pattern by treating AI infrastructure as composable, distributed, and on-demand. GPU clusters spin up when needed and scale to zero when idle. Inference workloads deploy globally, close to users, with auto-scaling based on demand. Developers consume pre-built stacks through self-service portals without understanding GPU memory management or network optimization. Governance happens at the platform layer, before workloads reach production.



Vultr's AI-native reference architecture

The three-layer architecture separates concerns and enables scale. Each layer has a distinct purpose, clear ownership boundaries, and well-defined interfaces. This separation allows platform teams to abstract complexity while giving application teams the flexibility to innovate on business logic.

The reference architecture highlighted in this report shows a scalable, AI-native inference platform running on Vultr Kubernetes with NVIDIA GPUs utilising VAST as the data platform.

It includes a centralized master node providing enterprise IAM, Slurm workload management, observability (Grafana, Prometheus, Loki), and automation via Ansible. NVIDIA Dynamo (including Planner, SmartRouter, and Events) orchestrates and optimizes inference workloads.

Across multiple racks, NVIDIA GPU nodes (GB300 NVL72 systems with

NVIDIA Vera CPUs and Rubin GPUs) run Slurm, monitoring agents, and GPU management components. High-performance memory and storage (HBM, NVMe SSD) connect to a shared VAST data platform for globally distributed data and model access. The architecture supports elastic scale-up and scale-out inference clusters for production AI applications.

It's important to remember that the power of this platform engineering reference architecture is in its composability. Organizations can choose NVIDIA GPUs and NVIDIA's AI software stack, or other stacks and GPUs for their desired use cases. The VAST data platform can be replaced by NetApp or other AI storage solutions. The specific vendors are interchangeable. What matters is the architectural principles like composability, centralized governance, and platform-managed complexity, not the infrastructure within each layer.



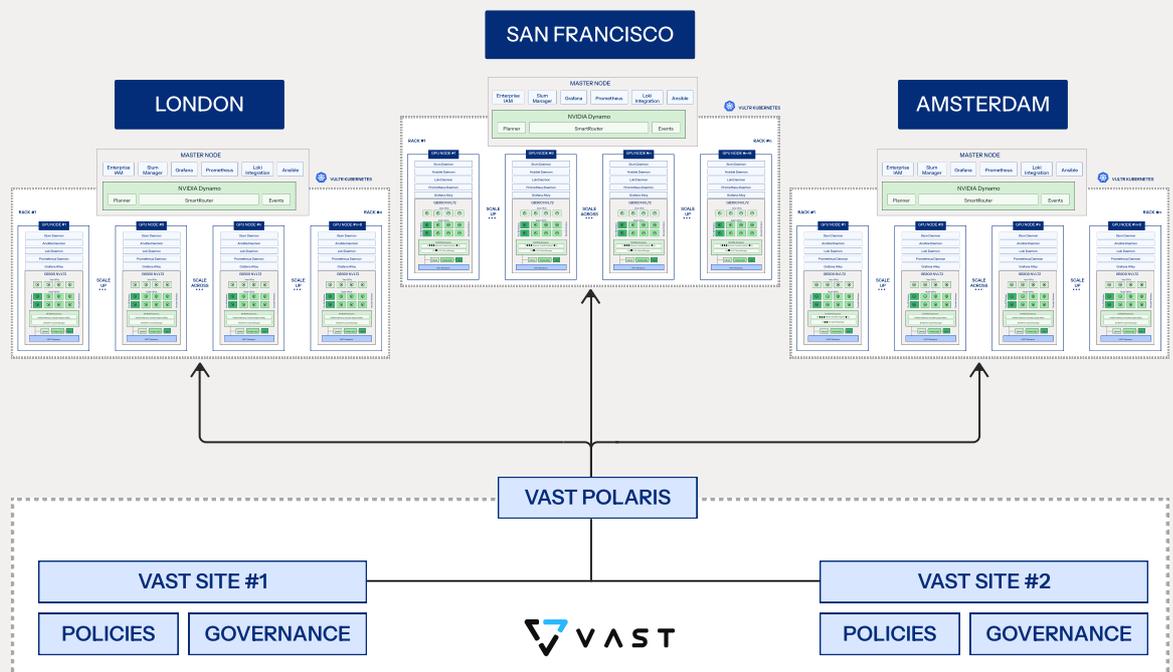
Layer 1: Governed data foundation

The data layer provides versioned, encrypted, and compliant datasets that power AI workloads. It serves as a governed foundation where every dataset is tracked, auditable, and approved for specific use cases. Platform teams, often working alongside data governance and compliance officers, certify datasets, establish access controls, and maintain data lineage.

In this example, containerized datasets flow through the architecture just like containerized microservices.

A hotel booking dataset utilised at multiple hotel locations globally might be versioned monthly, encrypted at rest and in transit, and tagged with compliance metadata (GDPR-compliant, PCI-DSS certified).

When a demand forecasting model needs training data, it pulls a specific dataset version with a full audit trail. When a RAG pipeline needs real-time updates, it subscribes to a data stream with guaranteed freshness and quality.



Without governed data, AI models lack accuracy and trust. No matter how high quality or valuable your data might be. Revenue teams won't act on pricing recommendations if they don't trust the underlying data. Compliance teams won't approve deployments if data lineage is unclear. Platform engineering solves this by treating data as a first-class platform service with the same rigor as APIs and microservices.

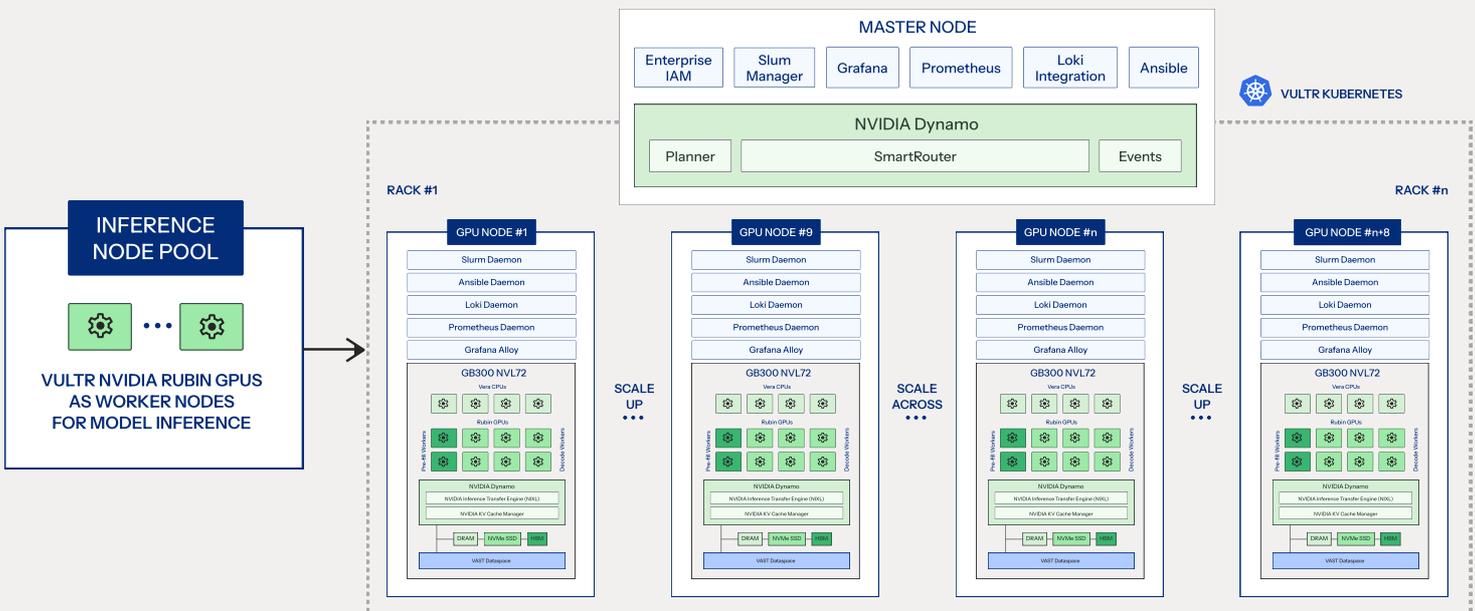


Layer 2: Accelerated compute

The compute layer provides GPU clusters optimized for two distinct workloads, scale-up training and scale-out inference. Training requires massive parallel compute on large datasets. Inference requires distributed compute close to users - think hundreds of smaller clusters globally, each handling real-time requests with low latency.

Platform teams provision and manage these clusters through the aforementioned composable infrastructure templates.

A template might define the GPU type and count, networking configuration, storage optimization, auto-scaling policies, and cost thresholds. Developers would then consume the template simply through the developer portal, customize parameters (region, capacity, budget), and deploy with a single click. The platform would abstract and handle Kubernetes configuration, GPU memory management, and cluster orchestration. This keeps the cognitive load, and the skill knowledge away from many of your developers.



Another key element of this is auto-scaling. Autoscaling to zero is critical for cost efficiency. Inference clusters scale up when demand spikes and scale down to zero when idle. A hotel chain doesn't pay for GPU capacity at 3 AM when booking traffic is minimal. The platform would monitor demand, join nodes to the cluster as needed, and decommission them when traffic drops. This is cloud-native economics applied to AI infrastructure.



Layer 3: Production orchestration

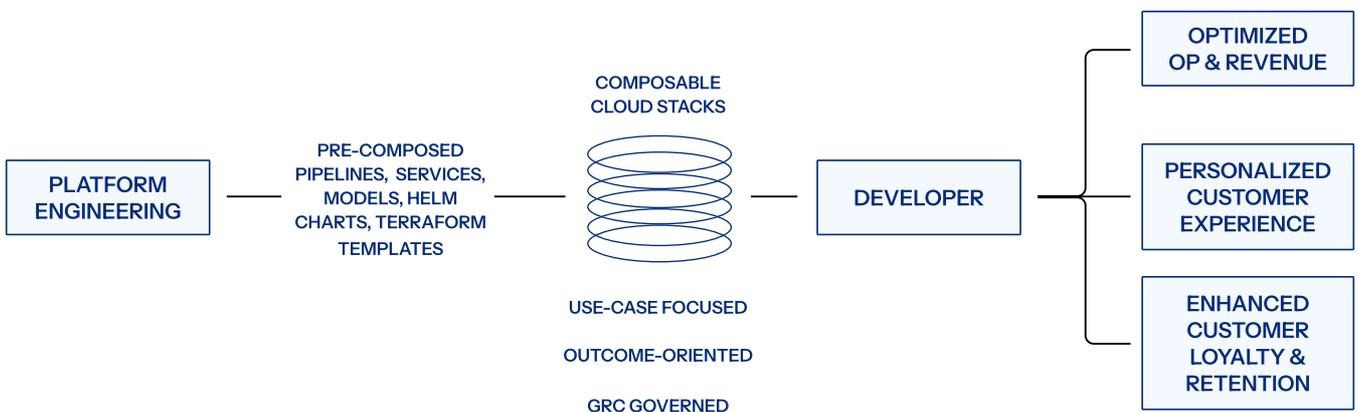
If layer one establishes trust in data and layer two provides elastic compute power, layer three is what turns AI into the repeatable production capability we are looking for.

The production orchestration layer operationalizes AI workloads the same way cloud-native platforms operationalized microservices.

Models, inference endpoints, and AI apps components are treated as versioned, deployable artifacts.

They move through controlled build, test, and production environments with the same focus as traditional software.

Platform teams do not expect AppDev teams to understand GPU scheduling, distributed inference routing, or K8s configuration. Instead, they provide golden paths for deployment encapsulating best practices for global inference deployment, observability, security policies, networking, and failover, all embedded into certified templates





A developer building a pricing optimization service or an AI-native guest experience application does not write low-level manifests or tune cluster parameters. They consume a production-ready deployment template that includes:

- ↳ Containerized model and application packaging
- ↳ Integration with centralized model hubs and registries
- ↳ Multi-region deployment patterns
- ↳ Observability and monitoring instrumentation
- ↳ Security controls and policy enforcement
- ↳ Rollback and version management capabilities



Remember that composability is the key to handling enterprise heterogeneity. A large organization will run various GPUs, NetApp and DDN storage, multiple cloud providers, and on-premises infrastructure. Platform teams can't standardize on a single vendor stack, they need to provide flexibility just as much as they need to maintain governance.

The complexity of orchestration remains within the platform boundary. The simplicity of consumption is what enables the AI driven velocity that enterprises are hungry for.



Hospitality case study: Operational AI *for* revenue optimization

As we navigate through this case study, think about how these same concepts could be executed inside your organization. A global hotel chain faces the classic hospitality challenge where occupancy has stabilized, revenue per available room (RevPAR) gains are modest, and labor costs continue rising.

Their profitability now depends on operational efficiency and margin optimization, not volume growth.

The chain needs to price smarter, staff more efficiently, and personalize guest

experiences to drive incremental revenue per room.

The business requirements are clear. They must increase RevPAR, improve gross operating profit per available room (GOPPAR), reduce labor cost per occupied room, and strengthen demand forecast accuracy. This isn't something that traditional analytics dashboards, and weekly strategy meetings can deliver. The chain needs operational AI that makes real-time decisions during live moments of opportunity for revenue.

From analytics to operational AI

The transformation starts with the three-layer architecture. Layer one unifies booking, revenue, loyalty, and operational data into a governed foundation. Guest profiles, stay history, room service orders, and booking patterns are versioned, encrypted, and made available through secure APIs. This operates as a real-time data platform where updates flow continuously to inference systems.



Layer two provides GPU clusters for demand forecasting, pricing optimization, and scenario simulation. Models run on NVIDIA-accelerated infrastructure, training on historical data and inferring on real-time booking patterns. The platform handles model deployment, version control, and API endpoint creation. Revenue teams don't manage infrastructure - they consume forecasting services through dashboards and applications.

Layer three orchestrates deployment across properties. AI-native applications such as front desk reservation system, in-room tablet experience, employee clienteling tools are able to be deployed consistently whether the property is in San Francisco, London, or Amsterdam. Kubernetes handles container orchestration, auto-scaling, and failover as the platform ensures compliance with regional data residency requirements.

This is where platform engineering becomes the enabler of outcomes. They standardize the discussed reference architecture into reusable blueprints with single-click infrastructure templates, pre-composed software stacks, approved models and endpoints, and secure

deployment patterns that plug directly into existing hotel systems.

That gives developers a paved road to build the AI applications they actually need, whether for dynamic pricing, staff scheduling, front desk upsell recommendations, in-room concierge experiences, or employee clienteling, without becoming experts in Kubernetes, GPU tuning, or distributed inference.

In practice, the platform team absorbs the complexity so application teams can focus on business logic and guest outcomes, turning the architecture into a repeatable engine for higher RevPAR, stronger GOPPAR, lower labor cost per occupied room, and more accurate real-time demand forecasting.



Beyond productivity improvements, the platform team also provides the secure guardrails that allow AI applications to access the data they need to succeed without creating security, governance, or policy risk.

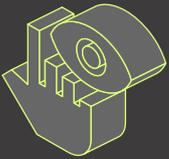
Sam Barlien

RESEARCHER, WEAVE INTELLIGENCE



The AI-native hotel guest experience

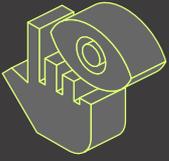
Let's break down some of the potential AI-native applications that could be utilized in this way. In this case, they are agentic apps focused on radically improving the guest experience and increasing RevPAR.



A traveler checks in late at 10 PM after a delayed flight. In a traditional hotel, the front desk (human) agent processes the reservation and hands over keys. In an AI-native hotel, the front desk application includes an embedded (AI) agent that has analyzed the guest's profile, detected the flight delay, and identified a pattern... this guest often forgets to order dinner when checking in late, and has remarked as such to staff.

The agent prompts the guest, "Mrs. Malone, I see you've arrived late. Would you like me to order your usual hamburger and onion soup for dinner? We can have it sent up in 15 minutes." The guest, focused on getting to the room, hadn't thought about dinner. But now that it's offered - yes, that would be perfect. The hotel just made \$75 in incremental revenue that wouldn't have happened without the AI-native system.

In the room, a tablet application offers three options, control lights, control TV, or talk to your agent. The agent, knowing the guest is attending a conference the next morning via the booking information, recommends car service to the venue and offers to pre-order breakfast



to arrive 30 minutes before pickup. The guest books both through the tablet. Another \$275 in incremental revenue.

Ten minutes later, room service arrives with the hamburger. The staff member, prompted by the employee clienteling application, brings a glass of the guest's preferred wine and an ice cream sundae. "I took the liberty of bringing these - I thought you might want to spoil yourself after the long travel day." The guest accepts both. Another \$50 in incremental revenue.

The room rate was \$400. The hotel just generated an additional \$400 through AI-native personalization, doubling revenue for that night. And this is just the first night of a multi-night stay. The pattern repeats, and the guest profile updates with new preferences that inform future stays in different locations across the globe.



How did they do it?

This hotel chain didn't become AI-native because it "used agents app." It became AI-native because it mobilized a composable AI stack through platform engineering allowing developers to easily and safely build the exact AI applications they needed to deliver the best results for their organization.

Behind the scenes, the infrastructure is intentionally "overwrought with GPU clusters, inference routing, storage tiers, global distribution, observability, security controls. The unlock is that downstream developers don't need to understand or rebuild any of it. The platform team packages that complexity into certified composable stacks.

That's how the applications were built quickly. One example organization built over 40 agentic workflows in just 3 weeks during a series of Hackathons. The platform and DevRel teams

pre-built the stack, then teams built on top of it: a front desk reservation app, an in-room tablet app, and an employee clienteling app. All backed by a fourth capability, a unified guest profile and memory.

While ensuring that the stack is vendor-composable with NVIDIA GPUs and software or AMD accelerators, and NetApp, DDN, or VAST for the data platform. This means that enterprises are heterogeneous by default. The same inference-backed applications also must run across properties in different regions, and the guest memory must travel with the customer.

That means globally distributed inference infrastructure and a globally distributed data layer, so a preference learned in San Francisco improves the experience in London and Amsterdam.



From single agents to agent swarms

AI-native velocity also comes from how the “agent” work is structured. These experiences aren’t powered by a single agent, they require a swarm:

Brand-to-agent

The business (the brand) defining how the AI agent should behave, what it is allowed to do, what it should optimize for, and how it represents the company.

Agent-to-agent

One agent summarizes historical behavior while another optimizes for the best offer and highest revenue likelihood.

Agent-to-employee

The system generates the staff talk track (“they arrived late, here’s what to offer, here’s how to phrase it”).

Agent-to-customer

The tablet/voice experience confirms car service, recommends add-ons, and closes the loop.



This is the great unlock in action, when platform engineering absorbs the complexity of infrastructure, governance, and agent orchestration, application teams can focus purely on business outcomes turning AI from experimentation into scalable, revenue-generating operations.



How to get started

The hospitality use case demonstrates what's possible with \$400 incremental revenue per guest stay through coordinated AI agents, deployed globally in weeks rather than months. The ability to deliver operational AI at rapid speed is the competitive differentiator.

START NOW WITH FOUR IMMEDIATE ACTIONS:

01

— **Audit your current AI initiatives**

Identify which projects are stuck in experimentation and which have clear paths to production. Prioritize the latter.

02

— **Establish platform engineering ownership for AI infrastructure**

Don't let every business unit build their own GPU clusters. Centralize governance while distributing capability.

03

— **Build your first composable template**

Start with one use case, customer service agent, or forecasting model and create a template that bundles GPU clusters, data pipelines, and AI frameworks. Version-control it. Test it. Deploy it.

04

— **Partner with CISO and compliance**

Embed security and governance into templates from day one. Don't bolt them on later.



Measure velocity and outcomes

Track time-to-production for AI applications. Track business outcomes (revenue, cost reduction, accuracy improvements). Use these metrics to justify continued investment. Can developers stand up a complete AI stack in 48 hours? Can they deploy operational AI in weeks rather than months? Track metrics like these:

→ **Time to first inference**

How long from project kickoff to first production inference request?

→ **Template adoption rate**

What percentage of AI projects use platform-provided templates?

→ **GPU utilization**

Are clusters efficiently used or sitting idle?

→ **Cost per inference**

What's the unit economics of AI workloads?

→ **Deployment frequency**

How often do teams deploy new model versions?

→ **Governance coverage**

What percentage of AI workloads have complete audit trails?



Start now: Unlock the AI revolution

The AI-native rebuild is underway. Organizations that extend platform engineering principles to AI infrastructure will compress deployment timelines from months to weeks, maintain governance and compliance, and deliver operational AI that drives the incredible business outcomes that AI promises. Those that treat AI as a special case requiring new organizational patterns will remain stuck at the implementation plateau, unable to move beyond experimentation.

Understand that platform engineering is the discipline that operationalizes AI. The same centralized teams that built developer portals and golden paths for cloud-native applications must now provide composable stacks for AI workloads. The same governance frameworks that enabled secure, compliant microservices deployment must now cover models, datasets, and inference endpoints. The same separation of concerns - platform teams abstract complexity, application teams focus on outcomes - must now apply to GPU clusters and vector databases.

The three-layer architecture shared here provides a key blueprint. Governed data foundation, accelerated compute, production orchestration. Each layer has clear ownership, well-defined interfaces, and composable components. Platform teams provide certified infrastructure and data services. Application teams build business logic and deploy through golden paths.

The gate is ready to be unlocked. We are at the dawn of a new technology cycle. Just as enterprises rebuilt infrastructure for web servers and cloud-native applications, they must now rebuild for AI-native operations. The technology differs, but the principles remain constant. Centralize governance. Abstract complexity. Provide golden paths. Enable distributed innovation.

The organizations that move first, that extend platform engineering remit to AI infrastructure now, will establish unfathomable competitive advantages that compound over time. And the organizations that wait, or do not start at all, will likely cease to exist.