

# Federator.ai Cortex™

The AI Ops Full-Stack Solution for Next-Generation GPU Data Centers

<p><b>2x</b></p> <p>GPU Efficiency</p>	<p><b>+50pp</b></p> <p>Utilization Gain</p>	<p><b>90%</b></p> <p>Downtime Reduction</p>
<p><b>19/19</b></p> <p>NCP API Coverage</p>	<p><b>PUE 1.15</b></p> <p>Cooling Efficiency</p>	<p><b>3 mo.</b></p> <p>Deployment Time</p>

- Autonomous AI Operations:** 12 AI agents perform predictive remediation, causal root-cause analysis, and self-healing—reducing unplanned downtime by up to 90%. When a CDU cooling failure occurs, Cortex traces the root cause across GPU workloads, network fabric, and cooling systems in real time—something no single-layer tool can detect.
- Every Week of Delay = \$4.3M Lost Revenue:** At 10,000 GPUs running \$18 per hour, each month of delayed operations wastes \$130M+ in potential compute revenue. Cortex compresses AI Factory deployment from 12 months to 3 months, unlocking 9 months of accelerated revenue.
- Full NVIDIA Certification (NCP) in Weeks, Not Months:** 19/19 API categories pre-built for DGX Cloud compliance. Competitors require 12–18 months to achieve what Cortex delivers out of the box—and NCP status commands a 15–20% pricing premium.
- Smart Liquid Cooling, AI-Driven:** Workload-aware thermal management with PID + feedforward control achieves  $PUE \leq 1.15$ , extends GPU lifespan by 30%, and delivers 45% higher cooling throughput than manual BMS systems.
- 16 US Patents + 15 Pending:** Multi-Layer Correlation (US Patent 11,579,933), Spatial & Temporal GPU Optimization, Predictive Self-Driving Autoscaling—a deep technology moat that cannot be replicated by open-source alternatives.
- GPU Lifespan Extension of Up to 30%:** Proactive thermal management maintains GPUs within 65–75°C optimal band, reducing thermal cycling stress and electromigration that shortens component life.

## The Cost of Inaction

GPU hardware failures cause 50% of all AI training interruptions (Meta Llama 3 study, 16,384-GPU cluster). When GPU monitoring and facility management operate in separate silos, a single CDU cooling failure can trigger a massive waste of compute resources long before the root cause in the cooling loop is ever identified.

The math is unforgiving: at \$25–\$40K per GPU-day on H100/GB300, a 512-GPU cluster spending 6 hours chasing the wrong root cause burns \$384K. Multiply this across a 10,000-GPU facility, and the annual exposure exceeds \$50M.

Industry GPU utilization averages below 50%. Static scheduling, thermal throttling, and fragmented tooling leave half of the most expensive compute hardware on earth sitting idle. Federator.ai Cortex exists to close this gap.

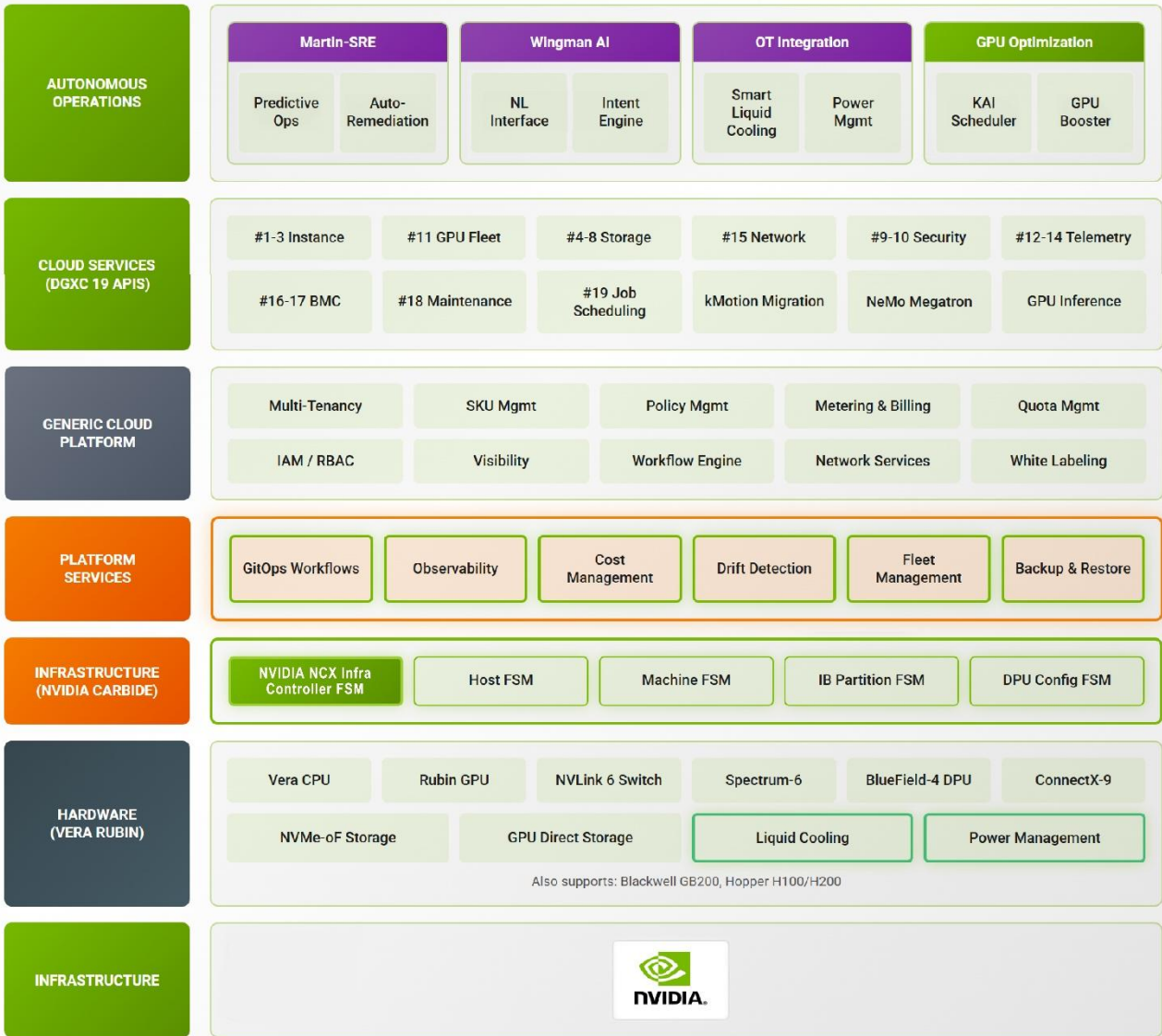
## Legacy GPU Cloud vs. Federator.ai Cortex

The table below quantifies the operational transformation that Cortex delivers across every dimension of AI factory management:

Capability	Legacy GPU Cloud	Federator.ai Cortex
<b>GPU Utilization</b>	✗ 30–50% (industry avg)	✓ 75–95% sustained (+50pp)
<b>GPU Scheduling</b>	✗ Static allocation, manual	✓ Predictive 4D scheduling (patented)
<b>Failure Detection</b>	✗ Reactive—after incident	✓ 48-hour advance prediction, 94% accuracy
<b>Training Failure Rate</b>	✗ 50% caused by GPU/HBM	✓ 50% reduction via proactive remediation
<b>Cooling Management</b>	✗ Manual BMS, thermal guesswork	✓ AI-driven PID + feedforward, PUE 1.15
<b>Cooling Throughput</b>	✗ Baseline manual	✓ +45% with workload-aware optimization
<b>Maintenance Downtime</b>	✗ Scheduled windows, SLA impact	✓ Zero downtime (kMotion live migration)
<b>NCP Certification</b>	✗ 12–18 months, custom build	✓ Weeks—all 19 APIs pre-built
<b>IT + OT Integration</b>	✗ Siloed—separate dashboards	✓ Unified cross-layer causal analysis
<b>Operator Interface</b>	✗ CLI/dashboards, manual runbooks	✓ Wingman AI natural language copilot
<b>Root Cause Analysis</b>	✗ Single-layer, manual	✓ 12-agent Bayesian DAG, 16 failure modes
<b>Financial Modeling</b>	✗ Spreadsheets, external tools	✓ Built-in ROI/IRR/Monte Carlo



End User Self Service Portal + Wingman AI Assistant



Federator.ai Cortex

AI Factory: NVIDIA Carbide FSM Infrastructure Control + Rafay K8s Operations

Federator.ai Cortex — Full-Stack AI Ops Solution: From Self-Service Portal to NVIDIA Carbide Infrastructure

## Platform Architecture

Federator.ai Cortex is the only platform providing cross-layer causal analysis and optimization across the entire AI Factory stack. Its patented Multi-Layer Correlation engine (US Patent 11,579,933) discovers real-time causal relationships between GPU workloads, network fabric, cooling systems, and power distribution. By seamlessly integrating with NVIDIA’s NCX Infra Controller and Data Center Infrastructure Management (DCIM) tools, Federator.ai Cortex accelerates "out-of-the-box" adoption. The full-stack architecture is shown in the figure above.

## DCOO — AI Factory Lifecycle Management

Cortex manages the complete AI Factory lifecycle. No other platform covers design through optimization in a single integrated system:

Phase	Key Capabilities
<b>Design</b>	Omniverse digital twin, CFD thermal simulation, ROI calculator, what-if scenarios, rack layout optimization
<b>Construct</b>	UL-certified prefab building blocks, 14-week rapid deployment, modular buildout
<b>Operate</b>	Autonomous operations, 12-agent Martin-SRE, Wingman AI copilot, 19/19 NCP API compliance
<b>Optimize</b>	Kaizen continuous improvement: Measure, Analyze, Improve, Validate, Standardize

### Quantified Business Impact

- \$130M+/month Revenue Acceleration:** Every month Cortex compresses time-to-production is \$130M+ in GPU compute revenue at 10,000 GPUs. Speed is not just a feature—it is the product.
- 60–80% OpEx Reduction in Engineering Headcount:** Every 10MW of AI infrastructure requires 20–40 specialized SRE engineers. Martin-SRE and Wingman AI eliminate the talent bottleneck—autonomous agents replace round-the-clock human teams.
- 15–20% Revenue Premium via NCP Certification:** NVIDIA NCP-certified AI factories command premium GPU-hour pricing. Cortex’s pre-built 19 APIs deliver certification in weeks, not months—unlocking premium economics from day one.
- \$40M+/year Cooling Energy Savings:** Smart Liquid Cooling v2 reduces cooling energy by 30–40% through workload-aware flow control, eliminating overcooling during idle periods. At scale (80MW facility), savings exceed \$40M annually.
- Single Platform Replaces 8–12 Point Solutions:** One system covers monitoring, scheduling, cooling, billing, compliance, incident management, capacity planning, and financial modeling. Unified event correlation eliminates blind spots from tool fragmentation.
- GPU Lifespan Extension of Up to 30%:** Proactive thermal management maintains GPUs within 65–75°C optimal band, reducing thermal cycling stress and electromigration that shortens component life.

## Deployment Specifications

Requirement	Specification
<b>Supported GPUs</b>	<ul style="list-style-type: none"> <li>NVIDIA H100, H200, GB200, GB300 (NVL72, NVL576, DGX SuperPOD)</li> <li>Vera Rubin (roadmap) — multi-gen unified management</li> </ul>
<b>Orchestration</b>	<ul style="list-style-type: none"> <li>Kubernetes v1.24+ (vanilla, OpenShift, Rancher)</li> <li>Helm Charts for automated deployment, GitOps-ready</li> </ul>
<b>Cooling Integration</b>	<ul style="list-style-type: none"> <li>Redfish v2.0+, IPMI v2.0, Modbus TCP/RTU</li> <li>MG Cooling AC250, Supermicro SCC, generic Redfish CDUs</li> </ul>
<b>Networking</b>	<ul style="list-style-type: none"> <li>NVIDIA InfiniBand 400Gb/s+, RoCE v2</li> <li>NVSwitch topology-aware scheduling, NCCL optimization</li> </ul>
<b>Telemetry</b>	<ul style="list-style-type: none"> <li>Prometheus, DCGM Exporter, VictoriaMetrics, OpenTelemetry</li> <li>Per-GPU: SM clock, HBM temp, NVLink, power draw</li> </ul>
<b>Backend</b>	<ul style="list-style-type: none"> <li>Python 3.12+, FastAPI, SQLAlchemy 2.0 async</li> <li>PostgreSQL, NATS JetStream, Redis</li> </ul>
<b>Frontend</b>	<ul style="list-style-type: none"> <li>Next.js, React, TypeScript, Tailwind CSS v4</li> <li>Real-time WebSocket, glassmorphism dark-theme UI</li> </ul>
<b>AI/ML Runtime</b>	<ul style="list-style-type: none"> <li>LangGraph multi-agent orchestration, Google Gemini</li> <li>NVIDIA NIM, Ollama, OpenAI, vLLM (pluggable LLM provider)</li> </ul>
<b>Security</b>	<ul style="list-style-type: none"> <li>JWT + API key auth, RBAC, SecretStr credential handling</li> <li>Rate limiting: 120 rpm/IP, 60 rpm/tenant, TLS encryption</li> </ul>
<b>Infra Management Software</b>	<ul style="list-style-type: none"> <li>NVIDIA NCX Infra Controller and DCIM (e.g., Netbox)</li> </ul>

## Industry Applications

- AI Cloud Service Providers:** Multi-tenant GPU-as-a-Service with per-second billing, SLA enforcement, and elastic scaling.
- Sovereign AI Programs:** National AI infrastructure with data residency, government-grade security, and national LLM training support.
- Semiconductor & EDA:** GPU-accelerated chip design, process simulation, and yield prediction with workload-aware scheduling.
- Healthcare & Life Sciences:** Drug discovery, genomics, clinical imaging with HIPAA-compliant multi-tenant environments.
- Financial Services:** Low-latency risk modeling, quant trading, fraud detection with SOC2-ready compliance.
- Research & Academia:** Foundation model pre-training, climate modeling, materials science with fair-share scheduling.