

Unlocking AI Performance and Efficiency with Federator.ai GPU Booster and Smart Liquid Cooling

Meeting the GPU Performance and Thermal Challenges of the AI Revolution



Table of Contents

Optimizing GPU Performance and Liquid Cooling Efficiency in Modern AI Data Centers1

Key AI Infrastructure Players2

Federator.ai GPU Booster: Maximizing GPU Utilization and Energy Efficiency.....3

 Updated Key Features3

Complementary Functions with Run:ai: Enhancing AI Workload Orchestration.....5

Smart Liquid Cooling: A Must-Have Solution for GPU-Dense AI Data Centers6

A Powerful Combination: Federator.ai GPU Booster + Smart Liquid Cooling8

Conclusion: The Future of AI Data Centers9

References.....10

Optimizing GPU Performance and Liquid Cooling Efficiency in Modern AI Data Centers

Controlling Compute and Energy Costs and Extending GPU Life with Federator.ai GPU Booster and Smart Liquid Cooling from ProphetStor

The rapid growth of artificial intelligence (AI) and machine learning (ML) applications, especially those leveraging large language models (LLMs), is fueling an unprecedented demand for high-performance GPU computing. Cutting-edge GPUs like the NVIDIA H100, B100, and GB200 deliver the massive computational power needed to train these models—but they also introduce substantial challenges regarding GPU costs and thermal management. Traditional air cooling often fails to meet the demands of such high-density environments.

[Federator.ai GPU Booster](#) and Smart Liquid Cooling from ProphetStor Data Services directly address these challenges by offering dynamic resource allocation and real-time thermal management. GPU Booster uses AI-driven insights to optimally distribute GPU resources and manage multi-tenant workloads. At the same time, Smart Liquid Cooling adaptively controls coolant flow and temperature to ensure peak performance and reduce cooling power consumption (which accounts for up to 40% of the AI data center power consumption) by up to 30%. This improved efficiency not only minimizes energy waste but also enhances overall GPU performance—ensuring that every compute cycle is maximized, enabling faster, more reliable AI model training.

Together, these solutions form a comprehensive optimization framework for high-density GPU clusters running compute- and thermal-intensive AI/ML workloads.

Key AI Infrastructure Players

The key players responsible for delivering AI infrastructure to support advanced model training, inference, and application development are AI Data Center Operators, AI Cloud Service Providers, and GPU Server Vendors.

1. AI Data Center Operators

| Who They Are | Key Needs |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Organizations that own and operate data centers dedicated to AI and HPC workloads. These operators lease server space and compute resources to businesses and researchers, emphasizing reliability, scalability, and energy efficiency. | <ul style="list-style-type: none">• Maximizing ROI on expensive GPU investments• Achieving high utilization and reducing operational expenses• Delivering sustainable, energy-efficient solutions to meet ESG targets and much reduced OpEx |

2. AI Cloud Service Providers

| Who They Are | Key Needs |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Providers of cloud-based platforms tailored for AI development and deployment. They offer access to GPUs and related infrastructure on a subscription or pay-as-you-go model, enabling customers to build, train, and run AI models without heavy upfront investments. | <ul style="list-style-type: none">• Dynamic resource allocation that maximizes cloud ROI• Fair multi-tenant management ensuring customer satisfaction• Streamlined operations with reduced energy consumption improved scalability and much reduced OpEx |

3. GPU Server Vendors

| Who They Are | Key Needs |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Companies that design, manufacture, and sell high-performance servers optimized for AI and HPC workloads. Their offerings are built around cutting-edge GPU technology, focusing on performance, efficiency, and scalability. | <ul style="list-style-type: none">• Differentiation through advanced, integrated solutions for much reduced OpEx• Value-added features that extend hardware lifespan with guaranteed SLA and lower total cost of ownership• Seamless integration with orchestration and cooling systems to support high-density deployment trends |

Federator.ai GPU Booster: Maximizing GPU Utilization and Energy Efficiency

[Federator.ai GPU Booster](#) is an AI-powered optimization platform designed to enhance GPU performance by dynamically adjusting resource allocation based on real-time workload demand. It not only optimizes scheduling but also leverages advanced, patented technologies that provide a multi-dimensional view of the entire AI compute stack.

Updated Key Features

- **Predictive Analytics and Dynamic Resource Allocation:**

Federator.ai GPU Booster leverages advanced machine learning—including its proprietary [CrystalClear Time Series Analysis Engine](#)—to analyze both historical and real-time data. This enables accurate forecasting of workload demands and dynamic adjustment of GPU allocations ahead of workload spikes, ensuring that every GPU is optimally provisioned.

- **Patented Multi-Layer Cascade Causal Analysis:**

At the core of Federator.ai is its patented technology that performs [multi-layer causal analysis](#). By examining correlations across application metrics, container statistics, node telemetry, and environmental conditions, the Booster makes precise resource recommendations. This holistic approach ensures that the allocation meets not only the compute needs but also aligns with power and cooling requirements.

- **Multi-Tenant and Application-Aware Optimization:**

Designed for complex, multi-tenant environments, the solution intelligently manages resource contention among diverse AI/ML workloads. By analyzing each application's behavior and performance characteristics, it assigns the appropriate Multi-Instance GPU (MIG) profiles and resource partitions to minimize idle time and drive utilization toward maximum capacity.

- **Kubernetes Integration and Automated Scaling:**

Federator.ai GPU Booster integrates seamlessly with Kubernetes, enabling effortless deployment within containerized environments. Beyond GPU allocation, it also provides recommendations for autoscaling of CPU and memory resources, reducing manual management and ensuring consistent performance across diverse workloads.

- **Cost Optimization and ESG Alignment:**

By continuously analyzing workload patterns and resource usage, the Booster delivers actionable insights that reduce over-provisioning and shift workloads to more cost-effective configurations. This not only drives significant cost savings but also supports ESG initiatives

by reducing energy consumption and carbon footprint.

- **Open API and Unified Management:**

With an open REST API and a single-pane-of-glass management console, Federator.ai simplifies integration into existing IT operations. This allows organizations to customize and extend its functionality across multi-cloud or on-premises environments with minimal effort.

Business Benefits:

Federator.ai GPU Booster has been shown to increase average GPU utilization by up to 2.5× (150%) and reduce workload runtimes by 50%, enabling faster training cycles and better ROI on GPU investments.

The table below summarizes the value these benefits offer to the key AI infrastructure players:

| Key Feature / Benefit | AI Data Center Operator (Cost savings, reduced OPEX, improved SLA) | AI Cloud Service Provider (Improved resource efficiency, reduced idle time, enhanced customer satisfaction) | GPU Server Vendor (Higher margins, competitive differentiation, increased customer value) |
|-----------------------------------------------|---------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------|
| Predictive Analytics & Dynamic Allocation | Optimizes GPU usage to lower energy consumption and reduce operational costs. | Ensures resources are provisioned when needed, enhancing SLA performance. | Superior resource utilization drives competitive product performance. |
| Multi-Layer Cascade Causal Analysis | Precisely matches resources to workload and environmental conditions, maximizing ROI. | Enables fair, real-time resource sharing among diverse workloads, unlocking revenue streams. | Differentiates product offerings with advanced resource management capabilities. |
| Multi-Tenant & Application-Aware Optimization | Increases capacity by effectively supporting more tenants on the same hardware. | Improves service delivery by minimizing idle time and balancing load dynamically. | Enhances overall system performance, providing a competitive edge. |
| Kubernetes Integration & Automated Scaling | Simplifies integration and accelerates provisioning, reducing IT overhead. | Streamlines deployment and scales AI/ML workloads efficiently. | Boosts product compatibility with orchestration tools for easier customer adoption. |

Complementary Functions with Run:ai: Enhancing AI Workload Orchestration

ProphetStor's Federator.ai is designed to complement established GPU schedulers like Run:ai. While Run:ai manages basic GPU scheduling and resource sharing, Federator.ai adds an additional intelligence layer through:

- **Patented Multi-Layer Causal Analysis:**

It analyzes data across applications, containers, nodes, and environmental conditions to ensure that both IT and cooling resources are optimally allocated.

- **Behavior-Aware Resource Allocation:**

The system adapts resource assignments based on real-time application performance and workload behavior. This fine-tuned approach ensures that every job receives precisely the resources it needs, minimizing waste and maximizing performance.

- **Enhanced IT and Environmental Synergy:**

In addition to GPU scheduling, Federator.ai factors in environmental metrics such as power consumption and cooling efficiency. This holistic view allows the solution to rebalance both compute and cooling resources as workloads shift, ensuring optimal performance without unnecessary energy use.

- **Prescriptive Analytics:**

Actionable recommendations are provided when workloads underutilize resources or when cooling resources are strained. This ensures that every GPU cycle—and every watt of cooling—is effectively used to drive peak performance.

Together, these functions allow Run:ai to manage baseline scheduling while Federator.ai fine-tunes resource allocation for both IT and environmental factors, creating a comprehensive optimization framework.

Smart Liquid Cooling: A Must-Have Solution for GPU-Dense AI Data Centers

ProphetStor’s Smart Liquid Cooling, integrated with Supermicro’s Direct-to-Chip (D2C) Cooling, delivers superior thermal optimization by dynamically adapting to AI workload intensity.

Key Features:

- **AI-Driven Cooling Management:**
Adjusts coolant flow rate and temperature in real time based on workload demands, ensuring that cooling capacity follows GPU compute intensity.
- **Dynamic Cooling Optimization:**
Continuously adapts liquid cooling parameters to prevent overheating and minimize energy waste.
- **Up to 30% Reduction in Liquid Cooling Energy Use:**
Significantly lowers data center power consumption by ensuring that only the necessary cooling is applied.
- **Scalability for High-Density GPU Deployments:**
Supports multi-GPU clusters, enabling higher density deployments without thermal issues.
- **Seamless Integration with Supermicro’s SuperCloud Composer (SCC):**
Provides real-time thermal monitoring and optimization, enhancing overall operational efficiency.

Business Benefits:

Smart Liquid Cooling addresses the limitations of traditional air cooling, offering the following advantages:

| Key Feature / Benefit | AI Data Center Operator (Lower OPEX, increased sustainability) | AI Cloud Service Provider (Enhanced performance reliability, reduced service disruptions) | GPU Server Vendor (Improved server reliability, extended GPU lifespan, competitive differentiation) |
|------------------------------|--------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------|
| AI-Driven Cooling Management | Lowers operating costs and improves energy efficiency, driving sustainability. | Reduces compute and energy costs, translating to improved margins and performance. | Enhances server reliability and extends GPU lifespan by reducing thermal throttling. |

| | | | |
|--------------------------------------------------|---------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------|
| Dynamic Cooling Optimization | Prevents downtime and reduces maintenance expenses through proactive thermal control. | Delivers consistent performance with fewer disruptions, increasing customer satisfaction. | Lowers maintenance costs and improves reliability, enhancing customer experience. |
| Up to 30% Reduction in Cooling Energy Use | Cuts electricity bills, reduces carbon footprint, and aids regulatory compliance. | Enables competitive pricing through reduced OPEX and improved profit margins. | Lowers total cost of ownership (TCO) and supports ESG initiatives for a competitive edge. |
| Scalability for High-Density Deployments | Maximizes space utilization with higher density GPU deployments. | Facilitates scaling of AI services to meet growing demand and support complex models. | Future-proofs server offerings with high-density configurations for evolving workloads. |
| Seamless Integration with Supermicro SCC | Simplifies management through real-time monitoring and proactive control. | Accelerates time to market with streamlined integration into existing systems. | Enhances compatibility and ease of use, reducing customer support costs. |

A Powerful Combination: Federator.ai GPU Booster + Smart Liquid Cooling

The synergy of Federator.ai GPU Booster and Smart Liquid Cooling delivers an unmatched AI infrastructure solution that empowers the key AI infrastructure players to lead in the competitive AI landscape.

Combined Key Features & Business Benefits:

| Combined Feature / Benefit | AI Data Center Operator | AI Cloud Service Provider | GPU Server Vendor |
|---------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| AI-Driven Thermal & Workload Optimization | <ul style="list-style-type: none">Reduced energy costsLower TCOIncreased sustainability | <ul style="list-style-type: none">Improved energy efficiencyReduced operating expensesGreater ability to meet ESG targets | <ul style="list-style-type: none">Differentiated product offeringHigher customer valuePotential for premium pricing |
| Adaptive Coolant Flow Rate & Temperature Control | <ul style="list-style-type: none">Dynamic cooling adjustmentsReduced energy wasteExtended hardware lifespan | <ul style="list-style-type: none">Optimized cooling for varying workloadsImproved resource utilizationCost savings | <ul style="list-style-type: none">Enhanced server reliabilityReduced maintenance needsIncreased customer satisfaction |
| Real-Time GPU Load Balancing & Predictive Scaling | <ul style="list-style-type: none">Efficient resource allocationMaximized GPU utilizationImproved workload throughput | <ul style="list-style-type: none">Dynamic resource allocation based on real workload demandEnhanced service scalabilityImproved customer experience | <ul style="list-style-type: none">Optimized server performanceIncreased workload capacityCompetitive advantage |
| Multi-Tenant AI Workload Optimization | <ul style="list-style-type: none">Fair and efficient resource allocation for all usersIncreased data center revenue | <ul style="list-style-type: none">Optimized allocation for diverse AI workloadsEnhanced service performanceIncreased customer satisfaction | <ul style="list-style-type: none">Improved server versatility and efficiencyExpanded customer baseIncreased sales potential |
| Seamless Integration with Supermicro SCC | <ul style="list-style-type: none">Simplified management and monitoringImproved operational efficiencyReduced IT overhead | <ul style="list-style-type: none">Streamlined integration with existing infrastructureReduced deployment complexityFaster time to market | <ul style="list-style-type: none">Enhanced product compatibility and ease of useIncreased customer adoptionReduced support costs |

Conclusion: The Future of AI Data Centers

In today's rapidly evolving AI landscape, balancing performance with efficiency is paramount. The escalating demands of GPU-based workloads, combined with rising energy requirements and cooling costs, require a dual-pronged approach that maximizes compute output while minimizing resource waste. ProphetStor's Federator.ai GPU Booster and Smart Liquid Cooling deliver a future-proof solution that not only enhances GPU utilization and reduces operational expenses but also optimizes both IT and environmental resource management.

By integrating advanced, patented technologies—such as multi-layer causal analysis and behavior-aware resource allocation—Federator.ai, a management plane optimization solution, complements solutions like Run.ai. This synergy ensures that every GPU cycle and every watt of cooling contributes directly to superior performance, accelerating AI model training and ensuring peak operational efficiency. For organizations looking to lead the AI revolution, this integrated solution is key to achieving unmatched efficiency, reliability, and sustainability.

To learn more about how Federator.ai GPU Booster and Smart Liquid Cooling can transform your high-performance GPU infrastructure for AI, visit <https://prophetstor.com> or contact sales at info@prophetstor.com.

References

1. Super Micro Computer, Inc., “Supermicro and ProphetStor Enable Better GPU Utilization,” 2024. [Online]. Available: https://www.supermicro.com/solutions/Solution-Brief_ProphetStor.pdf.
2. Super Micro Computer, Inc., “GPU SuperServer SYS-821GE-TNHR,” [Online]. Available: <https://www.supermicro.com/en/products/system/gpu/8u/sys-821ge-tnhr>.
3. Super Micro Computer, Inc., “GPU Servers For AI, Deep / Machine Learning & HPC,” [Online]. Available: <https://www.supermicro.com/en/products/gpu>.
4. NVIDIA, “NVIDIA H100 Tensor Core GPU: Extraordinary performance, scalability, and security for every data center,” [Online]. Available: <https://www.nvidia.com/en-us/data-center/h100/>.
5. Zinovyev, “Managed Kubernetes with GPU Worker Nodes for Faster AI/ML Inference,” The New Stack, 2024. [Online]. Available: <https://thenewstack.io/managed-k8s-with-gpu-worker-nodes-for-faster-ai-ml-inference/>.
6. W. X. Zhao, et al., “A Survey of Large Language Models,” arXiv, 2023. [Online]. Available: <https://arxiv.org/pdf/2303.18223>.
7. University of Hong Kong, “SLURM Job Scheduler,” [Online]. Available: <https://hpc.hku.hk/guide/slurm-guide/>.