

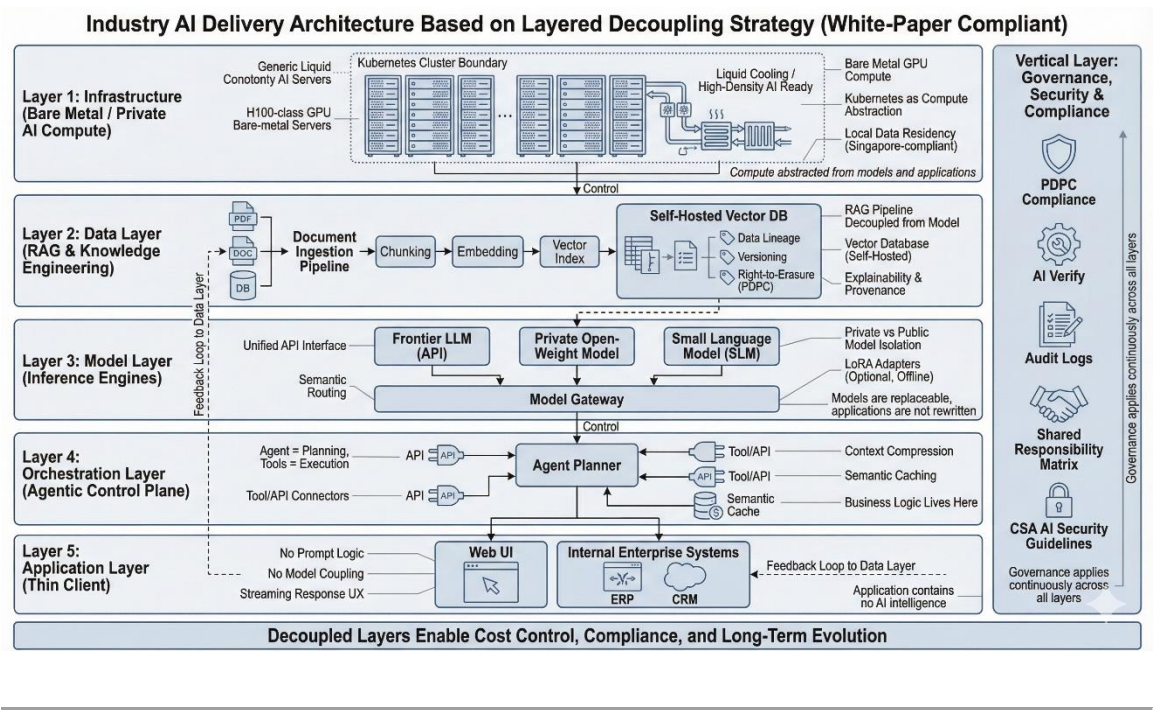
Industry White Paper: Integrated Delivery Strategy for Liquid-Cooled AI Workstations & Layered Decoupled Software

— Redefining the Physical Boundaries of Compute and Enterprise Sovereignty at the Edge

Author: dcdeeptech

Version: 3.0 (Comprehensive Edition)

Scope: Enterprise AI Infrastructure Implementation (2025-2026)



1. Strategic Background: The Convergence of Dual Critical Points

In 2025, enterprise AI infrastructure is encountering a convergence of two distinct "critical points" across physical and digital dimensions.

- The Physical Dimension: The Thermodynamic Critical Point.** With the advent of next-generation GPUs (e.g., NVIDIA Blackwell architecture), single-chip power consumption has breached the 1200W threshold. Traditional air-cooling technologies have hit a definitive physical ceiling. When rack power density exceeds 35kW, the airflow velocity required to dissipate equivalent

heat via air cooling approaches supersonic speeds; the resulting turbulence, vibration, and acoustic pressure cause physical damage to sensitive electronic components. Consequently, liquid cooling has transitioned from an "optional capability" to "**mandatory infrastructure**" for high-density compute.

- **The Digital Dimension: A Paradigm Shift in Software Delivery.** Early adoption models—characterized by direct API calls and prompt engineering—have proven fragile when facing challenges related to data sovereignty, cost containment, and regulatory compliance (e.g., Singapore's PDPC).

This white paper proposes a "**Software-Hardware Integrated**" solution: utilizing **silent liquid-cooled workstations** as the physical substrate, combined with a **Layered Decoupling** software architecture, to deliver data-center-grade compute capabilities within standard office environments.

2. The Physical Substrate: Technical Breakthroughs in Liquid-Cooled Workstations (Layer 1: Infrastructure)

The infrastructure layer is no longer confined to remote clouds; through liquid cooling technology, it "descends" directly to the enterprise edge.

2.1 Core Thermal Management: Direct-to-Chip (DTC) Liquid Cooling

We utilize **Direct-to-Chip (DTC)** technology—which commands approximately 90% of the global liquid cooling market—to resolve the thermal challenges of high-density compute:

- **Micro-Channel Heat Exchange:** Copper cold plates with micro-channel structures are mounted directly onto heat sources (CPU/GPU), utilizing a high-thermal-conductivity coolant (e.g., PG25, a water-propylene glycol mixture) to efficiently extract heat.
- **Internal Loop CDU System:** Unlike data centers reliant on facility water supplies, these workstations (e.g., the Toploong 8162-08 platform) feature independent, built-in Coolant Distribution Units (CDUs). A redundant multi-

pump design ensures balanced flow resistance across multiple GPU cold heads.


- Precision Thermal Logic:** Inlet coolant temperature is regulated between **35°C-45°C**, maintaining GPU core temperatures within the **55°C-65°C** range. This not only eliminates condensation risks in humid tropical environments but also yields a **17% increase in peak computational throughput** via sustained low-temperature operation.

Internal circulation CDU

- Built-in independent CDU system
- Self-developed data acquisition board
- Self-developed CDU management system
- Independent water distributor design
- Highly efficient and precise quick connector
- Multiple pumps for redundancy, high head, and large flow rate
- CPU and GPU dual-channel liquid cooling
- Self-developed unique microchannel design
- Self-developed unique full-water cooling design
- Approximately 35°C-45°C inlet water temperature
- Approximately 55°C-65°C GPU core temperature

Data regulation









- 7-inch touchscreen monitoring
- Displays GPU/CPU temperature and power consumption
- Displays water pump inlet and outlet water temperature, flow rate, and pressure
- Displays information about the mainboard, hard drive, memory, network, etc
- Displays noise and coolant level
- Can be used as an extended screen



A3745-01

High-efficiency noise reduction

- The machine features an ultra-quiet design, with the noise level at full load ranging from 45dB to 65dB

 Chassis size 670mm*620mm*300mm	 GPU Supports up to 8 GPU cards Compatible with 4090/5090/A100/H100/H200	 I/O expansion 11 full-height PCI/PCIE expansion slots, vertical mounting	 Hard drive backplate Front panel with 4 3.5" hard drive bays, compatible with 2.5" hard drive bays
 Compatible mainboard Standard EATX and below mainboards/Supports Supermicro/Gigabyte (15.1*13.2" 11-slot) mainboards	 Power adapter Two 3000W ATX power supplies working in parallel to solve power issues, offering high cost-performance	 AI Model Supports DeepSeek models 32B/70B	 Application areas Applications include cloud computing, AI, big data, rendering, and deep learning.

2.2 "Frictionless Deployment" in Office Environments

Liquid cooling enables high-performance compute nodes to migrate from the data center to the office (Edge):

- Acoustic Engineering:** By eliminating high-RPM industrial fans, system noise is suppressed to **45dB-65dB** (library-level silence), making it perfectly adaptable to R&D centers and office spaces.
- Physical Deployment Boundaries:**
 - Power Constraints:** The continuous safety load for a standard office 13A socket is approximately 2400W. Consequently, a **four-card configuration** represents the physical limit for non-disruptive deployment.
 - Industrial Expansion:** Configurations exceeding four cards (e.g., 8x GPU) necessitate industrial-grade electrical retrofitting.

- **Operational Assurance:** Given the complexities of fluid dynamics, deployment must be paired with localized Managed Services (MSP) that include "Electronic Equipment Liability Insurance" to mitigate user concerns regarding leakage.
-

3. The Digital Hub: Software Architecture Based on "Layered Decoupling" (Layers 2-6)

atop the liquid-cooled hardware, we architect a six-layer decoupled software stack to ensure technological agility and data sovereignty.

3.1 Data Layer: Constructing the Knowledge Moat (Layer 2)

Data is no longer egressed to public clouds; instead, it resides within **Private Vector Databases** (e.g., Milvus or Qdrant) hosted on the local liquid-cooled workstation.

- **RAG Pipeline Decoupling:** The architecture separates vectorization (Embedding) from inference. By establishing an independent **Feature Store**, enterprises can flexibly swap embedding models (e.g., transitioning from OpenAI to the open-source BGE-M3) without system refactoring.
- **Implementation of the "Right to be Forgotten":** To comply with PDPC regulations, user data can be effectively "forgotten" by physically deleting the corresponding local vector chunks. This achieves technical compliance without the prohibitive cost of model retraining.

3.2 Model Layer: The Commoditization of Inference (Layer 3)

- **Model Gateway:** A unified gateway (e.g., LiteLLM) is deployed to abstract the interface differences between underlying models (Llama 3, DeepSeek, GPT-4), exposing a standardized API to the upper layers.
- **Hybrid Routing Strategy:** The gateway dynamically routes traffic based on task complexity:
 - **High-Density Data/High-Frequency Tasks:** Routed to open-source models on the local workstation (zero data egress risk, fixed compute cost).
 - **Complex Logic/Low-Frequency Tasks:** Routed to frontier public cloud models (pay-per-token).

- **LoRA Adapters:** By loading lightweight Low-Rank Adaptation (LoRA) adapters locally (e.g., "Legal Contract Review," "Coding Assistant"), a single base model can support multiple business verticals, maximizing VRAM utilization.

3.3 Orchestration & Application Layers (Layers 4-5)

- **Agentic Orchestration:** Adopting an Agentic architecture separates task planning from tool execution. **Semantic Caching** prioritizes retrieving historical Q&A pairs locally; high hit rates can reduce external API calls by 30%-50%.
- **Thin Application Layer:** The application layer is streamlined to handle only UI interaction and streaming responses, excluding complex logic to facilitate rapid iteration and integration.

4. Economic Efficiency: Comprehensive TCO & ROI Analysis

The "Liquid-Cooled Hardware + Private Software" combination offers significant economic arbitrage, particularly in high-cost markets like Singapore.

Cost Dimension	Traditional Model (Public Cloud API + DC Colocation)	Integrated Model (Office Liquid-Cooled WS + Open Source)	Benefit Analysis
Compute Cost	High OPEX; scales linearly with business volume.	One-time CAPEX; negligible marginal cost.	ROI break-even in 3-6 months for volumes >5M tokens/day.
Facility Rent	Rack rental ~S\$3,000/month (Singapore).	S\$0 (Utilizes existing office footprint).	Liquid cooling eliminates expensive professional server room rentals.
Energy Efficiency	PUE > 1.5; high electricity overhead.	PUE 1.1 - 1.2; 40% reduction in cooling energy.	Significant long-term reduction in utility costs.

Cost Dimension	Traditional Model (Public Cloud API + DC Colocation)	Integrated Model (Office Liquid-Cooled WS + Open Source)	Benefit Analysis
Software Cost	SaaS Vector DB fees + API Token fees.	Open-Source Vector DB + Local Inference (Free).	Avoids runaway SaaS costs as data volume explodes.
Total Cost of Ownership	High (3-Year TCO).	Low (~30% of Public Cloud).	~70% reduction in TCO.

5. Vertical Governance & Compliance

- **Data Sovereignty:** Combining physical control of the workstation with localized storage ensures that core data from finance, healthcare, and legal sectors never crosses borders, fully satisfying Data Residency requirements.
- **AI Verify Integration:** The software delivery pipeline integrates Singapore's **AI Verify** toolkit. Automated testing for robustness and fairness is conducted on locally running models, generating compliance reports.
- **Supply Chain Security:** Adhering to CSA guidelines, all model weight files downloaded locally undergo strict hash verification to prevent supply chain poisoning attacks.

6. Conclusion

The "**Liquid-Cooled Layered Delivery Model**" defined in this white paper represents the ultimate morphology for enterprises facing future computational challenges.

Through **Liquid Cooling Technology**, we lock supercomputing power into silent chassis placed within arm's reach in the office; through **Layered Decoupled Software**, we reclaim control of intelligence from the cloud back to the enterprise. This is not merely a reconstruction of physical boundaries, but a comprehensive victory for enterprise economic efficiency and data sovereignty in the AI era.