

INTELLIGENCE UPDATE

AI embraces liquid cooling, but enterprise IT is slow to follow



Jacqueline Davis

13 Jan 2025

The enthusiasm for generative AI is attracting serious investment, and the associated power and cooling requirements will pose a significant challenge for the data centers that house it. Upcoming AI training clusters will escalate silicon power and rack density to unprecedented heights, upending infrastructure design conventions and accelerating the adoption of cold plate and immersion cooling systems.

AI training will act as a torture test for cooling designs and refine the engineering of cooling products that follow. It may also drive convergence and standardization. Lessons learned in the AI space may lead the way for conventional business IT to embrace densified and liquid-cooled designs as well. However, these workloads come with markedly different resiliency, networking and thermal needs — while still underpinning many businesses' primary revenue streams.

Data center infrastructure to support enterprise IT and AI training will diverge further in the coming years. AI almost stands alone in its urgency to densify its hardware, and a rising tide of liquid cooling may not lift all boats.

Liquid is a daunting change

Uptime Intelligence continually engages with data center operators to understand their experiences planning and deploying direct liquid cooling (DLC). The reasons why operators hesitate to switch to DLC remain consistent: it redefines the interface between facilities and IT teams, it can introduce unfamiliar and damaging equipment failure events, and standardization efforts will take time to develop (see [Resiliency considerations with direct liquid cooling](#)). While resilient air cooling frequently relies on cooling equipment redundancy, this approach is typically not economical or practical with current liquid-cooling designs.

In survey data and interviews, data center operators concur that the strongest incentive for liquid cooling is extreme heat output from very dense racks or powerful individual servers — with energy efficiency and sustainability objectives trailing behind. The workload dictates IT hardware selection, which, in turn, defines cooling system objectives. AI training hardware is

increasing in density much more rapidly than mainstream business IT — a trend that is likely to continue. Many enterprise operators still hesitate to consolidate workloads during their server refresh cycle.

AI is supercomputing

While AI software promises some novel outcomes, its infrastructure resembles the supercomputers that shaped today's cold plate and immersion systems. The density and workload characteristics of generative AI training clusters make them a good fit for commercial liquid cooling solutions. The task of processing billions (if not trillions) of parameters calls for high thermal design power accelerators, stacked in dense racks for low-latency communications between their many nodes. The model training workload also has more in common with high-performance computing (HPC) and supercomputing than with customer-facing, latency-sensitive business applications.

AI model training runs to completion without real-time user input, like other forms of batch processing such as engineering simulations. The task of AI training is also somewhat insulated from the revenue stream it will eventually serve. Inferences drawn from the trained model (on less powerful hardware) are where generative AI's "payoff" lies. Very large clusters, with tens of thousands of highly utilized nodes, are almost certain to see periodic hardware failures — and the training software needs to be able to accommodate this. In effect, losing a node or continuing training from a checkpoint does not cause an immediate loss of revenue.

Past forecasts of rapid densification have tended to overshoot, but product roadmaps for AI training clusters now promise 200 kW per rack and above within a few years. Training has such stringent latency requirements that close physical spacing is currently the only practical solution.

These characteristics of training clusters make the performance benefits of liquid cooling particularly attractive while also accommodating its drawbacks. Maximizing hardware performance justifies the expense of cooling equipment, and liquid coolants can handle power densities that would be difficult, expensive, or even insurmountable using air cooling alone. Since computation is performed well ahead of revenue generation, AI training can tolerate individual node failures — from cooling or otherwise — sidestepping the thorny problem of redundant coolant delivery paths. The training resiliency objectives, which inform cooling design, contrast starkly with those of mission-critical business IT.

Enterprises hesitant to take the plunge

Mainstream business IT has shaped the data center landscape thus far, placing and designing facilities to make digital services highly available and responsive to user requests. The density and thermal needs of generic IT have tended to progress more moderately over time. A large swath of critical IT applications lacks the incentive to densify — let alone reach rack densities

that demand liquid cooling. Outages often incur an immediate (and painful) loss of revenue. As a result, many organizations plan future DLC deployments with caution, evaluating the risk profile in light of DLC's resiliency differences and nascent standards.

Liquid cooling can introduce new concerns, such as reduced ride-through time on UPS battery power or cooling failures affecting multiple chassis. When cooling equipment cannot be made redundant for concurrent maintainability or fault tolerance, a software resiliency solution can compensate — though implementing this is not a trivial task. Techniques such as thin provisioning and dynamic workload allocation can maintain availability during cooling maintenance or failure by moving applications to another server, or even another data center.

While these software solutions are commonplace in cloud platforms, they remain rare in enterprise data centers. Re-architecting an application for software resiliency demands significant investment for development and testing, and risks introducing new single points of failure.

Enterprises unwilling to invest in such software resiliency might actively resist densifying their IT. Organizations may hesitate to consolidate applications on fewer, more powerful servers, fearful they will face greater disruption in the event of a single hardware failure. Data center operators that support business-critical workloads are unlikely to accept compromises to application availability to deploy a new mode of cooling, regardless of technical, economic, or sustainability benefits.

When faced with uncertainty in rack density design specifications, many organizations tend toward conservative estimates, as the financial risks associated with stranded power or cooling capacity typically outweigh those of unused space.

The enterprise IT applications best suited to aggressive densification are the most "cloud-like" — that is, dynamically allocated to compensate for liquid cooling's redundancy limitations. Otherwise, liquid cooling products would need to innovate in ways that are difficult to imagine today and close the gap to meet mission-critical IT's gold-plated resiliency standards.

Outlook

Uptime Intelligence research has previously described a key challenge facing the widespread adoption of DLC: adapting cooling designs, resiliency expectations, and operational procedures from the HPC and supercomputing world to that of mainstream business workloads. This barrier remains, and the liquid cooling revolution is likely to bypass enterprise IT for several years in favor of AI training — where liquid cooling is practical and profitable.

Generative AI training offers little insight into software resiliency for mission-critical IT, as the applications and objectives are not comparable. Further, if generative AI installations tie up much of the production capacity for liquid cooling equipment, they will likely steer the design of next-generation liquid cooling products as well. As a result, data center operators supporting conventional business-critical IT may continue to shy away from DLC if cooling manufacturers

prioritize engineering resources to improve outcomes for AI training instead.

Growth in AI training will be a boon for liquid cooling as a whole — in the sense that liquid-cooled IT will scale in the commercial data center and will no longer be confined to niche HPC, engineering simulation and cryptocurrency applications. However, even the widespread deployment of liquid-cooled AI clusters is unlikely to make liquid cooling “mainstream” as some headlines suggested; its goals and lessons learned largely do not translate to mission-critical IT applications.

The Uptime Intelligence View

AI training clusters are leading the way to extreme densification and deployment of DLC at scale, but not all can follow. The resiliency standard of mission-critical business IT is still not a fit for the liquid cooling products coming to market in the next few years, and the incentive to densify is lacking. Even widespread deployment of liquid cooling for AI training will do little to close this gap. Further, the preferential engineering focus on AI hardware may delay mainstream IT’s liquid cooling adoption even more.

Other related reports published by Uptime Institute include:

[*Resiliency considerations with direct liquid cooling*](#)

ABOUT THE AUTHOR



Jacqueline Davis

Jacqueline is a Research Analyst at Uptime Institute covering global trends and technologies that underpin critical digital infrastructure. Her background includes environmental monitoring and data interpretation in the environmental compliance and health and safety fields.

jdavis@uptimeinstitute.com

About Uptime Institute

Uptime Institute is the Global Digital Infrastructure Authority. Its Tier Standard is the IT industry's most trusted and adopted global standard for the proper design, construction, and operation of data centers – the backbone of the digital economy. For over 25 years, the company has served as the standard for data center reliability, sustainability, and efficiency, providing customers assurance that their digital infrastructure can perform at a level that is consistent with their business needs across a wide array of operating conditions.

With its data center Tier Standard & Certifications, Management & Operations reviews, broad range of related risk and performance assessments, and accredited educational curriculum completed by over 10,000 data center professionals, Uptime Institute has helped thousands of companies, in over 100 countries to optimize critical IT assets while managing costs, resources, and efficiency.