

ENERGY EFFICIENT COMPUTING: LIQUID COOLING IN THE ERA OF AI AND HPC

The demand for liquid cooling soars as companies aim to lessen the impact of increasing energy demands from AI and HPC workloads.

Produced by TCI Media Custom Publishing in conjunction with:



The Lenovo logo, consisting of the word "Lenovo" in a white, sans-serif font centered within a red rectangular background.

High energy consumption by Artificial Intelligence (AI) and High Performance Computing (HPC) workloads require innovative cooling systems to keep data center temperatures in optimal operational range. Failure to keep racks of systems sufficiently cool can lead to performance degradation and potential hardware failures. But keeping things cool is only one objective. To ensure sustainability goals are met, the heat resulting from such workloads shouldn't just be dissipated into HVAC systems to be re-cooled. Instead of diverting waste heat from the data center, the energy within that heat could be captured and repurposed in meaningful ways such as heating offices or even to make cold water through a technique called adsorption chilling.

This dual mission of balancing heat extraction from IT equipment coupled with energy reuse and being able to meet sustainability goals requires more than traditional data center air cooling techniques. Further, the coolants used to achieve high heat capture should not be toxic or have the potential to negatively impact the environment. Finally, the obvious aspect of liquid cooling in a data center means quality must be high to minimize leaks that could cause problems with electronic equipment.

Finding a liquid cooling system that performs well and meets these objectives is paramount for both business and sustainability success. There are key questions IT decision makers must ask when considering who to select to deliver their system:

- How long has the company been delivering high end liquid cooled systems?
- Do they use quality materials to prevent leakage?
- How much heat capture to liquid does a system provide?
- What type of fluids do they recommend?

The Liquid Cooling Landscape

For decades, data centers have leveraged the straight-forward method of using cold air to remove heat to keep the systems cool. However, this traditional air-cooling methodology can't keep pace with the high-power consumption and dense configurations of racks and components like central processing units (CPUs), graphics processing units (GPUs) and tensor processing units (TPUs) used for AI, machine learning, HPC and enterprise workloads. These types of workloads are growing rapidly in number and complexity, which leads to skyrocketing levels of energy consumption, heat generation and ultimately, carbon emissions.

Liquid cooling, once considered niche and a distant second in cooling options, is now in high demand as systems push the boundaries of data center power envelopes. As the energy requirements for advanced IT skyrocket, air cooling is not keeping up.

With liquid cooling, heat is dissipated by circulating a liquid coolant close to the heat-generating components, such as CPUs, GPUs or TPUs, memory, storage & PCIe. Liquids like de-ionized water or ethylene glycol have high thermal capacities that can easily absorb heat at significantly higher thresholds than air. That heat can be pumped to a heat exchanger and fed into the return loop of a data center to be harmlessly and cost-efficiently released outside the data center.

Liquid cooling had its roots in mainframe computing until it fell out of favor in the mid 1990s. It made a comeback in the mid 2010's in HPC systems running at universities conducting research. HPC systems drive compute performance to extremely high levels. Now, new workloads such as generative AI are driving Thermal Design Power (TDP) — a metric measuring the heat output that must be managed to maintain safe operating temperatures — even higher. AI workloads not only use high energy consuming components (high-end CPUs & GPUs), they use a lot of them. Keeping them cool allows them to run at their peak performance for extended periods of time.

But as with any technology, not all solutions are created equal. It is imperative to compare many factors in maintaining stable operating temperatures to ensure peak computing performance, lower hardware failures, environmental sustainability, and a reasonable total cost of ownership (TCO).

Lenovo's Approach to Liquid Cooling

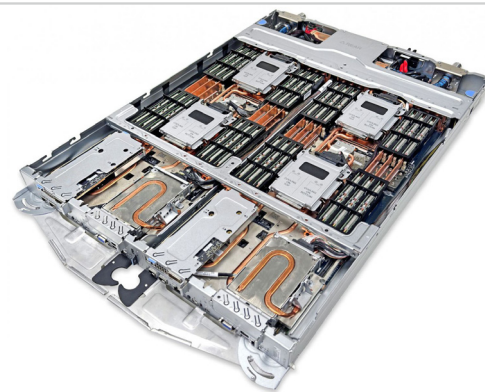
In 2012, Lenovo developed its first x86-based liquid cooled system for Leibniz Supercomputing Centre (LRZ) of the Bavarian Academy of Science and Humanities, one of the largest academic data centers in Europe. Dubbed SuperMUC, it was a revolution at the time and answered the challenge of creating a supercomputer that delivered groundbreaking power while maintaining energy efficiency. LRZ installed their follow-on system named SuperMUC Phase 2 to improve even more on their groundbreaking energy efficiency.

Lenovo broadened this liquid cooled installation and created Neptune®. Lenovo Neptune® liquid cooling technology is specifically designed to improve performance, heat removal, energy efficiency, and environmental sustainability in data centers, particularly for those running AI and HPC workloads. There are three different types of liquid cooled solutions: Neptune®, Neptune® Air and Neptune® Core.

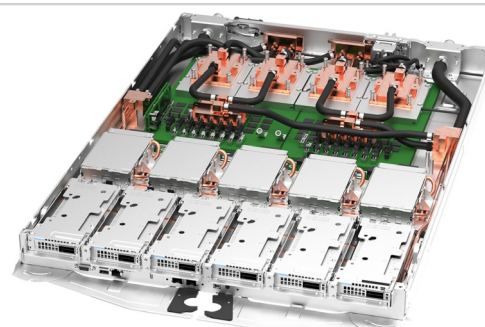


Lenovo Neptune® liquid cooling boosts AI and HPC performance

- Neptune® is Lenovo's original flagship open loop liquid cooling technology that uses copper cold plates throughout the system, removing the heat from the CPUs, memory, PCIe Storage and voltage regulators, negating the need for system fans, and achieving 100% heat capture compared to liquid.
- Neptune® Air uses liquid to augment air-cooling in a standard system. A closed loop cooling system within the server cools the CPUs, while using air to cool the remaining components. Admins won't know it's a liquid cooled system unless they open it.
- Neptune® Core uses liquid cooling in an open loop environment where liquid is brought from outside the system, across cold plates on the CPUs. Similar to the Neptune® Air solution, only the CPUs are liquid cooled, the remaining components are air cooled. The open loop then sends the discharge water back to the coolant distribution unit (CDU).



Lenovo ThinkSystem SC750 V4 Neptune®
Engineered for large-scale cloud infrastructures and HPC, this system excels in intensive simulations and complex modeling. Features Intel® Xeon® 6 processor.



Lenovo ThinkSystem SC777 V4 Neptune®
Engineered specifically for HPC, this system excels in accelerated computing for intensive simulations and Hybrid AI. Features NVIDIA GB200.



ThinkSystem N1380 Enclosure
This Neptune® chassis is the core building block for the SC750 V4 and SC777 V4 servers and represents the next generation of Neptune® liquid cooling technology. The N1380 enables exascale-level performance while maintaining a standard 19-inch rack footprint. It is a 13U enclosure and supports 8 trays mounted vertically.

Taking Sustainability to a Higher Level by Going Totally Fan-less

A main differentiator for Lenovo is its innovative approach to cooling that recovers more heat from the system. While most system vendors look to use liquids that use a maximum of 35 degrees Celsius inlet water, Lenovo pushes the envelope even higher by establishing the metric that their systems could absorb 100% heat from IT systems using inlet temperatures of 45 degrees Celsius. Higher inlet heat simply means there's no energy used to chill the water thus leading to higher energy efficiency in the data center. Additionally, the traditional approach is to cool only the components that use the most power and Lenovo does that with Neptune® Core. With the flagship Neptune® however, the company goes further and provides liquid cooling for every component that generates heat eliminating the need for internal fans. Everything that generates heat is touched by a cold plate that facilitates heat transfer to the liquid flowing through the system.

"The goal is to get rid of enough of the heat that you no longer have to have fans. Lenovo's interesting innovation is getting to the point where you're at 90% plus heat removal."

– Ian Fisk, Scientific Computing Core Director at Flatiron Institute.

"Lenovo decided very early on in the process, that they were going to do more than that. About 60% of the power is in the CPUs, a little more if you have CPUs and GPUs. But then the other components, like the memory, network cards, hard disks, all of these things have heat. And if you don't get rid of all of it, you still need airflow to get rid of the remaining heat. The goal is to get rid of enough of the heat that you no longer have to have fans. Lenovo's interesting innovation is getting to the point where you're at 90% plus heat removal," said Ian Fisk, Scientific Computing Core Director at Flatiron Institute.

Part of the innovation is Neptune® can enable a fully fan-less operation. Being fan-less is advantageous in several ways. For one, fans require power and fan power is additive to the overall power requirements of a server. Removing fans also leads to fewer moving parts inside the system and higher system availability. Finally, a fan-less server is quieter which could reduce the harmful effects that high decibels can cause on individuals' hearing.

"Moving parts like hard disk platters and fans eventually break. Things that don't have moving components have higher availability. And if there are no fans to move, there are no fans to break," says Fisk.

"The advantage of water over glycol is that you don't have to have a disposal plan... It is essentially nontoxic and if you need to get rid of it, you can flush it down the drain without having to have a disposal plan."

Becoming More Sustainable Using Water Versus Glycol PG25

Another Lenovo innovation is the use of liquid flowing through a water loop over all heat producing components to increase heat absorption. The fact that Lenovo prefers to use de-ionized water inside its liquid loops instead of glycol (commonly referred to as PG25) is also a major environmental win.

"The advantage of water over glycol is that you don't have to have a disposal plan. The water simply is water with a little bit of a biocide. It is essentially nontoxic -- de-ionized water with a little bit of anti-corrosive and anti-bacteria treatment --and if you need to get rid of it, you can flush it down the drain without having to have a disposal plan," said Fisk.

By comparison, liquid cooling with PG25 necessitates further effort and expense in disposing of it in an environmentally friendly way and in accordance with regional laws. One disposal plan will not suffice, rather disposal plans must be developed for each region where datacenters reside.

"One of the benefits we've seen from the beginning is that water cooling is very steady. Water has about 4000 times the heat capacity of air, and so it's able to pull a lot more heat away..."

However, there are other environmental advantages to using water in the liquid cooled system. Most data centers rely on evaporative air cooling which uses cold water to chill the air that is forced back into the data center. That water is evaporated, consuming thousands of gallons in the process. Only if the subsequent return line temperatures are higher than the outside ambient temperature, the data center could employ dry cooling instead of evaporative cooling. Because Neptune's water loops are enabled to run fairly warm and systems can be cooled with warm water, (think hot tub temperature water going into the system) this reduces the need to use power and energy for water chillers, and the water does not get hot enough to evaporate.

Lenovo is also focused on meeting The American Society of Heating, Refrigerating and Air-Conditioning Engineers (ASHRAE) standard TC 9.9 guidelines which outline standard practice to providing cooling to IT equipment. Not only does Lenovo adhere to the TC 9.9 standard, by leveraging de-ionized water, Neptune® can exceed the typical standard inlet temperature of 30-35°C and push inlet temperatures up to 45°C.

"The nice thing about water is that you know what it's supposed to look like, and there are standardized measures to ensure there isn't anything growing in it."

"For example, the CPU temperature might be at 80°C, and so it can be cooled by water that's quite warm at 45°C and allow it to run at full turbo. One of the benefits we've seen from the beginning is that water cooling is very steady. Water has about 4000 times the heat capacity of air, and so it's able to pull a lot more heat away, even

at high temperatures, and it allows us to operate the systems essentially in turbo mode all of the time and they remain much more stable at a high clock," said Fisk.

Increasing Sustainability and Hardware Safety by Physical Design

Another Neptune® innovation lies in the cooling system's physical design. The entire cooling loop is a single piece of brazed copper. The attached high-pressure, high-quality EPDM hoses, from vendors who supply F1 race car manufacturers with hoses, test up to 500 psi and are treated interiorly with peroxide to prevent corrosion which erodes the tubing. In short, there are no places where water can leak within these systems when water quality is maintained.

"One of the nice things about water is how much energy it can carry away. That means that the flow rates can be relative much lower than the air. ...the pumps and the pressures don't have to be incredibly high because of the cooling efficiency of the water."

"There's no O-ring to degrade either. There's simply a brazed piece of copper. The one challenge is in maintenance. Because it's a single piece of copper to remove, you must use a maintenance yoke and carefully slide it off in one piece. This also means that all of the components have to come off simultaneously. There's a little bit more work to do in maintenance, but those things are only rarely done, if at all, because CPUs rarely fail. Weigh that against the fact that every day you have a more reliable water loop," said Fisk.

Deploying liquid-cooling in a data center does mean adding another component to the data center itself: a cooling distribution unit (CDU) which circulates the water and provides heat exchanging from the Technology Cooling System (TCS) and the Facility Water System (FWS). These CDUs prevent mixing of fluid from the liquid going in and out of a server and that of the data center all the while rejecting the heat from servers to the data center cooling system. Additionally, Lenovo recommends that quarterly water quality checks be done of the TCS loop to prevent failures that can occur from biological growth, corrosion or scaling.

"The nice thing about water is that you know what it's supposed to look like, and there are standardized measures to ensure there isn't anything growing in it. Water cooling is a well understood technology. And so far, we've found the water quality remains very good," says Fisk.

Lower flow rates and pressure drops are another benefit to using Neptune®.

"One of the nice things about water is how much energy it can carry away. That means that the flow rates can be relative but much lower than the air. With very high wattage systems, like on the GPU rack, where it can be running up to 75 or 80 kilowatts on rack, the flow rate of water that needs to go through it is more than, for instance, on the CPUs, but it's also still relatively low. And so, the pumps and the pressures don't have to be incredibly high because of the cooling efficiency of the water," said Fisk.

Liquid Cooling Designed to Meet IT Scaling Challenges

Lenovo's manufacturing and delivery approaches provide additional value by adding built-in options and complementary technologies and innovations. Examples include:

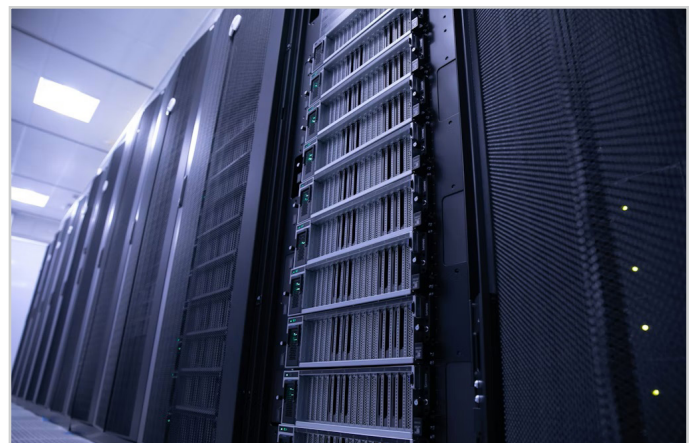
- Pre-tested equipment with de-ionized water
- Empty water loops are charged with nitrogen to prevent damage to expensive components during shipping
- Inspection and verification of cooling loop integrity prior to filling the loop with liquid at data centers installation site

The combined innovations, features and manufacturing acumen leads to industry accolades and recognition. For example, Lenovo's supercomputing pedigree is repeatedly verified by third-parties and demonstrates the company's commitment to its products.

- *HPCwire* Readers' and Editors' Choice Awards: including multiple awards for best HPC or AI product or technology
- Sustainability, Environmental, Achievement and Leadership (SEAL) Sustainable Product Award 2024
- Business Intelligence Group (BIG) Sustainability Product of the Year

Another industry list that tracks the most energy efficient supercomputers, the Green500, ranks supercomputers based on performance measured in double-precision floating point operations per watt of energy. The Flatiron Institute is a leader in energy efficiency with its Henri system holding the number one spot for the world's most energy efficient supercomputer for three consecutive lists from 2022 to 2023.

To make these benefits available to customers who don't have liquid cooling infrastructure, Lenovo also partners with leading colocation companies, like Digital Realty, where Neptune-ready infrastructure is readily available for private AI, HPC and enterprise workloads.



Flatiron Institute's Henri Supercomputer

A User's Perspective on Neptune® Liquid Cooling

The Simons Foundation is a private foundation dedicated to advancing research in mathematics and basic sciences. Its computational research division is called The Flatiron Institute. The Institute is focused on advancing scientific research through computational methods, including data analysis, theory, modeling, and simulation. It is comprised of five different centers, each focused on a specific area of computational science: astrophysics, biology, quantum physics, mathematics and neuroscience.

The Flatiron Institute [operates a significant HPC infrastructure](#) with extensive computational capability. It has an on-site cluster with 91,000 cores, 300 GPUs, and 34 petabytes of raw storage, with plans to expand by adding 20,000 more cores. Additionally, the Institute co-manages a second cluster at the San Diego Supercomputer Center, which includes 41,000 cores, 128 GPUs, and 16 petabytes of storage.

"We are now running on average between 40 and 70 kilowatts per rack, depending on the kind of equipment, which is much higher than we were able to sustain with traditional air cooling. That's allowed us to use fewer racks for our CPU systems with 72 nodes per rack."

Fisk says Flatiron has realized significant cost savings using Neptune® to cool the density of its 18 racks. He said the decision to move to Lenovo Neptune® liquid-cooling for Flatiron was based on the following factors:

- Power savings due to eliminating fans
- Replacement of fewer components with water cooled systems
- Amount of heat recovery
- Systems are designed for liquid cooling versus retrofitting existing servers for cooling
- Use of copper over plastic piping
- Ease of maintenance
- Water-cooled power supply makes the system quieter and more efficient
- Flatiron's experience with multiple generations of Neptune®

Fisk says he thinks Flatiron's San Diego installation of Neptune® may be the first installation in North America, or at least one of the earliest. Later Flatiron also did

one of the early installations of Neptune® with GPUs and two racks of 144 x 100 each using Neptune® which proved to be “workhorses for all sorts of GPU based applications.” In the Spring of 2024, Flatiron installed Neptune® with four NVIDIA H100S and 1.6 terabit of InfiniBand per node, which are intentionally designed for large language models.

“As we moved to high performance computing and machines specifically designed for AI training, the density has become very high. We are now running on average between 40 and 70 kilowatts per rack, depending on the kind of equipment, which is much higher than we were able to sustain with traditional air cooling. That’s allowed us to use fewer racks for our CPU systems with 72 nodes per rack. Each of those nodes is dual socket with a very high wattage CPU, and that would be impossible to do in air,” Fisk said.

New and expanding AI and HPC workloads will continue to add pressure to computing environments where both energy requirements and heat will soar at unprecedented rates. As a result, energy efficient computing and sustainability will be ongoing focus areas for businesses for the foreseeable future. Traditional air-cooling technologies can’t keep pace with the heat production and energy demands of the modern data center. IT decision makers must partner with technology providers who have expertise in deploying innovative cooling solutions. Lenovo Neptune® cooling technology is one of the first and among the few on the market that can do so.

For more information visit lenovo.com/Neptune.

Lenovo’s award-winning HPC solutions are helping researchers across the globe push the boundaries of what is possible and tackle some of the world’s greatest challenges. Lenovo is a global technology powerhouse, ranked at 217 in the Fortune Global 500, and a \$62 billion revenue company that employs 77,000 people worldwide focused on delivering “smarter technology for all”. Building on their success as the world’s number one PC maker, the company is expanding its research into growth areas to advance “New IT” technologies (client, edge, cloud, network, and intelligence).