# The industry's first 100% fanless direct liquid cooling system architecture

Future-focused cooling for optimizing compute-intensive AI and HPC workloads

**HPE GreenLake**

## Liquid cooling matters

- From 2007 to 2024[1]:
  - Transistors per CPU and GPU have risen by 500x
  - Single server power consumption has risen 33x
- As much as 94% cooling cost savings (over air cooling)[2]
- 84% reduction in carbon dioxide and energy when direct liquid cooling (DLC) is used combined with optimization of data center cooling infrastructure.[2]
- Only Hewlett Packard Enterprise has 100% fanless DLC that powers the top 2 supercomputers[3] on the TOP500 and 7 of the top 10 in the Green500[4]

## Introduction

As businesses drive growth from ever-growing data, more powerful IT systems are being deployed to create value more quickly. While these systems enable more complex workloads, the heat generated is outpacing traditional air-cooled infrastructures and needs an innovative 100% fanless direct liquid cooling system architecture.

In 2007, a typical server in a data center consumed around 330 watts of power. While inefficient, air-cooling could remove sufficient heat for the servers of this era, which had only 200 million transistors.[5]
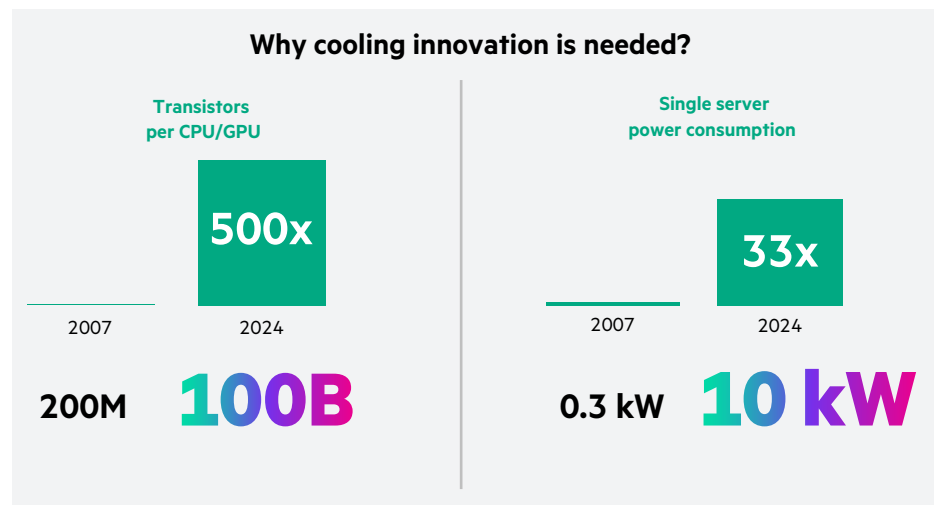


**Why cooling innovation is needed?**

Transistors per CPU/GPU

500x

2007 — 2024

200M 100B

Single server power consumption

33x

2007 — 2024

0.3 kW 10 kW

**Figure 1.** Comparison of transistor and power changes from 2007 to 2024

Figure 1 shows the massive increases in transistor density and the associated power required between 2007 and 2024. The 200M transistors in 2007 required 0.3 kW, while the 100 billion transistors packed into chips in 2024 can consume up to 10 kW of power.

Now, CPUs and GPUs contain many more transistors to process data and typically require over 330 watts individually. When you use two or more each in a system, the heat created has gone up tremendously.

It is now common for a data center to run 10,000 air-cooled servers, each consuming 10 kW each, to require 100 MW in total. On top of power and cooling for the processing of these servers, air cooling itself could require an additional 80 MW[6] of power which could cost $56M and produce some 268K tons of $CO_2$ each year![7] **That's just not sustainable**.

---

[1, 5, 6] live-lbl-eta-publications.pantheonsite.io/sites/default/files/lbnl-1005775_v2.pdf

[2] Modeling based upon deployed HPE systems; data compiled in April 2023. These data also reflect optimizations to the data center cooling infrastructure and compute performance as part of a holistic cooling upgrade.

[3] top500.org/lists/top500/2024/06/

[4] top500.org/lists/green500/2024/06/

[7] epa.gov/energy/greenhouse-gases-equivalencies-calculator-calculations-and-references

# So, what are the alternatives?

Liquids remove heat much more efficiently than air. Using similar volumes, liquids can remove more than 3000 times more heat than air, based on volume.[8] Heat can be removed to air, liquid can cool the air, or liquid can be pumped through heat plates on system components.

## DLC delivered reliably at scale

| **Air cooling** | **Liquid-to-air cooling** | **70% DLC** | **100% fanless DLC** |
|---|---|---|---|
| Fans, air conditioning, and vents circulate air and remove heat from computing equipment | Chilled water supply from the facility cools down the air-cooling system positioned close to the servers | Combined direct liquid cooling and air cooling | Coolant flows through a network of tubes and cold plates to extract heat directly from all components on the server |

**Cooling efficiency and capacity (kW/rack) increases from left to right**
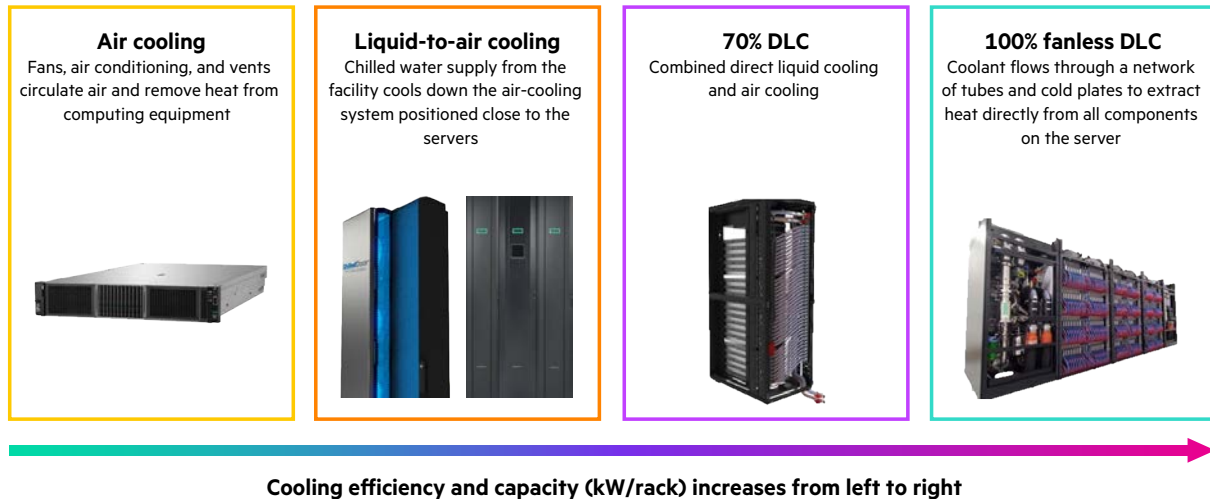
**Figure 2.** HPE expertise in liquid cooling

Figure 2 shows four types of system cooling, including air cooling, liquid-to-air cooling, and two types of direct liquid cooling (70% DLC and fanless DLC).

1. **Air cooling:** Heat sinks and server-level fans producing flowing air are used to extract heat from server components and systems, exhausting heat into the data center. Heat sink efficiencies can also be increased by using small server-level pumps to move coolants between a heat sink on the CPU or GPU to a radiator located in the server (Figure 3). The heat in the radiator is then exhausted to the air in the data center.

   **Figure 3.** Closed-loop cooling example with server-level pumps and radiator

2. **Liquid-to-air cooling:** Liquid coolants absorb heat from the air. Air-to-liquid cooling can occur either at the inlet or exhaust. Rear door heat exchangers (RDHX) neutralize exhaust air with cooled liquid at the back of the server rack. Cooled air can also be used to cool server inlet air in a self-contained rack system such as adaptive rack containment (ARC). In both cases, the heated air is cooled by a heat exchanger, with the heated fluid from the heat exchanger then cooled by facility water systems.
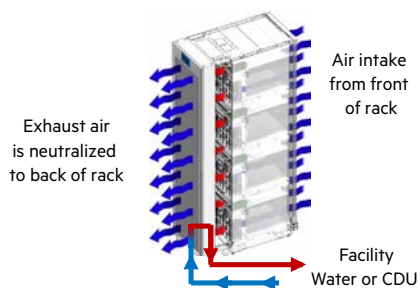
Exhaust air is neutralized to back of rack

Air intake from front of rack

Facility Water or CDU

**Figure 4.** RDHX airflow

This drawing of a rear door heating exchanger shows that the hot air exhausted by the server is neutralized by cold air through a cooling coil in the back door of the server rack.
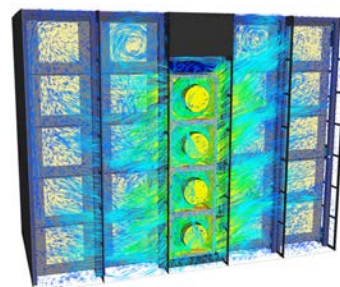
**Figure 5.** Liquid-to-air cooling air in an enclosed rack—HPE ARCS

This drawing of HPE ARCS shows how cool air is dispersed by fans throughout the enclosed rack system. Heat generated is transferred back into facility water (not shown)

[8] Using 25% propylene glycol mixture, compared to air at ambient conditions by volume; specific heats retrieved from Engineeringtoolbox.com, July 2024.

3. **70% DLC:** 70% DLC uses cold plates on GPUs and CPUs within compute systems. Fans, typically running at much lower speeds than air-cooled systems, extract the remaining heat. It is called 70% DLC, as liquid cooling typically extracts about 70% of the heat through the system via cold plates on the GPUs and CPUs. These hybrid methods cool some HPE ProLiant and HPE Cray XD systems.



Approximately **70%**
server heat goes to water (DLC)*

Approximately **30%**
server heat goes to air (fans)*

**Secondary side**
**(CDU to compute)**
**Coolant is PGW25%**

**Primary and secondary side fluids do not mix.**
**Heat is transferred through a plate-to-plate**
**heat exchanger in the CDU.**

**Primary side**
**(Facility water to CDU)**
**Coolant is water or similar**
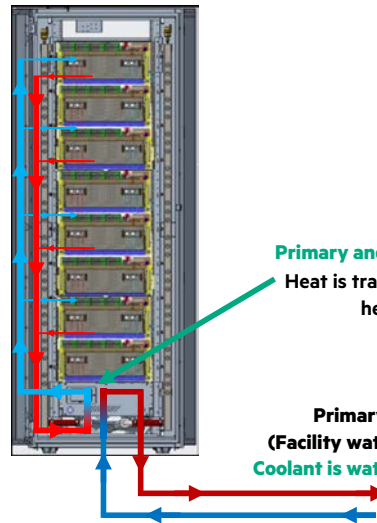
* These are typical values, not absolute

**Figure 6.** Combined liquid and air-cooling components with 70% DLC

Figure 6 shows the primary and secondary cooling pathways for a 70% DLC configuration. A cooling fluid extracts about 70% of the heat from the systems while the remaining is air-cooled. The cooling fluid from the servers is routed through a heat exchanger to discharge its heat through facility water.

4. **100% DLC (fanless):** HPE pioneered 100% fanless DLC for HPE Cray EX systems, the fastest supercomputers in the world.[9] It's also the same cooling technology that has helped us achieve 7 of the top 10 on the Green500 list.[10] This method pumps cooling fluid through cold plates through the GPUs, CPUs, memory, and rectifiers helping eliminate the need for fans. Switches and interconnects in these systems are also water-cooled. Fanless DLC can reduce the cooling portion of a data center's carbon footprint and utility cost by up to 90%. Why does it work so well? When you compare equal volumes of fluid and gaseous air, fluid has more than 3000 times the cooling capacity when compared to air.[11]

If you are still not convinced of the power of liquid cooling. Let's take an everyday example. When you burn your finger, do you blow on it? Or put it in the refrigerator? No. You run it under cool water because water transfers the heat more effectively. It works similarly in the data center.

With the increasing heat loads generated coupled with the heat-absorbing efficiency of liquid, liquid cooling will be essential for AI and compute-intensive workloads. In fact, many GPUs and CPUs require liquid cooling to achieve specified performance.

Remember that data center we talked about earlier, running 10,000 air-cooled servers, each consuming 10 KW and totaling 100 MW for the data center? That example is exactly why fanless DLC matters. A fanless direct liquid-cooled system based on HPE Cray EX technology could reduce[12] the cost of cooling power needed in this example from $56M per year[13] to $2.1M, and drop the annual $CO_2$ produced from 268K tons to around 10K tons.
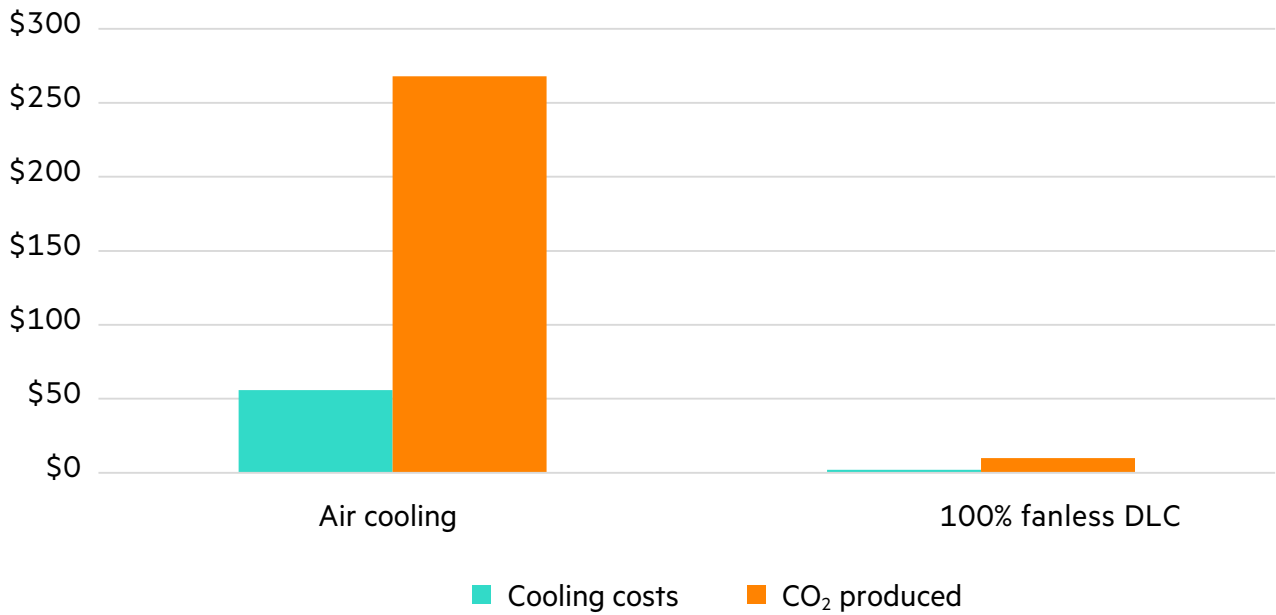
[9], [10] top500.org/lists/green500/2024/06/, retrieved October 8, 2024

[11] ctmmagnetics.com/general/air-cooled-vs-liquid-cooled-differences/#:~:text=It%20is%20striking%20to%20see,3%2C500%20times%20that%20of%20air

[12] hpcwire.com/2022/06/08/at-isc-the-green500-witnesses-a-new-frontier-in-efficient-computing/

[13] live-lbl-eta-publications.pantheonsite.io/sites/default/files/lbnl-1005775_v2.pdf

## Annual incremental cooling spending and emissions*



* Modeling based upon deployed HPE systems; data compiled in April 2023. These data also reflect optimizations to the data center cooling infrastructure and compute performance as part of a holistic cooling upgrade

**Figure 7.** Comparison of air cooling and 100% fanless DLC spending and emissions for 10,000 servers at a total of 100 MW.



**Secondary side**
**(CDU to compute)**
**Coolant is PGW25%**

**Primary and secondary side fluids do not mix.**
**Heat is transferred through a plate-to-plate heat exchanger in the CDU**

**Primary side**
**(Facility water to CDU)**
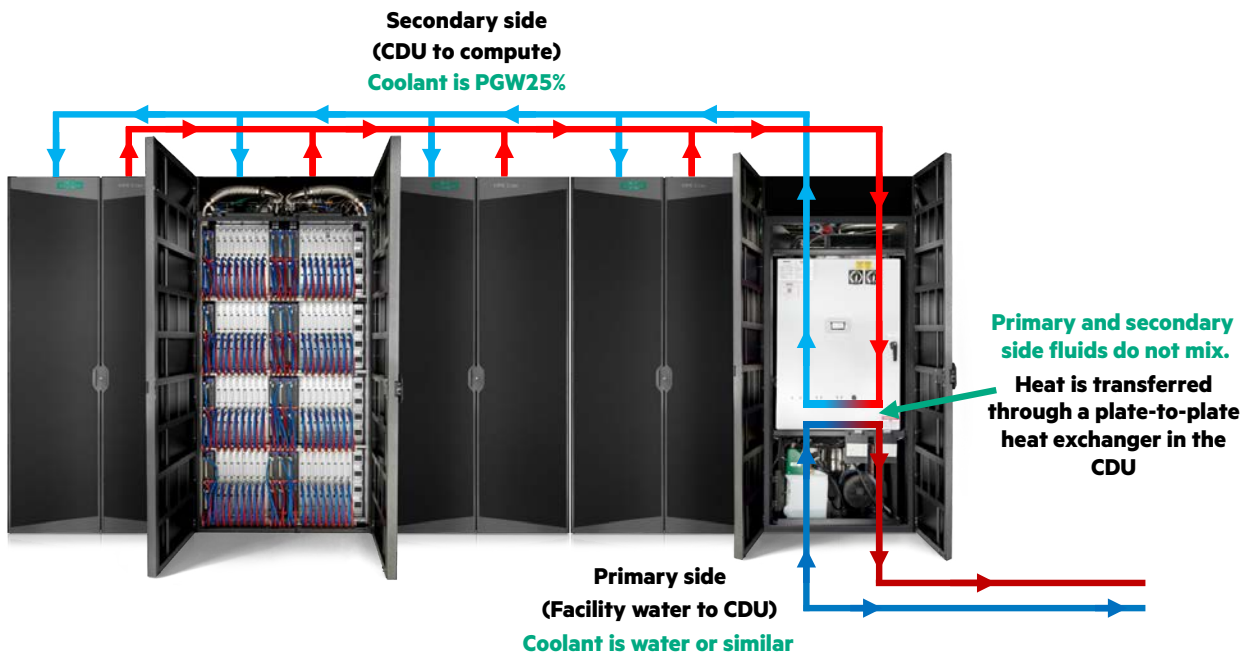**Coolant is water or similar**

**Figure 8.** Liquid cooling pathway and components with DLC

Figure 8 shows 100% DLC configurations, including a rack of servers, the coolant flow, the cooling distribution unit (CDU), and the primary cooling supply from facility water.

DLC and liquid-to-air cooling rely upon two stages for cooling. The first side is the primary, which uses facility to transfer heat from a heat exchanger. Heat from compute systems is transferred to the facility water through a plate-to-plate heat exchanger inside a CDU. The facility water and the server coolants do not mix, as the heat transfer takes place across a metal plate.

The HPE patented 100% fanless DLC innovative system design expertise enables us to create combinations of cooling technologies for the most compute-intensive applications.

**We pioneered this technology in the Frontier system, the world's first system to break exascale.[14] Since then, we have refined it and announced the HPE 100% fanless DLC system architecture built on four pillars.**

1.  8-element cooling design—includes liquid cooling for the full server blade, the network fabric, local storage, GPU, CPU, rack/cabinet, pod/cluster, and CDU

2.  High-density and high performance system design, complete with rigorous testing, monitoring software, and on-site services to support successful deployment of these sophisticated compute and cooling systems

3.  Integrated network fabric design based on dragonfly topology and connected directly with copper, because copper is not only less expensive, it doesn't consume additional power

4.  Open system design to offer flexibility of choice in accelerators

**The 100% fanless DLC system architecture delivers unique benefits to our sovereign AI and large enterprise customers.**

•   37% reduction in cooling power required per server blade, when compared to hybrid DLC alone

•   Lower utility costs and lower carbon production

•   Reduction of up to 100% data center fan noise

•   Reduced floor space from twice the server cabinet density of our competitors

•   50%–100% reduced connection power consumption through copper cables rather than depending on optical cables that require powered transceivers

We also have a dedicated team of AI experts in services and R&D with hundreds of DLC patents. While many vendors only use hybrid; only HPE is the leading provider in the world of 100% fanless direct liquid-cooled systems, and generative AI will drive this need further.

[14] Top500 List

# The right liquid cooling for the right compute workloads

While there are many options for types of liquid cooling, most organizations make compute decisions on the servers, processor, and memory requirements. To accommodate the cooling needs in your data center, HPE compute systems/servers are compatible with multiple cooling systems.

| Server | Description | Primary use | Configuration | Liquid cooling type |
|---|---|---|---|---|
| **HPE ProLiant DL325 Gen11** | The HPE ProLiant DL325 Gen11 Server is a 1U 1P solution that delivers exceptional value balancing compute, memory, and network bandwidth at 1P economics. Powered by 4th and 5th Generation AMD EPYC™ processors with up to 160 cores, increased memory bandwidth (up to 3 TB, 6000 MT/s), high-speed PCIe Gen5 I/O, and EDSFF storage, and supporting up to two DW GPUs at the front, this server is a superb low-cost, 1U 1P, performance solution for your virtualized workloads. The HPE ProLiant DL325 Gen11 Server is an excellent choice for virtualized workloads such as software-defined compute, CDN, VDI, and secure edge apps that require balancing processor, memory, and network bandwidth. | AI inference at the edge; data acquisition and preprocessing. | 1U, 1 socket, up to 2 single-wide or double-wide GPU | Liquid-to-air cooling |
| **HPE ProLiant DL360 Gen11** | The HPE ProLiant DL360 Gen11 Server is a rack-optimized, 1U 2P solution that delivers exceptional compute performance, upgraded high-speed data transfer rate, and memory depth at 2P compute capability. Powered by 4th and 5th Gen Intel® Xeon® Scalable processors with up to 64 cores, 8 TB of memory, and 20 EDSFF drives as well as increased memory bandwidth and high-speed PCIe Gen5 I/O, the HPE ProLiant DL360 Gen11 Server is a perfect solution for electronic design automation (EDA), CAD, and VDI. | Advanced natural language processing (NLP), real-time analytics, distributed machine learning (ML) | 1U, 2 socket, up to 3 GPU | DLC, liquid-to-air cooling |
| **HPE ProLiant DL365 Gen11** | The HPE ProLiant DL365 Gen11 Server is a rack-optimized 1U 2P dense solution that delivers exceptional compute performance, upgraded high-speed data transfer rate and memory depth at 2P compute capability. Powered by 4th and 5th Generation AMD EPYC processors that supports up to 160 cores and 400W per CPU, up to 6 TB of DDR5 memory, and 20 EDSFF drives as well as increase memory bandwidth and high-speed PCIe Gen5 I/O, the HPE ProLiant DL365 Gen11 Server is a perfect solution for workloads such as VDI, EDA, and CAD, and general-purpose virtualization workloads. | Advanced NLP, anomaly detection, real-time analytics, distributed ML | 1U, 2 socket, up to 2 single or double-wide GPU | DLC |
| **HPE ProLiant DL380 Gen11** | The HPE ProLiant DL380 Gen11 Server is a scalable 2U 2P solution that delivers exceptional compute performance, memory density with scalability and high-speed data transfer rate to run your most demanding applications. Powered by 4th and 5th Gen Intel Xeon Scalable processors with up to 64 cores, 8 TB of memory, and 36 EDSFF drives as well as increased memory bandwidth and high-speed PCIe Gen5 I/O, the HPE ProLiant DL380 Gen11 Server is a perfect solution for software-defined storage, video transcoding, and virtualized apps. | Deep learning (training and inference), NLP, AI-powered data warehousing, real-time data processing | 2U, 2 socket, 8 single-wide or up to 3 double-wide GPU | DLC |
| **HPE ProLiant DL385 Gen11** | The HPE ProLiant DL385 Gen11 Server is a 2U 2P solution that delivers exceptional compute performance, upgraded high-speed data transfer rate and memory depth at 2P compute capability. Powered by 4th and 5th Generation AMD EPYC 9004 and 9005 series processors with up to 160 cores, increased memory bandwidth and capacity, high-speed PCIe Gen5 I/O, enhanced GPU support, and EDSFF storage, the HPE ProLiant DL385 Gen11 Server is a superb accelerator-optimized 2U 2P solution. HPE ProLiant DL385 Gen11 Server is an excellent choice for compute and data storage demanding workloads requiring increased core count, and storage and I/O scalability. | Big Data analytics, ML, and deep learning | 2U, 2 socket, up 8 single-wide or 4 double-wide 350W GPU | DLC |

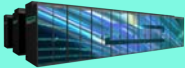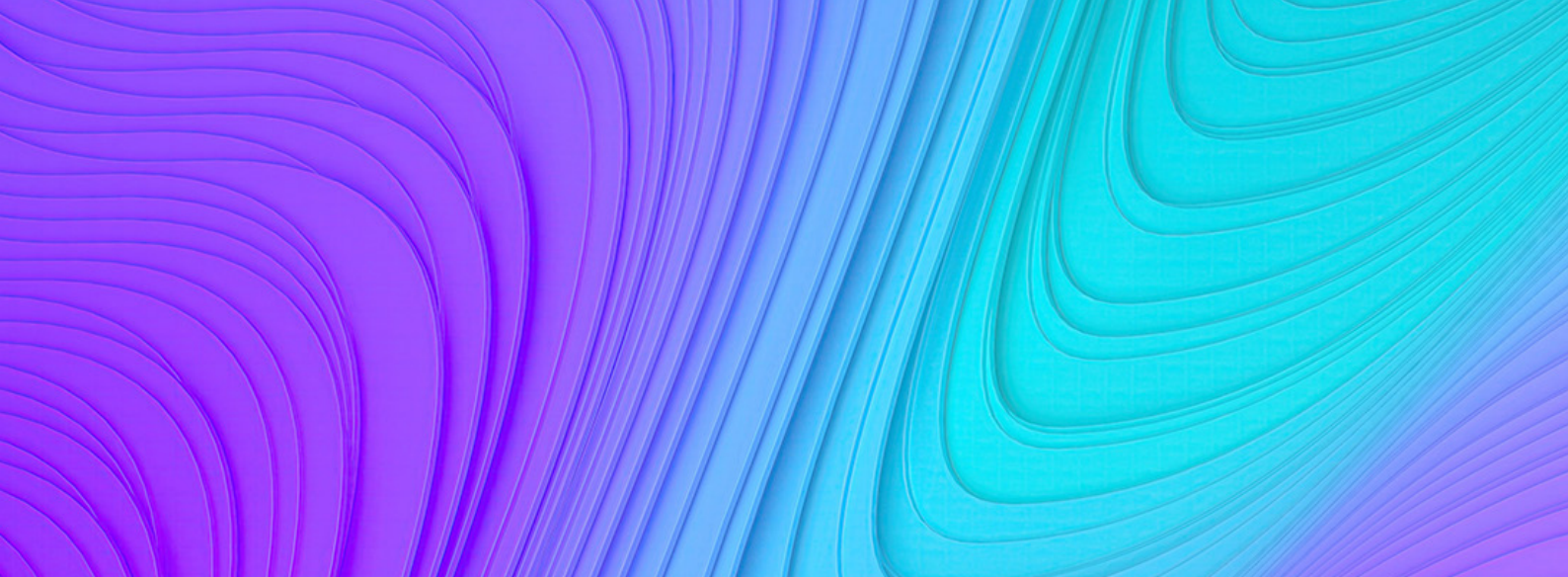| Server | Description | Primary use | Configuration | Liquid cooling type |
|---|---|---|---|---|
| **HPE ProLiant DL560 Gen11** | The HPE ProLiant DL560 Gen11 Server is a high-density, four-socket (4S) server with high performance, scalability, and reliability, all in a 2U chassis. Powered by 4th Gen Intel Xeon Scalable processors with up to 60 cores, the HPE ProLiant DL560 Gen11 Server offers greater processing power, up to 16 TB of faster DDR5 memory, I/O of up to six PCIe 5.0 slots, and up to 2 OCP slots. The HPE ProLiant DL560 Gen11 Server is an excellent choice for business critical, virtualization, server consolidation, compute, business processing, and in-memory database and data analytics workloads requiring maximum core count, memory capacity, and network and I/O bandwidth. | High performance deep learning, large-scale ML, BI and data warehousing, real-time analytics | 2U, 4 socket, up to x GPU | Liquid-to-air cooling |
| **HPE Cray XD2000** | HPE Cray XD2000 System is a dense, multiserver platform that packs incredible performance and workload flexibility into a small data center space while delivering the efficiencies of a shared infrastructure. | High performance parallel model training, Big Data analytics for finance, science, and market research | 2U chassis with choice of 1U or 2U server nodes | Liquid-to-air cooling, DLC |
| **HPE Cray XD665** | HPE Cray XD665 provides the perfect balance of GPU and CPU compute power for blended sequential and parallel HPC and AI workloads across verticals such as fintech, healthcare, research, and manufacturing. The system features four NVIDIA® H100 GPUs and two AMD 4th Generation EPYC processor platform CPUs. | Deep learning with large-scale model training, scientific research and simulation, AI-powered Big Data processing, ML model training | 4U chassis, 4 GPUs | Liquid-to-air cooling, DLC |
| **HPE Cray XD670** | HPE Cray XD670 is specifically designed and optimized for AI workloads that are heavily parallelized, requiring GPU acceleration for optimum performance. The system features eight NVIDIA H200 or H100 Tensor Core SXM5 GPUs and the latest advances in hardware and software to deliver a complete, scalable solution, purpose-built for AI. | Heavily parallelized AI workloads that require GPU acceleration | 5U chassis, 8 GPUs | Liquid-to-air cooling, closed-loop liquid cooling, DLC |
| **HPE ProLiant Compute XD685** | HPE ProLiant Compute XD685 is purposefully optimized for NLP, large language models (LLMs), and multimodal training. Powered by eight of the latest GPUs, with a DLC option and managed by HPE iLO, this system is designed to deliver unparalleled AI training performance sustainably and securely. | AI model training, NLP, LLM training, multimodal training | 5U chassis, 8 GPUs | Liquid-to-air cooling, DLC |
| **HPE Cray Supercomputing EX2500** | HPE Cray Supercomputing EX is for customers who desired to have the exascale technology of the world's largest supercomputers,[15] but in a more enterprise data center-friendly packaging. While the HPE Cray SC EX4000 supports up to 64 compute blades in a supercomputing cabinet, the HPE Cray SC EX2500 supports up to 24 compute blades in a rack. | Scientific research/simulation, large-scale deep learning training, predictive analytics | Supports the following compute blades<br>**HPE Cray Supercomputing EX4252**<br>• Four 2-socket CPU nodes<br>• Support for the 4th Gen AMD EPYC 9004 series processor stack<br>**HPE Cray Supercomputing EX425**<br>• Four 2-socket CPU nodes<br>• Support for the full 2nd Gen AMD EPYC 7002 or 3rd Gen AMD EPYC 7003 series processor stack<br>**HPE Cray Supercomputing EX420**<br>• Four 2-socket CPU nodes<br>• Support for the 4th Gen Intel Xeon Scalable processor stack (XCC and HBM)<br>**HPE Cray Supercomputing EX255a**<br>• Two 4-socket AMD Instinct™ MI300a Accelerator APU nodes<br>**HPE Cray Supercomputing EX254n**<br>• Two 4-socket NVIDIA GH200 Grace Hopper Superchip nodes | 100% fanless DLC |

[15] Top500 List

| Server | Description | Primary use | Configuration | Liquid cooling type |
|--------|-------------|-------------|---------------|---------------------|
| **HPE Cray Supercomputing EX2500 (continued)** | | | **HPE Cray Supercomputing EX235n**<br>• Two 4-socket NVIDIA A100 GPUsx1-socket CPU nodes<br>• Support for the full AMD 3rd Gen AMD EPYC 7003 series processor stack<br>**HPE Cray Supercomputing EX235a**<br>• Two 4x AMD Instinct M250X Acceleratorsx1-socket CPU nodes<br>• Support for the 3rd Gen AMD EPYC processor | |
| **HPE Cray Supercomputing EX4000** | The HPE Cray SC EX4000 solution has established itself as the world's leading supercomputing solution.[16] Today more than half of the aggregate compute power of the world's Top 100 verified supercomputers is delivered by it. | Large-scale deep learning and complex neural network training, AI-enhanced simulations, scientific discovery | Supports the following compute blades<br>**HPE Cray SC EX4252**<br>• Four 2-socket CPU nodes<br>• Support for the 4th Gen AMD EPYC 9004 series processor stack<br>**HPE Cray SC EX425**<br>• Four 2-socket CPU nodes<br>• Support for the full 2nd Gen AMD EPYC 7002 or 3rd Gen AMD EPYC 7003 series processor stack<br>**HPE Cray SC EX420**<br>• Four 2-socket CPU nodes<br>• Support for the 4th Gen Intel Xeon Scalable processor stack (XCC and HBM)<br>**HPE Cray SC EX255a**<br>• Two 4-socket AMD Instinct MI300a Accelerator APU nodes<br>• HPE Cray SC EX254n<br>• Two 4-socket NVIDIA GH200 Grace Hopper Superchip nodes<br>**HPE Cray SC EX235n**<br>• Two 4-socket NVIDIA A100 GPUsx1-socket CPU nodes<br>• Support for the full AMD 3rd Gen AMD EPYC 7003 series processor stack<br>**HPE Cray SC EX235a**<br>• Two 4x AMD Instinct M250X Accelerators 1-socket CPU nodes<br>• Support for the 3rd Gen AMD EPYC processor | 100% fanless DLC |
| **HPE Slingshot interconnect** | HPE Slingshot is a modern high performance interconnect for exascale era that delivers industry-leading performance, bandwidth, and low latency for HPC, AI/ML, and data analytics workloads. HPE Slingshot dramatically controls loaded latency, the true determinant of realized performance on tightly coupled HPC and AI workloads. Innovative congestion management reduces tail latency and run-to-run implementation time variability, and fine-grained adaptive routing helps ensure high utilization of the bandwidth available. | High performance interconnect for HPC, AI/ML and advanced supercomputing workloads | N/A | 100% fanless DLC |

**Table 1.** Compute system workloads and liquid computing compatibilities

This table shows HPE compute product names, photos, typical workloads, and types of liquid cooling that can be used. The compute products include HPE ProLiant, HPE Cray XD, and HPE Cray Supercomputing systems. Individual systems are mapped to the applicable cooling choices: RDHX, HPE ARCS, 70% DLC, and fanless 100% DLC.

[16] Top500 List

# Liquid cooling expertise

This is where HPE legacy of innovation uniquely positions us to deliver critical solutions for today and the future. As you saw earlier, HPE has over 50 years of expertise in solving some of the greatest computing challenges, including heat mitigation with liquid cooling.

## History informing future innovation



**1970s-2010s**

**Cray 1**
Refrigerant cooled

**Cray XT**
Vertical refrigerant cooling

**Cray 2**
Pumped single phase immersion cooled (Fluorinert)

**2010s**

**HPE Apollo 8000**
Liquid cooling

**2020s**

**HPE ProLiant DL365 Gen11**
DLC

**HPE ProLiant DL385 Gen11**
DLC

**HPE ProLiant DL360 Gen11**
DLC

**HPE ProLiant DL380 Gen11**
DLC

**HPE Cray XD670**
DLC

**HPE Cray XD665**
DLC

**HPE Cray 2000**
DLC

**Slingshot interconnect**
100 % fanless
DLC

**HPE Cray EX2500**
100% fanless DLC

**First to exascale**
100% fanless DLC

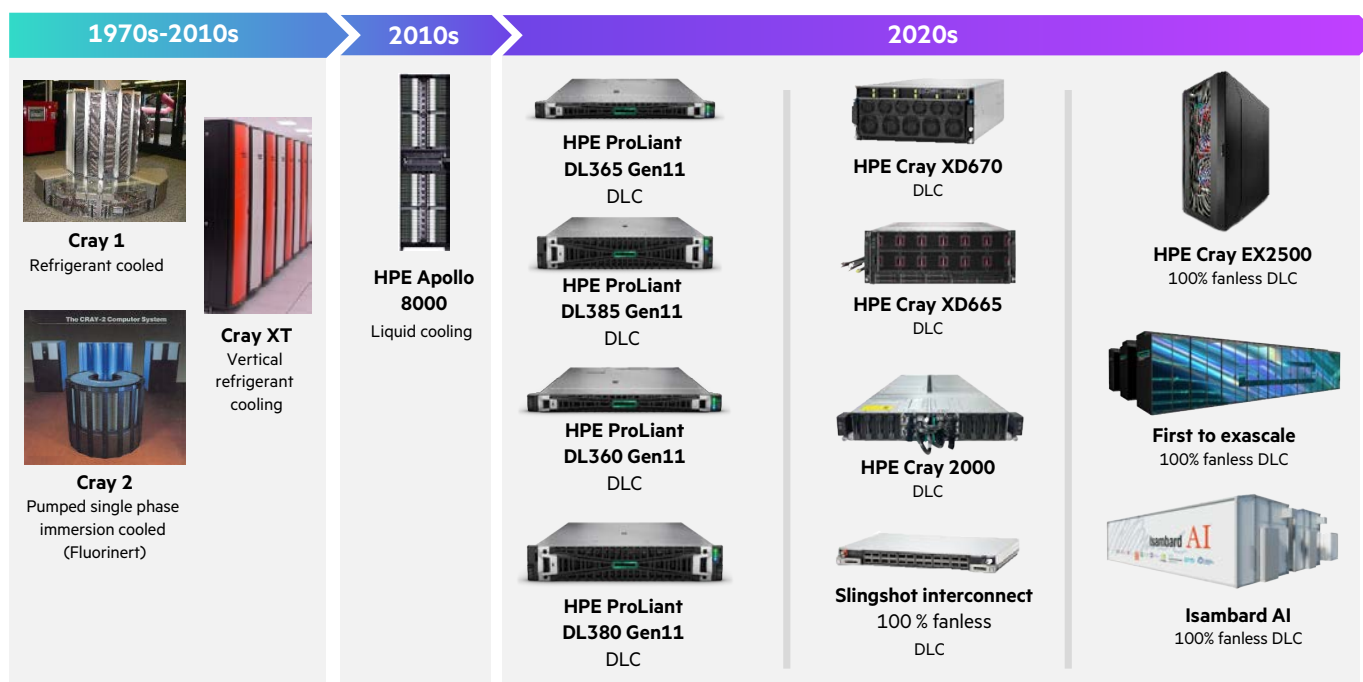**Isambard AI**
100% fanless DLC

**Figure 9.** Timeline of HPE liquid cooling expertise

Figure 9 shows the time span of HPE liquid cooling, beginning with Cray 1 in the mid-1970s, to Cray XT then HPE Apollo, and finally to today's HPE ProLiant, HPE Cray XD, and HPE Cray Supercomputing EX systems.

# Conclusion: Partner with HPE to expertly deploy liquid cooling

AI and business-intensive computing increasingly require higher performing GPUs and CPUs with liquid cooling required to cool those systems. In addition to improving system performance, efficiency, and density, liquid cooling can help achieve up to 94% cost savings over air cooling, over $300 per node per year.[17]
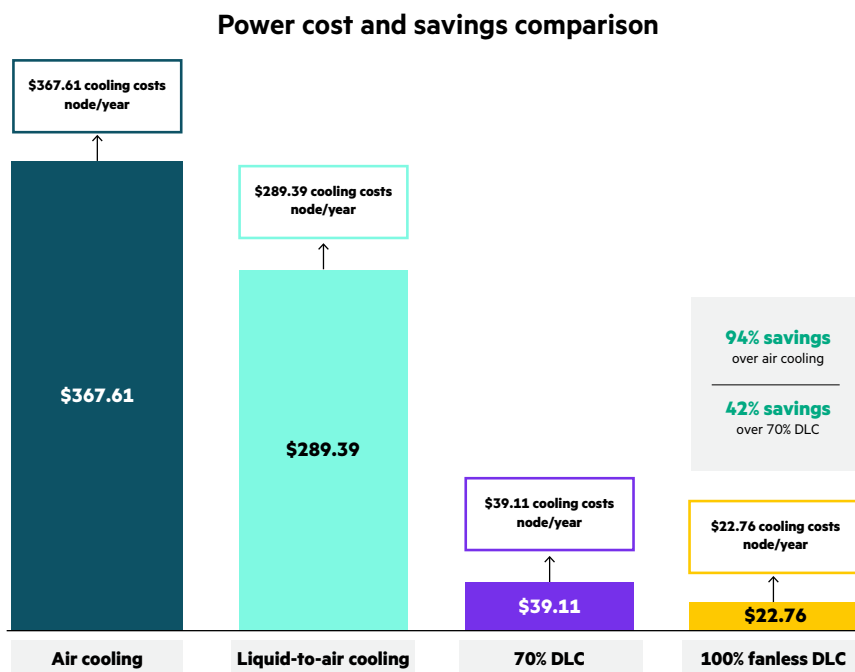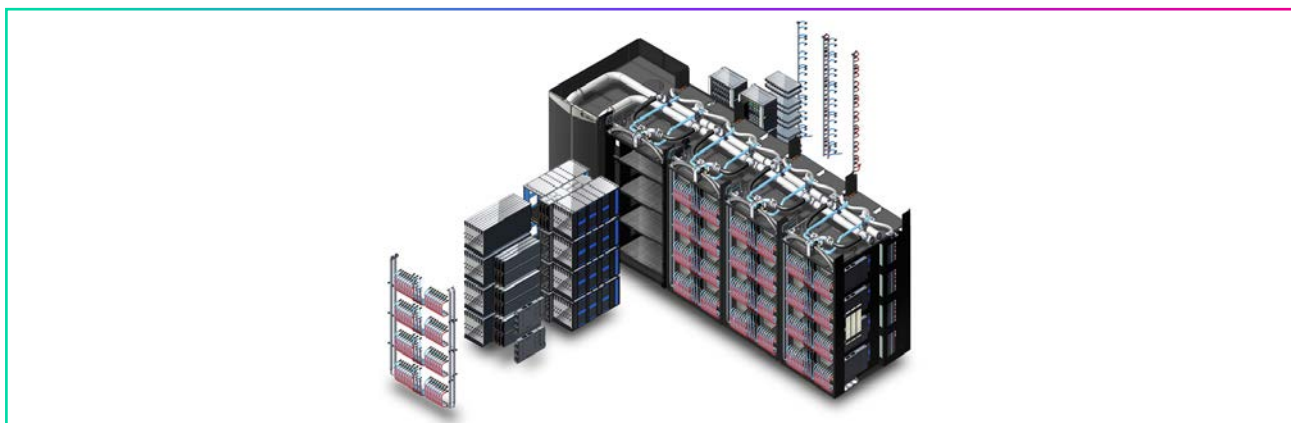
## Power cost and savings comparison



**$367.61 cooling costs node/year**

**$289.39 cooling costs node/year**

**94% savings** over air cooling

**42% savings** over 70% DLC

**$39.11 cooling costs node/year**

**$22.76 cooling costs node/year**

$367.61

$289.39

$39.11

$22.76

| Air cooling | Liquid-to-air cooling | 70% DLC | 100% fanless DLC |

**Figure 10.** Comparison of the average incremental cooling costs per node per year across the four types of cooling.

Only HPE combines 100% fanless DLC expertise, a wide range of liquid-cooled platforms, and a top 10[18] data center services team to successfully tackle today's and tomorrow's AI workloads.



[17] Modeling based upon deployed HPE systems; data compiled in April 2023. These data also reflect optimizations to the data center cooling infrastructure and compute performance as part of a holistic cooling upgrade.

[18] bdcnetwork.com/top-60-data-center-engineering-ea-firms-2022

Visit **HPE GreenLake**

**Chat now (sales)**

**Hewlett Packard Enterprise**