# Optimizing AI Infrastructure:
## The Critical Role of Liquid Cooling

Life Is On | Schneider Electric

# KEY TAKEAWAYS

**1** **Liquid cooling is essential for accelerated computing.**

Graphics Processing Units (GPUs) used in accelerated parallel computing (e.g., AI training) and high-performance computing (HPC) are being introduced with ever increasing power demands. As a result, liquid cooling has become a crucial solution to dissipate heat and maintain optimal performance.

**2** **Direct liquid cooling is the current industry preferred method.**

A coolant distribution unit (CDU) is a necessary component in single-phase direct liquid cooling systems and serves as the link between the data center's cooling infrastructure and liquid-cooled IT equipment. Liquid-to-air and liquid-to-liquid are two common CDU types. CDUs can exist at the rack, row, or room level.

**3** **Careful planning and integration can address concerns related to liquid cooling.**

In most data centers with AI, liquid cooling will co-exist with traditional air cooling for the foreseeable future. Early planning, including infrastructure assessment, is central to successful deployment by addressing concerns related to downtime risk, skill sets, and equipment damage.

**4** **You'll want a diverse ecosystem of partners spanning the entire IT stack.**

Focus on building the right ecosystem of partners. This should include technology vendors, systems integrators, designers, consultants, and service providers, to implement successfully, overcome challenges, and provide ongoing support.

**5** **Plan the cooling in parallel with the IT.**

To avoid delays and problems when deploying AI workloads, cooling architecture planning should occur simultaneously with IT planning. Otherwise, you've invested in very expensive IT gear, only to have it sit unproductive while you figure out a plan to support it.

**6** **Air-cooling is still a necessary part of the architecture.**

While liquid cooling covers most of the heat rejection, it does not cover the entire thermal load. Air cooling is still required and, with the expected IT densification, air-cooled capacity is forecasted to grow even more.

## RATE THIS REPORT
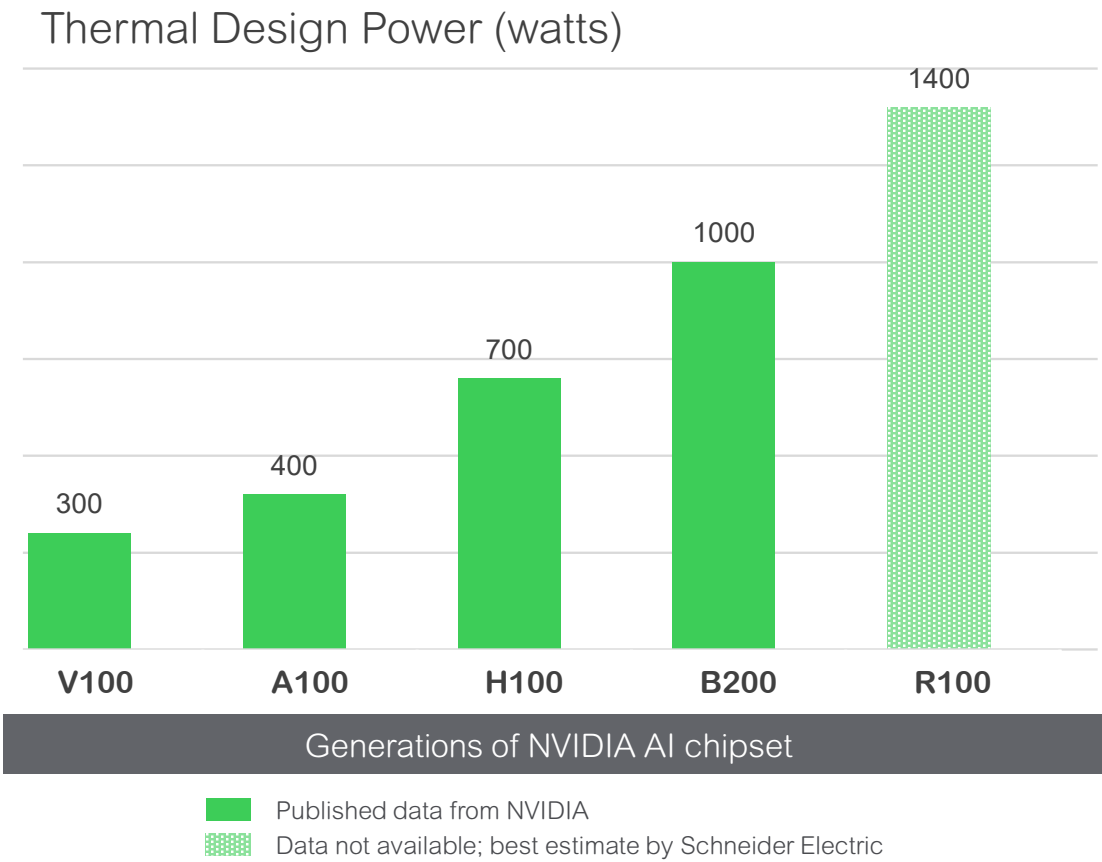★★★★★

LET'S GET STARTED
→

se.com

# The need for liquid cooling

Liquid cooling has become a central topic of conversation as companies explore and validate their AI use cases and begin to integrate it into their operations. Why all the buzz?

The growing demand for AI – particularly the computation-intensive training workloads associated with generative (gen) AI models – is driving the adoption of specialized servers with processors known as accelerators. GPUs are the most common type of accelerator for gen AI, with their parallel processing capabilities, versatility across applications, high memory bandwidth, and scalability. GPUs are delivering performance gains with every new generation, but power consumption grows as well (specified as thermal design power or TDP). See **Figure 1**.[1]

**Figure 1:** TDP is rising with every new generation of NVIDIA AI chipset

Thermal Design Power (watts)



| Generation | TDP (watts) |
|---|---|
| V100 | 300 |
| A100 | 400 |
| H100 | 700 |
| B200 | 1000 |
| R100 | 1400 |

Generations of NVIDIA AI chipset

■ Published data from NVIDIA
▦ Data not available; best estimate by Schneider Electric

That brings up the question: how do we remove heat from these high-powered GPUs? A traditional air-cooled server relies on heat sinks to dissipate the heat from its chips, that have lower TDPs than GPUs. But with TDPs over 700 W per GPU, using a heat sink would require increasing the server's height. That design choice would limit the number of servers that can fit in a rack. So, we hit a critical constraint, since maximizing rack density reduces latency (between GPUs), which reduces training time – an essential factor for AI training investments.

# This is why we're witnessing a fundamental shift toward liquid cooling solutions.

Simply put, liquid removes heat more effectively than air.[2] It enables data centers to support higher compute densities while reducing energy consumption and operational costs.

Liquid-cooled versions of AI training servers are becoming more prevalent, with some models designed to be exclusively liquid-cooled. There are two primary approaches to liquid cooling of servers: direct and immersion.[3] Note, both of these approaches are available with single-phase or two-phase fluids. See White Paper 265, *Liquid Cooling Technologies for Data Centers and Edge Applications*, for more on these architectures.

Direct liquid cooling (single-phase), also referred to as direct-to-chip, has emerged as the preferred method throughout the industry today. *Cold plates* are used to remove heat from server components like GPUs without any contact between the fluid and the server. This method minimizes, or in some cases, eliminates, the reliance on server fans and optimizes space utilization within racks. Other reasons for industry preference include adaptability with existing air-cooled configurations, simple implementation with pre-embedded cold plates, and regulatory advantages over two-phase fluids.

Most direct liquid-cooled servers require a hybrid or blended cooling approach (air plus liquid) since some components in the server still require air cooling. Even if some of your servers are 100% direct liquid cooled, other IT equipment in the data center like storage and networking still require air cooling. Therefore, in most AI data centers, liquid cooling will co-exist with traditional air cooling for the foreseeable future.

It's important to recognize liquid cooling as an architecture rather than a standalone solution. It represents a comprehensive system-of-systems, aimed at optimizing thermal efficiency across different components. This tight integration with IT devices requires careful coordination with existing infrastructure, including traditional heat rejection systems.

se.com

# Navigating the "fear of the unknown"

Although liquid cooling has been around for decades, it has been concentrated in high performance computing (HPC) environments. It has not been in mainstream data centers. And with all things new comes fear and uncertainty. Based on our experience, we believe that executives' apprehension about the unknown boils down to five main concerns:

## 1 Will deploying liquid cooling delay my AI implementation?

Given that liquid cooling is new for most companies, there is uncertainty around how best to plan for it. Finding experienced vendors and installers who can deliver and deploy a liquid cooling system on schedule is a concern. The bottom line is nobody wants their cooling infrastructure to be the bottleneck for getting their AI investment up and running.

## 2 Can I obtain the right skill sets for operations and support?

There is a common concern that existing data center staff, both on the facilities and IT side, won't have the right knowledge and skills to operate and maintain a liquid cooling architecture. Some of the components will be unfamiliar to them, and tasks like performing maintenance on, or replacing a server with water pipes appears daunting. The increasing complexity, coupled with the specialized skills required to operate and maintain liquid cooling systems (MOPs, SOPs),[4] makes it difficult for organizations to find qualified personnel.

## 3 Will liquid in my IT space create added downtime risk?

Water has always been part of data center cooling architectures – pipes under the raised floor, pipes in the mechanical room, etc., but we have a deep-seated fear of bringing it closer to the IT gear. Anything that appears to increase the risk of data center downtime is something we've always discouraged.

Water pipe or connection leaks are a prevalent concern that often comes to mind, driven by the potential extended downtime risk it poses.

There are additional concerns around corrosion and the compatibility of materials used in liquid cooling systems. Corrosion can compromise the integrity of components over time, leading to early failures and eventual downtime. For this reason, stainless steel piping is typically used, along with a water-based mixture of propylenic glycol and other corrosion inhibitors to prevent biological growth and avoid corrosion.

Adding liquid cooling to a data center environment can also bring increased maintenance. Adding any new maintenance step to your data center operations represents an increased risk of downtime caused by human error.

Equipment compatibility and interoperability are also tied into this fear of downtime, as existing systems may not be designed to work with new cooling technologies. This lack of standardization can result in increased integration and maintenance complexity, which heightens risk of system failures, contributing to potential downtime.

## 4 Does liquid cooling increase the risk of damaging my expensive GPU servers?

With investments in servers on the order of $200k or more per server, it is no surprise there is concern about any potential to damage the servers. Years of conditioning have taught us fluids and computing equipment end in disaster.

When liquid is pumped into servers, it can cause rapid temperature fluctuations that stress the components. Thermal shock occurs when there is a sudden change in temperature, which can lead to physical damage or failure of sensitive electronic parts, such as the GPUs. Additionally, there is fear that fluid leaks inside the server can damage GPUs and other hardware, resulting in malfunction or complete failure.

## 5 Will liquid cooling void my server warranty?

Warranties for liquid-cooled IT equipment generally cover defects in materials and workmanship, including cooling system components inside the server. Damage from aftermarket modifications, poor maintenance, and non-approved cooling fluids can void the warranty. Some server manufacturers could be restrictive about what types, brand, and configuration of cooling system they would accept for their IT equipment. This includes manifolds, piping from manifold to server, connectors, regulation valves, and pumps. So, it's essential to adhere to the manufacturer's guidelines and warranty restrictions.

## No need to panic.

While the transition to liquid cooling may seem challenging, there is no need for concern. By addressing initial fears through proper execution and an ecosystem of trusted partners, organizations can position themselves for success. No doubt, the data center industry must adapt and become comfortable with these advancements, as it offers a pathway to meet the demands of modern technology.
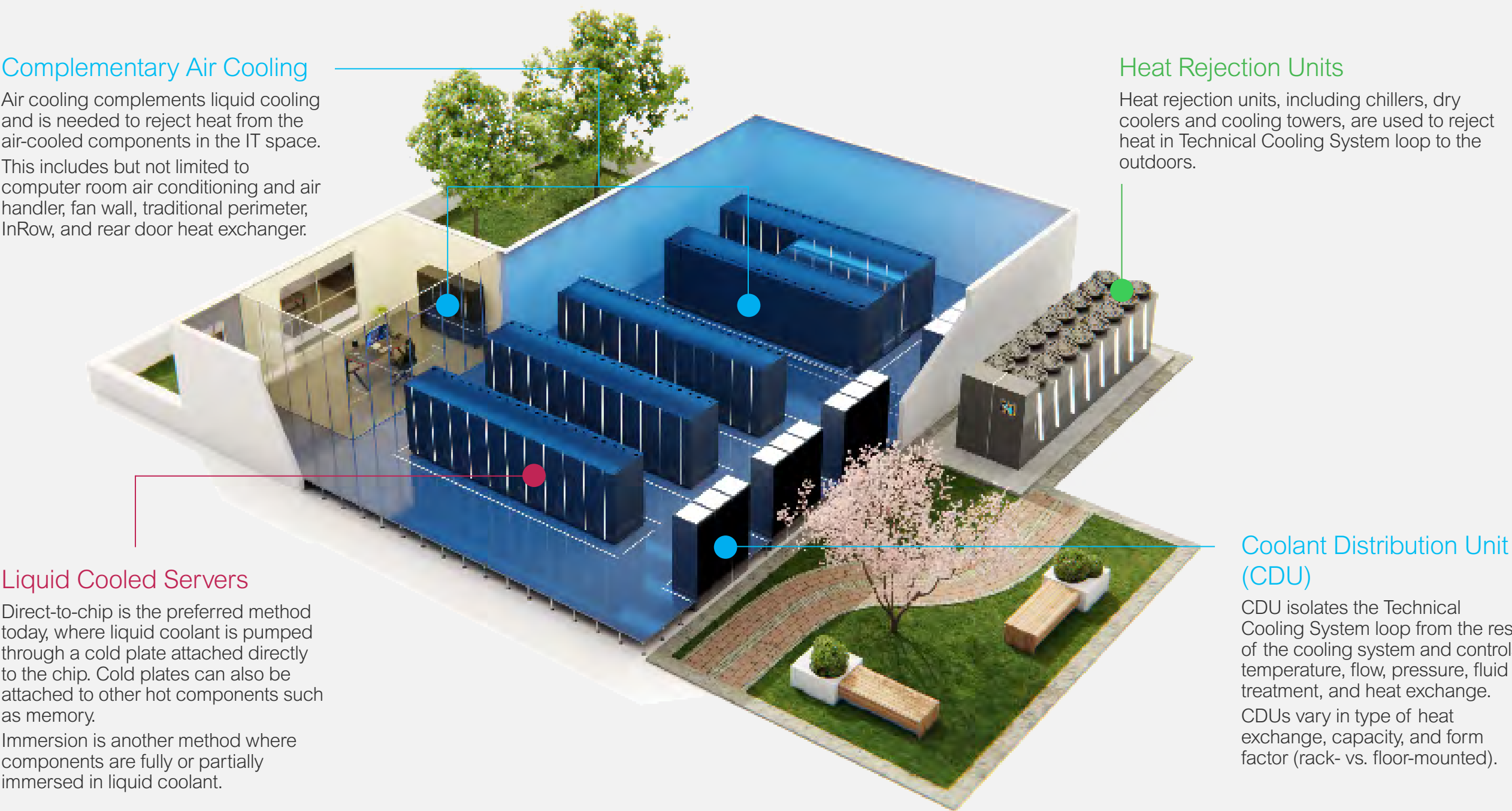
# Key elements to a direct liquid cooling architecture

While liquid cooling architectures are different than traditional chilled-water cooling architectures, there are also some things that stay the same. In both architectures, heat is transferred from the IT equipment to the outdoors using a heat rejection system located outside. In fact, in many retrofit cases, you can use much of the existing physical infrastructure, as opposed to installing a new dedicated heat rejection system for liquid-cooled servers.

*Note, if your project is a new / greenfield data center, you have additional freedom to specify a chiller that provides the best balance in efficiency and performance for both air-cooled and liquid-cooled loads.*
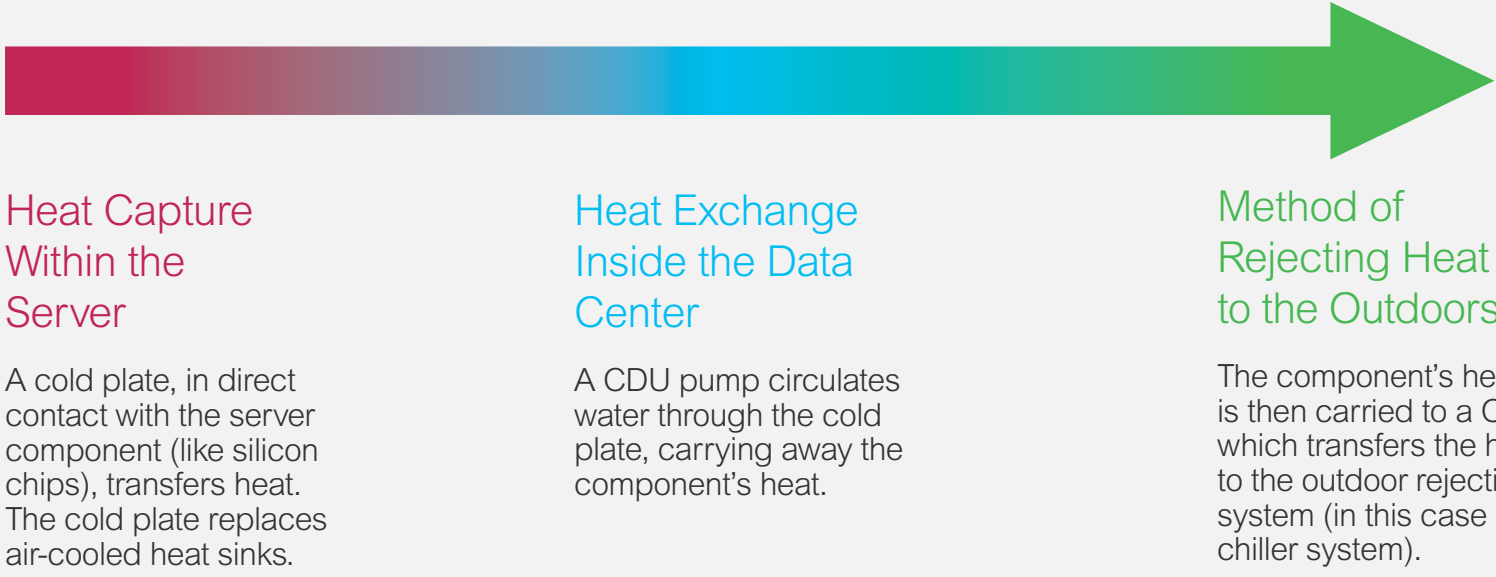
So, what's really different between air-cooled and direct liquid-cooled architectures? Liquid cooling architectures use coolant distribution units (CDUs), that transfer the heat from the cold plates to another part of the cooling system. **Figure 2** describes the three main elements of this direct liquid cooling ecosystem.

**Figure 2:** Understanding chip-to-chiller



**Complementary Air Cooling**
Air cooling complements liquid cooling and is needed to reject heat from the air-cooled components in the IT space.
This includes but not limited to computer room air conditioning and air handler, fan wall, traditional perimeter, InRow, and rear door heat exchanger.

**Heat Rejection Units**
Heat rejection units, including chillers, dry coolers and cooling towers, are used to reject heat in Technical Cooling System loop to the outdoors.

**Liquid Cooled Servers**
Direct-to-chip is the preferred method today, where liquid coolant is pumped through a cold plate attached directly to the chip. Cold plates can also be attached to other hot components such as memory.
Immersion is another method where components are fully or partially immersed in liquid coolant.

**Coolant Distribution Unit (CDU)**
CDU isolates the Technical Cooling System loop from the rest of the cooling system and controls temperature, flow, pressure, fluid treatment, and heat exchange.
CDUs vary in type of heat exchange, capacity, and form factor (rack- vs. floor-mounted).

**Chip-to-Chiller**
3 main elements in a liquid cooling ecosystem

**Heat Capture Within the Server**
A cold plate, in direct contact with the server component (like silicon chips), transfers heat. The cold plate replaces air-cooled heat sinks.

**Heat Exchange Inside the Data Center**
A CDU pump circulates water through the cold plate, carrying away the component's heat.

**Method of Rejecting Heat to the Outdoors**
The component's heat is then carried to a CDU which transfers the heat to the outdoor rejection system (in this case a chiller system).

CDUs can transfer heat in two ways. One uses a liquid-to-air heat exchanger (like a radiator). The second uses a liquid-to-liquid heat exchanger. CDUs can exist at the rack, row, or room level.
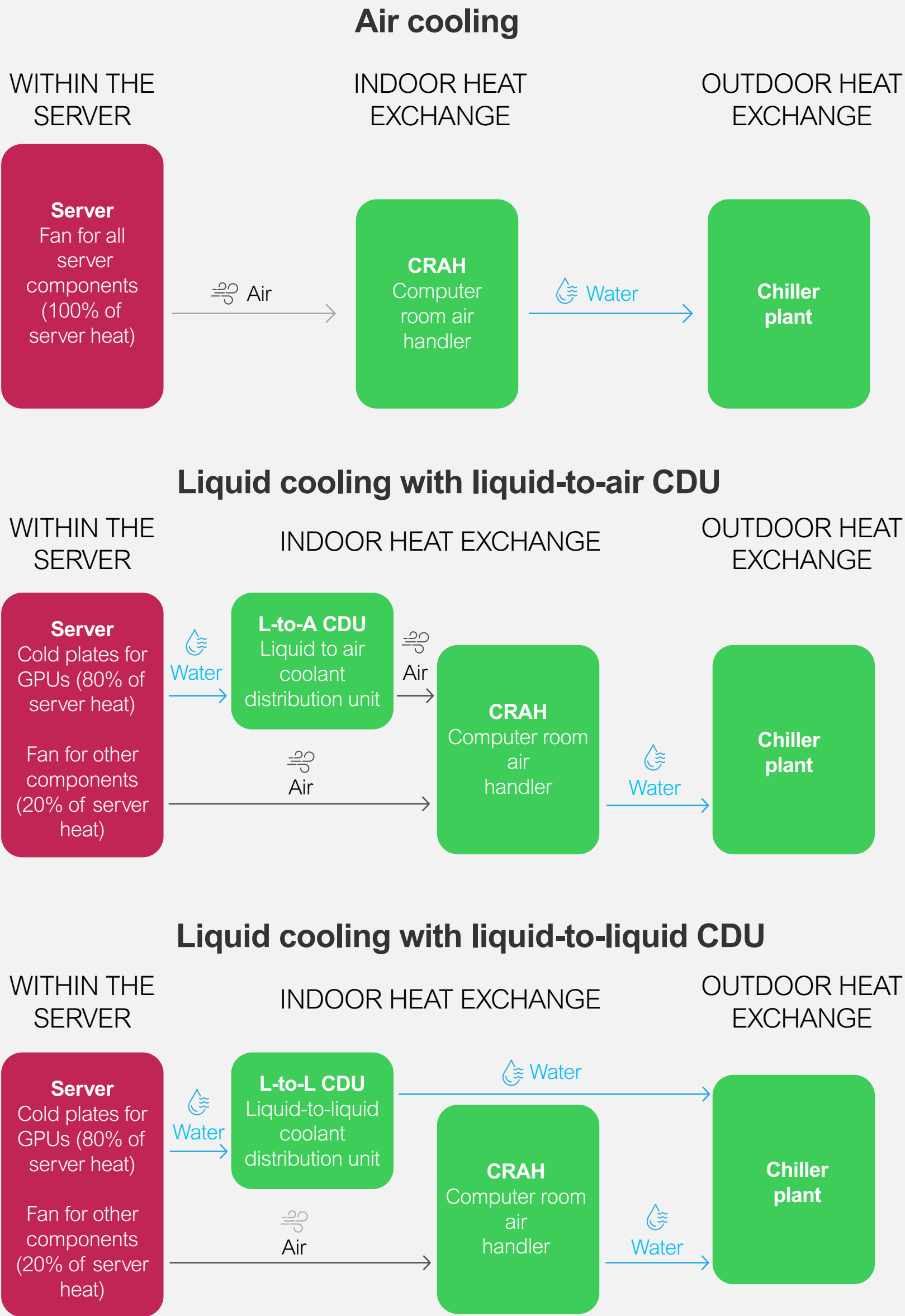
**Figure 3** helps illustrate where a CDU becomes part of the heat transfer.

The top diagram (3a) represents traditional air-cooled servers using a computer room air handler (CRAH) to send heat to a chilled-water plant. No CDU is required.

The second (3b) shows direct liquid-cooled servers using a liquid-to-air CDU that sends heat to a CRAH, which then sends heat to a chilled-water plant.

And the last (3c) shows direct liquid-cooled servers using a liquid-to-liquid CDU that sends heat to a chilled-water plant. The CRAH is still needed for the air-cooled components within the servers.[5]

**Figure 3:** Simplified diagrams of air cooling vs. direct-to-chip liquid cooling (3 parts)

## Air cooling



## Liquid cooling with liquid-to-air CDU



## Liquid cooling with liquid-to-liquid CDU



From these diagrams and descriptions, a reasonable question arises: *Instead of using a CDU, why not just use the IT water directly from the chiller plant?* One key answer lies in understanding that the IT water flows through tiny channels in the cold plate which are susceptible to clogging if the water isn't filtered and treated (**Figure 4**).

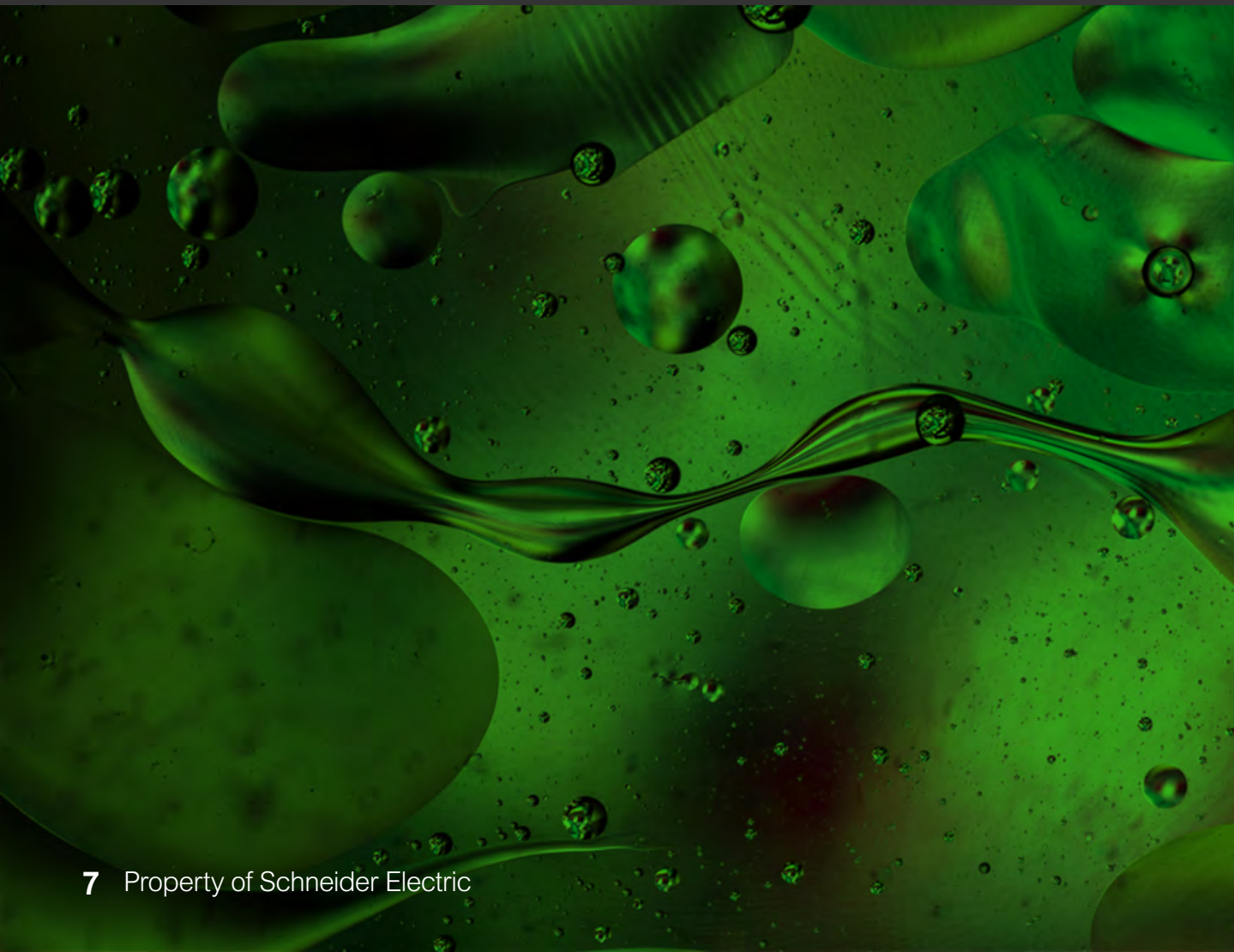**Figure 4:** Microchannels inside the cold plate in the server



Source: Motivair

Filtering the water removes particles greater than about 25-50 microns [6] (less than half the thickness of a human hair), and chemically treating the water [7] prevents biological growth and fouling. Therefore, the need for the CDU… In essence, the **CDU uses a heat exchanger to isolate the facility water system from the IT system**.[8] The CDU also provides other key functions including regulating water flow, water temperature, and water pressure. The CDU's control system manages these functions that keeps the temperature across the chips steady, so as not to thermally shock and damage the chips.

Having explored the architecture distinctions of liquid cooling, we now shift our focus to the planning guidance necessary for successful implementation in your AI data center, so that the common concerns associated with deploying liquid cooling are addressed.

# Planning for liquid cooling

Given the power, cooling, and rack challenges[9] that go along with deploying AI workloads, the physical infrastructure planning must occur in parallel with the IT planning. Doing this planning early will reduce risks of problems and delays. This should be a collaborative effort between IT and facilities. Otherwise, you've invested in very expensive IT gear, only to have it sit unproductive while you figure out a plan to support it. Below we provide guidance to address the five common concerns noted earlier.

## Avoiding delays in your AI implementation

We know of instances where an organization purchased multiple NVIDIA DGX servers only to find out that they didn't have enough power capacity to install them. This same scenario can occur with cooling. If you're building a new data center, the design process should account for and optimize around liquid cooling requirements. In the case of existing data centers, a site assessment is one of the first things your team should do to identify risks to your AI deployment schedule.

Here are 10 checklist items that will help reduce the likelihood of project delays:

### Capacity and redundancy

》 Make sure your cooling system has enough spare capacity to support the planned AI servers.

》 Determine if your uninterruptible power supply (UPS) has spare capacity to support the CDU pumps. A loss of cooling could force the server to shut down or even worse, damage the chips.

》 Understand the redundancy in your chiller plant. Having redundant components makes deployment less disruptive to your existing operational critical loads.

### Structural considerations

》 Check if you have enough ceiling clearance to install air containment around the AI rack cluster. This will re-duce the likelihood of hot spots and improve efficiency.

》 Validate that the data center slab or raised floor is rated to support the weight of the proposed AI rack(s). There have been real cases of expensive AI training racks crashing through floors that weren't rated properly.[10] When validating the weight, make sure you use the "wet weight", meaning the weight is inclusive of fluid in the IT equipment and CDUs.

### Cooling system details

》 Determine if you need a chiller plant. If you don't have one, it becomes less realistic to deploy liquid cooling as the number of liquid-cooled servers increases.

》 Determine if the chiller plant can supply water temperature that is compatible with both air-cooled and liquid-cooled IT equipment. Note, if you have an economizer mode, check that it is compatible with the proposed chilled-water temperature.

》 Check if you have spare valves on the chilled-water loop that can be used to supply a CDU.

》 Assess if your data center infrastructure management (DCIM) system can monitor the liquid cooling system, including leak detection. If you don't have DCIM, you should consider implementing it.

》 Confirm the cooling distribution equipment is on the server manufacturer's approved vendor list and doesn't violate warranty policy. This could include manifolds, piping from manifold to server, connectors, regulation valves, pumps, and CDUs.

se.com

## Gain the right skill sets for operations and support

It's probable that your data center facilities team doesn't have direct experience with CDUs and water distribution to liquid-cooled servers. But they should have experience with chiller plants, which gives them a good foundation upon which to build their knowledge. Use an ecosystem of partners to build upon this foundation. Vendors who design and manufacture CDUs understand the complexities of liquid cooling and can help mitigate risks. When it comes to installation, vendors should also have a list of recommended installers across different trades that understand these systems. Below are just a few examples of how experienced installers de-risk the installation and operation of your liquid-cooled system:

**Example 1:**
Pressure test preparation: They flush new piping to ensure cleanliness before pressure testing and chemical treatments.

**Example 2:**
Isolation of components: They isolate circulating pumps, cooling coils, heat exchangers, and existing piping during initial flushing and drainage to prevent contamination.

**Example 3:**
Final cleaning: They use a non-foaming liquid detergent oil-dispersing cleanser to remove oil and foreign matter from the piping and equipment prior to the final filling of the systems.

## Reduce downtime risk

There was a time when chilled-water piping was installed in large troughs, formed into the data center's concrete floor. This, now antiquated practice, mitigated the risk of flooding the data center in case the steel pipe catastrophically failed.

Over time, experience allowed operators to became comfortable with this tiny risk. Designers would later build data centers without raised floors and with water piping overhead. The water risk has become so well managed that operators began installing row-based chilled-water CRAHs in between production IT racks. Just as chilled water in the IT space presents managed risks to data centers worldwide, so does water inside racks and servers.

The reason why there's little concern for chilled water in data centers today is the experience gained over decades. This same experience applies to liquid cooling architectures. The key is to manage the risk by using a combination of design, installation, and operation best practices, and to work with trusted vendors, installers, and operators.

## Reduce the risk of damaging expensive GPU servers

Like cooling vendors, server vendors have had decades of experience when it comes to water inside their servers. That experience has led to design, manufacturing, and process improvements. And because of this, they accept what minimal risk of failure remains, as evidenced by their equipment warranties.

For example, in 2014, IBM released the water-cooled IBM NeXtScale System. In 2016, Dell introduced its third-generation liquid cooling solution, Triton. Properly installing and maintaining liquid-cooled servers is critical to mitigating these risks and protecting valuable equipment in data center environments. Documents like IBM's "Water cooling system specification" provide detailed information for deployment teams.

Once again, the ecosystem of partners you use, server vendors included, form the basis of risk mitigation around your AI deployment.

se.com

# NEXT STEPS

The evolution to liquid cooling in data centers presents both challenges and opportunities. To navigate this shift successfully, businesses should take the following **four proactive steps**:

## 1 Educate and inform decision teams.

Familiarize your business with the technology. Learn the lingo. Your teams can leverage materials like white paper 133, Navigating Liquid Cooling Architectures for Data Centers with AI Workloads. Understand the role a CDU plays in a liquid-cooled architecture, and what attributes to look for. This will help your organization participate in liquid cooling conversations and make informed decisions. AI and liquid cooling technology are evolving fast, so stay up to date!

## 2 Assess current infrastructure

You should audit and analyze your existing infrastructure to develop a solid understanding of existing cooling capacities, current loads, and future requirements. Evaluate how integrating liquid cooling will affect the existing systems, including the compatibility of current equipment with liquid cooling solutions, potential modifications needed for piping and power distribution, and how new heat loads will interact with existing systems.

## 3 Engage and plan with an ecosystem of partners

Collaborate with technology partners, including IT vendors, cooling specialists, systems integrators, and service partners. They can help plan, design, implement, and maintain your liquid cooling architecture, in alignment with your business and operational goals. There's no need to do it alone. And when engaging with your cooling vendor, ask specific questions to plan for a successful integration. Inquire about how their CDU solution will integrate with your existing data center infrastructure, as well as the expected lifespan and maintenance requirements. Additionally, seek clarification on the support and service offerings available post-installation to understand how they will assist you in maintaining optimal performance and reliability over time.

## 4 Remain focused on sustainability

As discussed, liquid cooling of high-density AI racks is more energy efficient than traditional air cooling. Compare the specific systems you choose (like the CDU) with your target energy and sustainability standards, if applicable. By tightly integrating controls, you can automate temperatures to further reduce your energy consumption. Software monitoring, management, and reporting also helps your data center track energy consumption and cooling efficiency in real-time, allowing for adjustments to further optimize and ensure compliance with sustainability standards and regulatory requirements. There's no need to compromise on sustainability goals; liquid cooling can reduce your carbon footprint and support your overall sustainability objectives.

# Endnotes

1    Reference links to TDPs: V100, A100, H100, B200

2    "Water is a better conductor of heat by a factor of more than 23 compared with air, and can hold far more heat than air, around 3,243 times more by volume.", Schneider Electric white paper 265, Liquid Cooling Technologies for Data Centers and Edge Applications, page 12

3    "With immersive liquid cooling, the liquid coolant is in direct physical contact with the IT electronic components. The servers are fully or partially immersed in a dielectric liquid coolant covering the board and the components, which ensures all sources of heat are removed.", Schneider Electric white paper 265, Liquid Cooling Technologies for Data Centers and Edge Applications, page 4

4    Schneider Electric white paper 178, A Framework for Developing and Evaluating Data Center Maintenance Programs, page 4

5    Note, this diagram applies to data centers with chilled-water plants. But in cases where computer room air-conditioners (CRACs) are used, refrigerant would transport the server's heat to the outdoors via the outdoor condenser.

6    OCP, ACS Liquid Cooling Cold Plate Requirements Document, page 23

7    The water is treated by mixing it with about 25% glycol.

8    The IT system in this context is also known as the technology cooling system.

9    These challenges are explained in Schneider Electric White Paper 110, The AI Disruption: Challenges and Guidance for Data Center Design.

10    EMRC research, based on user conversations.

**RATE THIS REPORT**

★★★★★

se.com

# Authors

### Victor Avelar

Chief Research Analyst
Data Center Research &
Strategy
Schneider Electric
LinkedIn

**Victor Avelar** is a seasoned expert in data center energy efficiency and design, serving as the Chief Research Analyst at Schneider Electric's Data Center Research & Strategy group. With over 25 years of experience, Victor leads cutting-edge research and best practice development for sustainability, risk management, and next-generation data center technologies. He's a trusted advisor to clients globally, providing actionable insights on enhancing infrastructure performance through innovative solutions such as liquid cooling and energy modeling. Known for his clear, practical guidance, Victor helps organizations tackle the evolving challenges of sustainable and efficient data center operations. He is central to the development of technology adoption forecasts for data centers. He also leads the peer review process for all EMRC content. Victor holds a bachelor's degree in mechanical engineering from Rensselaer Polytechnic Institute and an MBA from Babson College. He is a member of AFCOM and a sought-after speaker on AI infrastructure.

### Wendy Torell

Senior Research Analyst
Data Center Research &
Strategy
Schneider Electric
LinkedIn

**Wendy Torell** is a Senior Research Analyst in Schneider Electric's Data Center Research & Strategy group bringing 30 years of data center experience. Her focus is analyzing and measuring the value of emerging technologies and trends: providing practical, best practice guidance in data center design and operation. Beyond traditional thought leadership, she championed and leads development of interactive, web-based TradeOff Tools. These calculators help clients quantify business decisions, while optimizing their availability, sustainability, and cost of their data center environments. Her deep background in availability science approaches and design practices helps clients meet their current and future data center performance objectives. She brings a wealth of experience across Schneider Electric's broad portfolio and with the market at large. She holds a BS in Mechanical Engineering from Union College and an MBA from University of Rhode Island. Wendy is an ASQ Certified Reliability Engineer.

Life Is On | Schneider Electric

se.com