

# What’s the Difference Between Immersion and Direct-to-Chip Liquid Cooling?

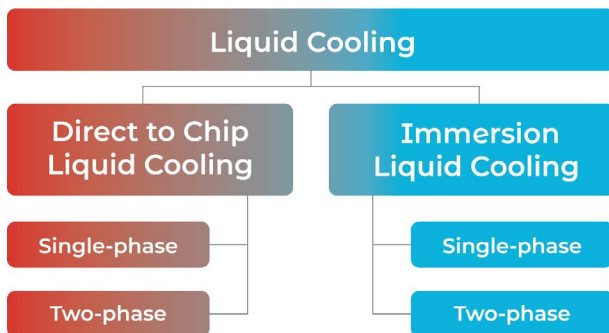
The next generation of high-performance CPUs and GPUs targeting the demanding needs of AI and HPC will generate a significant amount of heat, making efficient cooling solutions critical for optimal system performance and sustainability.

A massive shift is taking place in today’s data centers as rack power consumption skyrockets to levels never thought possible. Driven by the emergence of compute-intensive artificial-intelligence (AI) and high-performance computing (HPC) applications, data centers have quickly transitioned from needing to cool 10- to 20-kW racks with air-cooling strategies to cooling [120-kW racks](#) powered by NVIDIA’s Grace Blackwell superchip—and that’s just for one rack!

Air-cooling methods don’t stand a chance of cooling this amount of heat, and, as a result, has paved the way for novel liquid-cooling technologies that fall into one of two categories: “direct-to-chip” or “immersive.” Unlike traditional air-based approaches, these technologies use liquid, either water or dielectric fluid, to remove heat from the equipment.

As the industry marches toward a sustainable AI future where AI factories emerge to meet growing compute demands, liquid-cooling technologies are poised to become a critical enabler for data centers already challenged with controlling their heat dissipation, energy consumption, and footprint. In fact, with the introduction of next-generation GPUs that boast up to 1,200 W, liquid cooling has quickly gone from a “nice to have” to an “absolute requirement.” As the demand for this technology escalates across the globe, analysts such as [Mordor Intelligence](#) are estimating it will reach \$14.8 billion 2029.

In this article, we’ll explain the difference between both immersive and direct-to-chip liquid cooling—each consisting of a single phase and a two-phase option. We’ll then outline the pros and cons of each approach as it relates to overall sustainability, power consumption, ease of use, risk, scalability, and costs (*Fig. 1*).



1. Market segmentation showing variants under direct-to-chip and immersive cooling options.

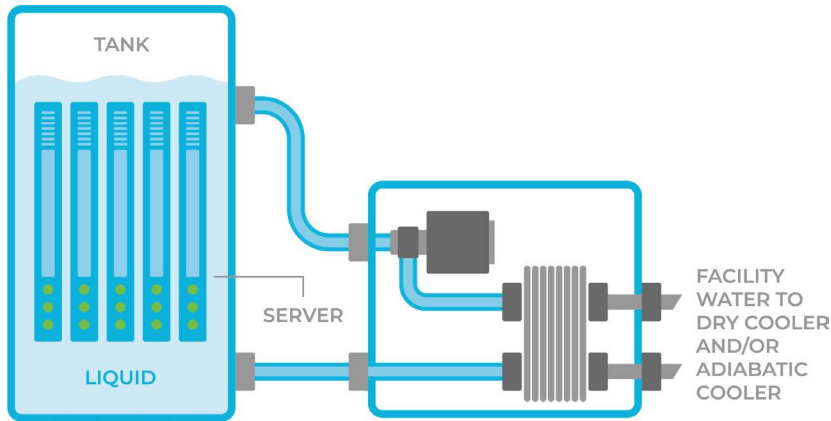
## Immersive Liquid Cooling Submerges Components

With immersive liquid cooling, servers and other components are completely submerged in dielectric liquid. The heat generated from the equipment is transferred to the surrounding coolant. As the heated coolant rises to the surface, it’s then sent to a cooling system to dissipate the accumulated heat and return it to the original tank housing the equipment.

There are two types of immersive liquid cooling:

### Single-phase immersion

With this method, all servers and other IT equipment are



**2. Single-phase immersion: Equipment submerged in dielectric fluid.**

immersed in dielectric fluid. As the CPU or GPU heats up, the fluid absorbs the heat. This heated fluid is then pumped to a heat exchange unit, where it cools the fluid and then sends it back to the tank housing the hardware (Fig. 2).

*Pros:*

- Complete server power absorption, which means that all of the heat coming out of the server (GPU, CPU, DIMS etc.) is collected and cooled.
- The use of dielectric fluid will not short-circuit the components and the servers.

*Cons:*

- The thermal design power (TDP) is limited. If a GPU has a TDP over 700 W, the single-phase immersion method is unable to effectively cool the hardware.
- Requires significant data-center infrastructure investment because large and heavy tanks are now required to house the equipment. This technology is better suited for new data centers or those with the space and ability to be significantly modified to fit the tanks. They should also have the structural support needed for the added weight.
- All equipment immersed in the tank—e.g., servers, connectors, and PCBs—needs to be compatible to the dielectric liquid so that it's not damaged by the fluid itself. This often requires the selection of specialized equipment or a modification to the servers.
- Because some of the components on the server, such as fiber-optics connectors, can't function when immersed, the servers require mechanical reconfiguration.
- The type of fluid used, which features hydrocarbons, are flammable and combustible. This could cause catastrophic damage if a fire were to occur inside the data center.
- Maintenance of the servers is difficult where any required server maintenance requires "pulling" the single server out of the tank using a crane and letting it drip for 30

minutes before starting to service.

- Any contamination of the fluid, such as, for example, with water, needs draining and cleaning of the tank, which may result in a full day of downtime.

**Two-phase immersion**

Similar to single-phase immersion, this method also involves completely immersing the servers and IT equipment in dielectric fluid. However, as the components on the board heat up, it boils the fluid, which creates vapor that rises from the liquid to the top of the tank. Located on the top of the tank is a

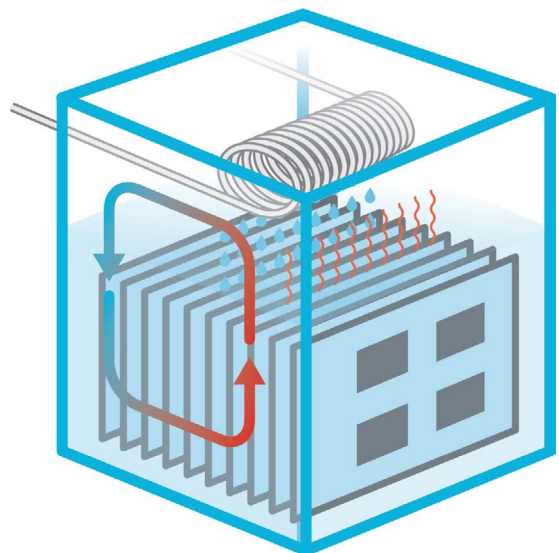
network of tubs that flow cooled water. The vapor from the tank touching the cold tubs condenses and drips back as liquid to the tank; the water in the tubs is heated, carrying the heat to the facility via the heated water and removed from the data center (Fig. 3).

*Pros:*

- Complete server power absorption, which means that all of the heat coming out of the server (GPU, CPU, DIMS, etc.) is collected and cooled.
- Can support very high TDPs.
- The use of dielectric fluid will not short-circuit the components and the servers.

*Cons:*

- All equipment immersed in the tank—e.g., as servers, connectors, and PCBs—needs to be compatible to the dielectric liquid so that it's not damaged by the fluid



**3. Two-phase immersion: Server equipment immersed in dielectric liquid.**



**4. ZutaCore's HyperCool cold plate uses dielectric liquid to remove heat on demand.**

itself. This often requires the selection of specialized equipment or a modification to the servers.

- As part of fluid boiling, the aggressive cavitation process causes damage to IT parts, the PCB, and windings.
- Requires significant data-center infrastructure investment because large, heavy tanks are now required to house the equipment, and added support is needed to carry the weight.
- Maintenance often involves long downtimes with the use of cranes due to the weight of the tanks and the equipment being immersed in liquid.
- Each time the tank is opened for service, vapors that contain PFAS (per- and polyfluoroalkyl substances) are released into the environment, resulting in a 10% loss of liquid (100s of liters) per year and the release of large amounts of PFAS vapors into the atmosphere.

#### **How Direct-to-Chip Liquid Cooling Works**

Unlike immersion liquid cooling that submerges the entire server and other IT equipment in fluid, the direct-to-chip process brings cooling liquid to a cold plate that's placed directly on top of the high heat flux, such as CPUs and GPUs. This liquid removes heat from the components. It's contained in the cold plate and never touches the chips or server equipment.

The direct-to-chip method is widely considered to be faster and more efficient than other types of cooling because it can target the areas generating most of the heat. In fact, at a recent [Omdia Analyst Summit](#), their analysts stated: "As rack sizes grow above 50kW, direct-to-chip technologies will dominate."

As cold plates are mainly located on the high heat flux

components, server fans are required to remove the excess heat from the low heat flux components.

There are two types of direct-to-chip liquid cooling:

#### **Single-phase direct-to-chip cooling**

This method of liquid cooling uses water as the coolant in the cold plate. Water always remains in a liquid state and the ability to take away heat with this method depends on water flow—the higher the heat, the more flow of water is required.

*Pros:*

- With a high flow rate of cold water, it can support cooling of high TDP components.
- Data-center infrastructure and server are much like air cooling; small changes are required to add the cold plate cooling.

*Cons:*

- The risk of leakage is significant and costly because water is being used instead of dielectric fluid. Not only can water leak destroy \$300K servers, but it also leads to corrosion, mold, residue, biological growth, and other environmental impacts.
- Requires the investment of larger pipes, tubs, and connectors that can resist leakage of water as flow and pressure increases. The water flow also requires larger, power-hungry pumps to continually carry the water through the system.

#### **Two-phase direct-to-chip cooling**

Unlike the single-phase direct-to-chip approach that uses water, the two-phase approach leverages dielectric liquid that is 100% safe to the IT equipment. The heat from the GPUs and CPUs boils the dielectric fluid at low temperature, absorbing the heat using the very efficient phase-change physical phenomena keeping the chip at a constant temperature. This is akin to how boiling water keeps the bottom of a pasta pot at 100°C, just at a lower temperature (*Fig. 4*).

The process of changing the state of the dielectric fluid from liquid to gas and then back to liquid again is done in a completely closed-loop system. As the liquid inside the cold plate boils, the liquid in the cold plate never passes the boiling temp even if the heat increases by 3X (aka hotter GPUs and CPUs). It thus makes the technique highly scalable for cooling higher-power chips in the future.

This is in contrast to the single-phase direct-to-chip approach that depends on the flow of water (and lots of it) to cool chips. To put this in perspective, a 100-kW rack using two-phase direct-to-chip technology will use less than four gallons of dielectric fluid, compared to immersion cooling that needs over 100 gallons per rack.

*Pros:*

- Waterless direct-to-chip liquid cooling: Dielectric fluid is 100% IT safe and will not damage the server even in

the unlikely event of a leak.

- Increases compute density: Allows for >150 kW per rack.
- Future-ready: Up to 2500 W per chip and beyond.
- Can save up to 80% power consumption compared to air-cooling.
- Because the liquid maintains a constant temperature, the heat from the servers can be harvested for heat-reuse applications such as heating adjacent offices, other parts of the data center, or even schools and office buildings in proximity.
- Requires little to no data-center infrastructure change, leading to low up-front investment costs, with a simple installation process.
- Low maintenance requirements: The dielectric fluid never needs to be replaced and unlike immersion, it doesn't get released into the atmosphere during server and rack maintenance. The liquid used has been engineered with an ozone depletion potential (ODP) of 0 and very low global warming potential (GWP).
- Maintains a 1U server form factor, even as the heat increases for next-gen GPUs.

*Cons:*

- Liquid cooling is used to dissipate heat only from CPU/GPU. Air cooling is still required for cooling other components such as memory, I/O, and others.

### **All Roads Lead to Liquid Cooling**

The future growth of AI will largely depend on the ability to expand data-center capacity, which will push the heat inside these facilities to unprecedented levels. As this article describes, many different types of liquid-cooling technologies are now available to help remove this heat, but each has its own set of pros and cons.

It's up to data centers and hyperscalers of the world to decide which solution works better for them. This should be weighted in terms of cost, power consumption, ease-of-use, scalability, and sustainability. Only then will the industry be able to arrive at true AI sustainability.

*Shahar Belkin has nearly 20 years of experience developing companies, new technologies, algorithms, and products for security and remote video. He has received a number of international patents and was the founder of OzVision and FST Biometrics. He served as VP R&D of ZutaCore for three years prior to his current position as EVP Product. Shahar holds a B.Sc. in electrical engineering.*