

REPORT REPRINT

Liquid cooling is about to go from trickle to torrent

MAY 03 2019

By Daniel Bizo

It's been a long time coming, but direct liquid cooling (techniques that deliver a coolant fluid to heat sources, as opposed to relying on air as a heat-transfer medium) may finally gain acceptance as pressure on datacenter performance grows. Here we lay out the major factors that we believe will drive the shift.

THIS REPORT, LICENSED TO ZUTACORE, DEVELOPED AND AS PROVIDED BY 451 RESEARCH, LLC, WAS PUBLISHED AS PART OF OUR SYNDICATED MARKET INSIGHT SUBSCRIPTION SERVICE. IT SHALL BE OWNED IN ITS ENTIRETY BY 451 RESEARCH, LLC. THIS REPORT IS SOLELY INTENDED FOR USE BY THE RECIPIENT AND MAY NOT BE REPRODUCED OR RE-POSTED, IN WHOLE OR IN PART, BY THE RECIPIENT WITHOUT EXPRESS PERMISSION FROM 451 RESEARCH.



Introduction

Direct liquid cooling (DLC), a collection of techniques that deliver a coolant fluid to heat sources, as opposed to relying on air as a heat-transfer medium, has been a long time coming to datacenters, but has yet to become the common choice. This is about to change, we believe, and our research coverage of liquid cooling will grow in 2019 to reflect this trend. In this report, we lay out the major factors that we believe will drive the shift.

451 TAKE

For many years, we've held a positive outlook on liquid cooling for its superior efficiency and tracked the evolution of suppliers, yet precious little happened outside the domain of high-performance computing. It would be all too easy to invoke the datacenter sector's notoriety for its resistance to change, but it would not be fair. Most operators had no strong incentives to jump into such a fundamental change to their infrastructure – one that also required the IT to change with it, and vice versa. Additionally, DLC suppliers tried more often than not to sell products that simply did not fit the bill for installations at scale, and used messages, such as extreme rack density, that did not resonate. We see this all changing, and a new constellation of technical and business factors appears to align in favor of liquid cooling.

Why it should be different this time

Direct liquid cooling in datacenters – a broad variety of techniques where a coolant fluid is delivered to the processor and other electronics for heat transfer – has been around for too many years to carry any sense of novelty. 451 Research has been covering the area for longer than we can remember. On technical merit, it should have always been the dominant form of cooling, rather than air, which is much less efficient at heat transfer, meaning lots of it needs to be circulated and cooled in order to control temperature – not to mention all the problems with humidity in the air, which is possibly an even bigger factor behind premature failures of some IT components, such as storage drives, than elevated temperature is.

Contrary to clear benefits, DLC is still not the standard in datacenters. DLC is always happening 'next year' to deliver lower costs, add efficiency and (ultimately) give back control over the technology stack – much like Linux desktops. Inertia in the datacenter ecosystem around air cooling means that, without solid reasons, customers and their suppliers will not invest time and money in a major change with an unclear return. Equally, air cooling systems have become much more efficient at their job over the last decade and allowed operators to cut build and operational costs. Unlike Linux desktops, however (and that's where the analogy breaks down), DLC can deliver tangible business benefits, as opposed to just reshuffling the deck to move the bottleneck elsewhere.

Super-efficient air-cooled datacenters still use considerable energy to power all the facility and IT system fans that move air even if the air is not cooled at all – and most datacenters have no chance of being anywhere close to the efficiency levels of the best sites. Energy consumption is not the entire efficiency story: massive amounts of power capacity are stranded to act as a reserve for peak cooling requirements, when the fans, pumps and compressors work the hardest. This can be as much as one-third of a site's power envelope. This is power capacity that could be reallocated for IT – a great boon to operators in locations where access to more power is either slow or expensive, and these days many multi-tenant operators and service providers want lots of new power, and fast.

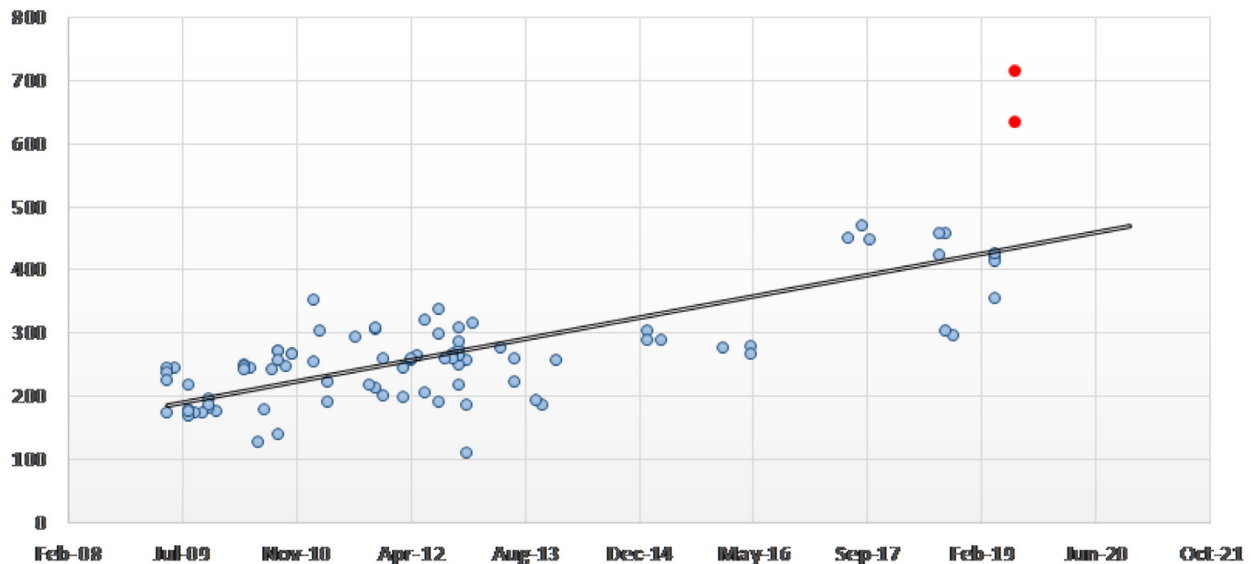
REPORT REPRINT

None of these advantages is a recent development, technically speaking. What's changed is the technical and business environment in and around the datacenter. Air cooling will hardly get much more efficient than it is today – datacenter designs are highly optimized and products are well honed. A hard limitation is supply air temperature: most IT tenants would not accept the wide temperature bands (15-32 C or 59-89.6 F) that would be needed in many locations to eliminate much of the air cooling infrastructure and related energy, regardless of savings. DLC offers an answer to the 'where to go next?' question on cost-performance optimization.

All the while, server processors have become more power-hungry with the growing density of microelectronics. At the start of the decade, mainstream server processors were rated under 100W for thermal design power (engineering speak to specify what cooling needs to handle at full load), while dual-processor servers consumed in the region of 200W when highly loaded, as witnessed by the power efficiency database (spec_powersj2008) of the Standard Performance Evaluation Corporation, a nonprofit industry body for server benchmarking.

Figure 1: Server power at full load – dual processor, 1U

Source: Standard Performance Evaluation Corporation



Fast forward to 2017, and Intel's mainstream server chips, which are in over 90% of server processors in datacenters, broke past the 200W barrier, bringing total server consumption close to 500W. The next doubling happened in April, less than two years later – Intel's latest-generation server platform supports processors up to 400W each (estimated server power consumption highlighted in red on the chart). AMD is also widely anticipated to launch high-power server products, while gaming GPU-maker-cum-AI-accelerator NVIDIA has been selling 300W parts into datacenters for years now, too.

At such thermal density (the surface area of processor packages is a few square inches at most), air cooling is still feasible, but becomes costly and impractical. And because recent processors dynamically maximize their speed within temperature and power limits, more cooling capacity translates into more performance, something that air cooling cannot match. There is always some lag between the introduction of new chips and them finding their way into datacenters in volume, and it's also true that not all servers will use the highest-power models. However, the direction of travel is clear: server power is on a steep rise.

451 Research views silicon power and performance density as a primary factor in making a shift to DLC, not rack density, which will be a secondary cost-reduction benefit to operators that decide to make their next-generation infrastructure denser. Even then, 451 Research does not expect operators that adopt DLC to go anywhere near supercomputing-like densities. This is diametrically opposite to the dynamics seen with extremely dense compute clusters of supercomputing (often 10x denser than a typical datacenter), compressed because of the need for high-speed data sharing at scale.

The year liquid cooling becomes a solid choice?

Adding it all up, DLC promises datacenter operators a path to improve the cost-performance profile for both the facility and IT infrastructure, including having the ability to run more IT capacity on any given site. It also offers a potentially lower-cost option to refurbish outdated facilities for next-generation IT.

Even so, mainstream datacenter operators will need DLC systems that work for them with minimal disruption to their design, build and operations regimes. DLC designs mandating that everything (including entire supply chains) has to change at once will, by nature, run into more resistance than those that offer the ability to integrate into existing practices and relationships. This has been a learning experience for makers of DLC systems, which all too often created products that were manifestations of an ideal – interesting for those few with a flair for something radical, but unpalatable for most.

451 Research sees promising changes in the approaches of DLC vendors – established and new entrants alike. Products are not simply improved based on customer feedback, and increasingly more customization to specific requirements is becoming possible. In some cases, this goes well beyond factory configuration options – full customization around a core DLC component is possible. With hyperscale and tier two service providers having gained a commanding weight in engineering decisions for datacenter infrastructure – not just for their own, but for colocation and wholesale facilities – the ability to customize DLC systems with speed and at cost will be crucial for success.

That this report hasn't discussed the various types of DLC on offer is deliberate. The key factors are the step change in heat transfer efficiency and the need for more integration and co-engineering between facilities and IT – it's a change of method. Everything else is an engineering decision made downstream. Just like air cooling is a collection of vastly different technologies, from chilled-water systems to evaporative and adiabatic economizers, liquid cooling is not a specific area of engineering. Typically, 451 Research subsegments DLC into cold plates and immersion systems as a classical take on classification to differentiate between partially liquid cooled systems from totally liquid ones, but the merit of this differentiation may erode as engineering efforts converge.

Finally, there are some pressing questions that remain unanswered. Even though 451 Research views DLC as ultimately inevitable, and we are aware of multiple proof-of-concept efforts at some major operators, the timing of any major shift remains highly uncertain – will mainstream take-up start in six months, two years or more than three? Who is going to pull the trigger first? Will operators go all-in on DLC or just use it for specific workloads in the mix? How this will affect the revenues and margins of datacenter cooling makers is not trivial either. If DLC goes big, most of the chillers, cooling towers, air conditioners and air handlers, and related power infrastructure will go away. Operator desire to cut cost is inescapable, and equipment makers will need to find ways to offset compression in their air cooling businesses.

In subsequent reports, we will explore these questions alongside market developments and revisit the vendor landscape, which has changed considerably over recent years.