

A decade of HPC in oil

Jean-Yves Blanc^{1*} and Laurent Clerc¹ give an overview of oil cooling systems for data processing centres.

Introduction

CGG has always been at the forefront of industrial High Performance Computing (HPC) architectures: we were operating vector supercomputers (Convex, Cray and NEC) in the early 1990s, and large parallel supercomputers (Convex SPP, IBM SP, Sgi Origin) by the end of that decade. At the turn of the millennium, we were pioneering the use of commodity clusters, and started to add accelerators a couple of years later, even before GPGPU programming languages formally emerged.

Our oil story started later, one day in mid-November 2009, at the SC09 SuperComputing conference in Portland. In the ‘New Emerging Tech Corner’ of the exhibition hall, a small

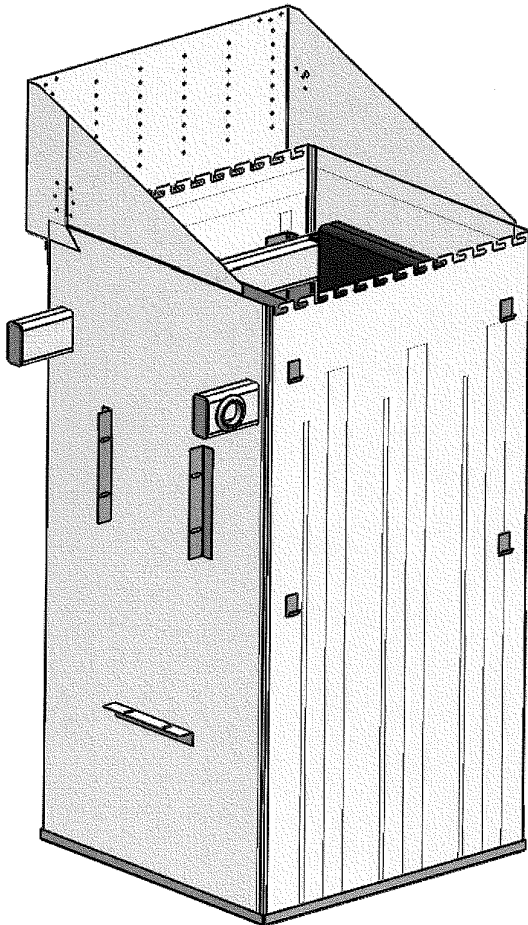


Figure 1 Drawing of the prototype small 15U tank, used for initial testing.

Austin-based startup had installed a tank full of liquid and servers. The company founder told us that their aim was to start manufacturing containers that would hold mineral oil and cool servers fully submerged in liquid. Mineral oil had the advantages of being nontoxic, inexpensive and easy to handle, would not conduct electricity, and is significantly better than air at cooling the components. Once we got over our instinctive reaction that putting electrical equipment in a liquid was not a good idea, we began to realize that this could actually be a very efficient idea for the kind of heat densities that multi-GPU nodes release.

Then we had to pitch the idea to our management, staff, hardware manufacturers, leasing companies, insurers and lawyers, etc. None of that was easy, but this was not the first time we were adopting cutting-edge technology, so we applied the same process we had used to adopt other products before and since: build a proof of concept (see Figure 1), then an industrial prototype, before deploying the product. This approach is effective at de-risking the solution and raising the confidence of everyone involved. In this case, the process worked well, rapidly demonstrating the workability of this idea. Of course, there were some issues, mostly related to all the moving parts, such as fans and spinning disks (or floating labels), but we moved to solid-state drive (SSD) and removed everything else that was a problem, ending up with a more elegant and functional solution.

Dipping our toes in oil

As we worked through our validation process, new algorithms began to emerge for processing data from the Gulf of Mexico, which required a significant increase in processing capacity. The challenge was to fit much more, and hotter, hardware into the same space for, ideally, the same cost (a recurring theme in our business). That situation, and the fact that the oil immersion technology provider was in Austin, convinced us that our Houston hub was the right place to try out the technology at scale.

Our initial cost model showed us that we needed to exceed 25 kW in 32 U (~800W/U, 1U being 1.7-inch vertical space in a standard rack) for the solution to become cost-effective over a two-year period. The idea was to compensate for the cost of the installation through the savings we would make in power, due to the reduced effort needed to move a lot less fluid with less temperature excursions and at greatly reduced velocities, as well as in space because of the higher system densities we could achieve. Since we did not know if oil immersion would be a

¹ CGG

* Corresponding author, E-mail: jean-yves.blanc@cgg.com

DOI: 10.3997/1365-2397.fb2021092

long-term solution for us, we wanted to break even over a short period of time to limit risks. Of course, we ended up keeping the tanks for more than 10 years and counting, eventually making this an outstanding proposition for our company. But we could not anticipate that then.

To achieve 800W/U, we needed GPU nodes, since CPU nodes did not generate that much heat at the time (see Figure 2). To start with, we acquired a small 15U tank, and started dipping servers in it to build our proof-of-concept system. We used old servers at first, to understand the interactions between the oil

and computer equipment, and then newer equipment, closer to the target configuration. This is how we learnt not to put spinning disks in oil (most have a little hole to balance pressure... so we used SSD instead). Incidentally, that also gave us the time to understand interactions between the oil and everything else: labels, cables (wicking), floor, people, clothes, etc. And after more than 10 years, the fact remains that oil is, well, oily: but that really is the only drawback of the solution. Everything else has worked as expected. To handle the oil, we developed processes to keep everything, computer rooms and people clean, by using absorbing mats, paper towels, gloves, draining boards, as and where needed. On the plus side, we quickly realized that liquid-cooled environments do not smell like a fast-food joint (unlike some early comments we received!) and are extremely quiet. There is almost no noise, even to this day when we cool 50 kW per tank (1.2kW/U). Readers familiar with air-cooled computer rooms know how 30kW+ racks sound: more like a plane taking off, and probably using as much energy too. This is not at all the case with oil. In fact, the only reason we even have air-conditioning in the computer rooms is so that our technicians can operate in a healthy environment. The tanks could (and probably should) reside in a warehouse, a container, or even outdoors when properly sealed.

Next came the first production installation: a set of four tanks, each roughly the equivalent of a 40U rack, sharing a common oil/water heat exchanger and pump module (see Figure 3). This happened to be the configuration proposed by the supplier we worked with, and it has performed well all these years. This system gave us the opportunity to design our industrial prototype. The difference between the proof of concept and the industrial prototype is that, while the proof of concept tells us if a solution works, the industrial prototype tells us how to operate it at scale. Industrial prototypes are very important in our business but often overlooked. For us, this was the important stage where we learnt how to install the systems in a computer room, connect them to existing infrastructure, and keep them operating safely at peak efficiency. Having a higher power density has an impact on power and cooling distribution, as



Figure 2 High-density GPU servers submerged in oil.

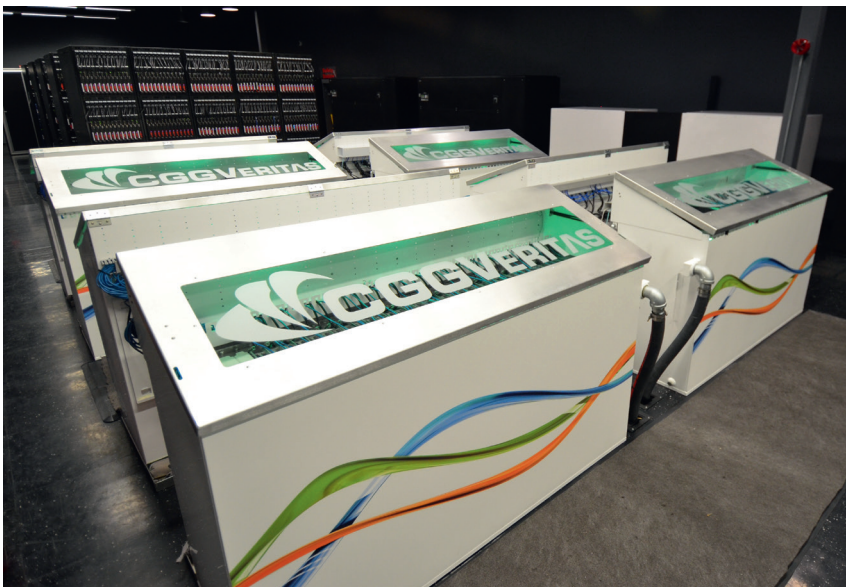


Figure 3 Two groups of quad tanks in their operational setup.

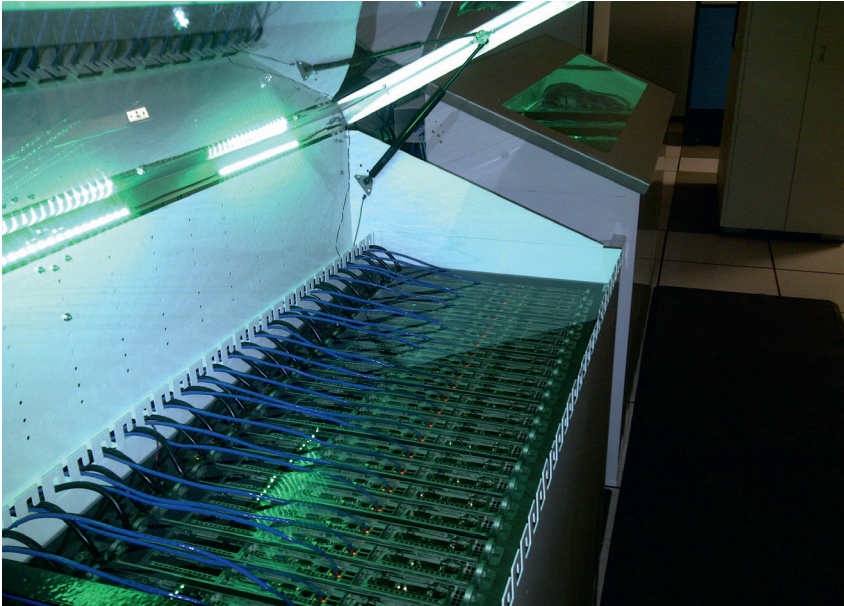


Figure 4 A row of high-density GPU servers in oil.

well as on how much time you have to deal with loss of cooling and other incidents: all the things you really need to know about in a large-scale industrial setting.

The industrial prototype quickly gained popularity with our users, and we started planning for a computer room dedicated to these systems. We deliberately chose the oldest room we had because its power usage effectiveness (PUE) was the worst in the building. It happened to be designed for air cooling, with a false floor and air handlers, so we left that equipment as it was since we did not need any of that infrastructure. The PUE was significantly reduced, so more and more quad systems were deployed, until the room was completely full (see Figure 4). We then started overflowing to the adjacent room where our operators used to be, and soon we were able to expand even more. Along the way, we learnt that oil is actually a very good environment for electronic equipment but more about that later.

While CGG started alone along this path, we have tried to get other parties interested in using or supporting this technology: this is because, although we see oil immersion as a significant technological differentiator, we also recognize that there needs to be an established ecosystem for the technology to thrive, and this requires a customer base. Interestingly, no matter what we shared, explained or demonstrated, only a handful of suppliers and very few potential customers took the plunge.

Nevertheless, we want to acknowledge the few manufacturers who agreed to maintain a warranty for their equipment after we shared failure rate statistics that, if anything, were better in oil than air. Along the same lines, a couple of manufacturers agreed to create oil-ready configurations, with (minimal) modifications to remove unused components (e.g. small high-speed fans), upside-down designs (so that all connections and supports are on the same side), and freer-flowing radiators and chassis.

Why bother, and what did we learn?

Our purpose with oil immersion cooling has always been to enable simple, stable, safe and efficient long-term operations of denser and higher heat density HPC systems.

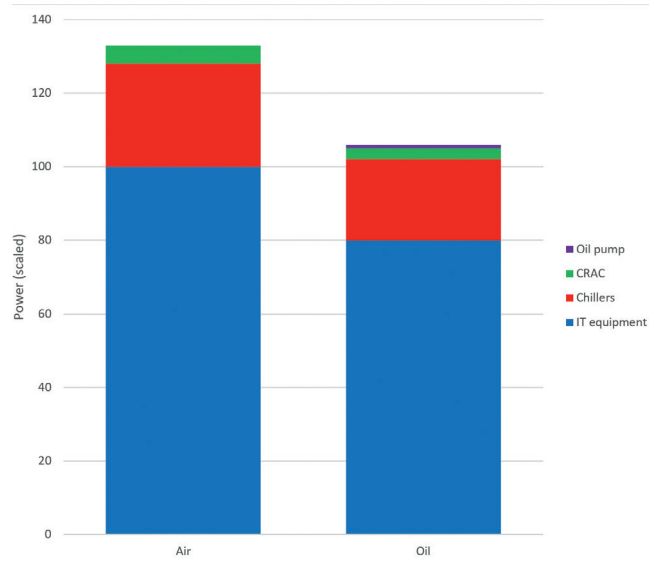


Figure 5 Power comparison between air and oil-immersion cooling.

Unlike conventional data centres, HPC systems use considerable power and dissipate a lot of heat because they require a large number of very tightly integrated systems using high-power components. As many HPC users have discovered by now, air cooling cannot get much past 30-35 kW/rack, because at that level, air flow and ΔT are maxed out and PUE deteriorates rapidly. Oil immersion is therefore an attractive option for reducing the environmental impact of data centres, as it ensures most of the energy is used by the IT equipment (see Figure 5). In other words; improved data centre cooling efficiency has a significant benefit from an environmental perspective.

To enable higher-power densities and improve PUE, we decided that it was worth looking at another way to cool our systems when we moved to GPUs, more than ten years ago. Today, this trend is accelerating: we started with 2 GPU per 1U node, and are now routinely using 6 GPU/1U and testing 8 GPU/1U, with GPU cards that are using more and more power. Our target

is > 4 kW/1U in dense environments (call it 150 kW per rack equivalent). This is way beyond what even the most optimistic proponents of air cooling would consider, and probably closer to the limits of what even a liquid can do.

In addition, and what is perhaps unique to CGG, we serve a business with an obligation to deliver the best price/performance ratio possible. We need to run 24/7 at full capacity for years and must deliver stability with systems that cannot be

over-engineered and should be built with commodity components wherever possible. We need robust, massively efficient, simple and practical solutions, not delicate works of art.

As an engineering strategy, we prefer to remove things rather than add them to solve a problem. Oil immersion does exactly that: there are fewer parts in the nodes, fewer parts in the cooling system, less fluid velocity and temperature differences or fluctuations, no humidity or static electricity considerations. Any other cooling technique adds components, cost, complexity and opportunities for failures.

As for the efficiency, the cooling effort of full immersion is reduced (leading to better PUE) because a liquid is significantly better than a gas at transporting heat: that means less fluid velocity and volume, less mechanical effort to circulate the fluid, less temperature differences between fluid input and output. The fact that a relatively small pump can replace hundreds of fans in nodes, racks and CRACs improves power usage, and the fact that we can run with a smaller ΔT allows us to use warmer water as an input, further reducing the cooling effort.

Obviously, we design, own and operate our data centres, define and install our O/S and middleware, write and optimize our own code, and control our jobs and priorities. That level of control gives us significant freedom to achieve a very deep optimization of our systems, up to and including our computer room environments. Not everyone is in that situation. But we were, and so we took advantage of what we saw as a better way to operate.

So, after ten years, was it really better?

First, there were issues, that broadly fall into two categories: chemical interactions and operational issues.

Chemical interactions usually occur relatively quickly and can be caught in the POC or industrial prototype. We mainly had issues with the following products melting away:

- Thermal paste, greases
- Glues and varnishes holding labels or components (lots of bits floating around)

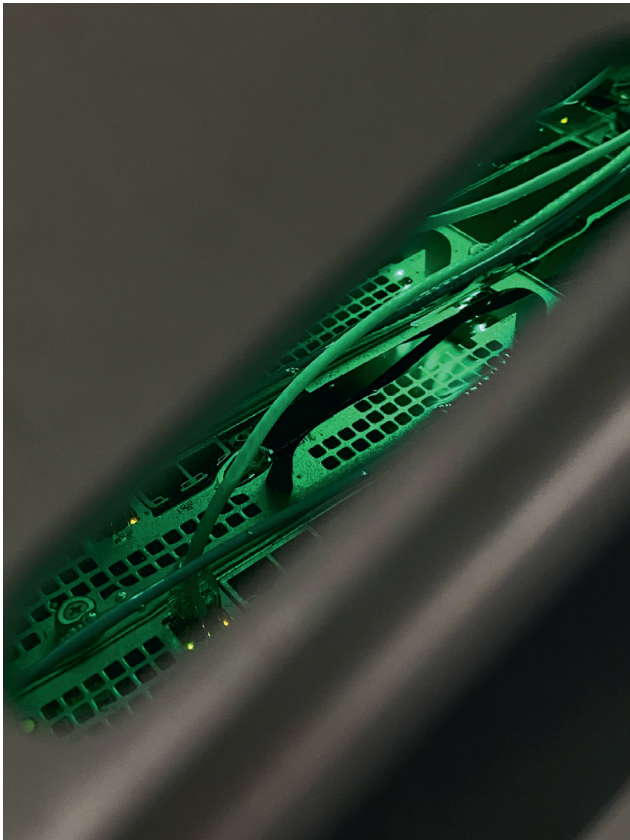


Figure 6 Copper ethernet connectivity to mitigate failure of optical cables in oil.



Figure 7 Pipes transporting hot oil to the chilled water heat exchanger.

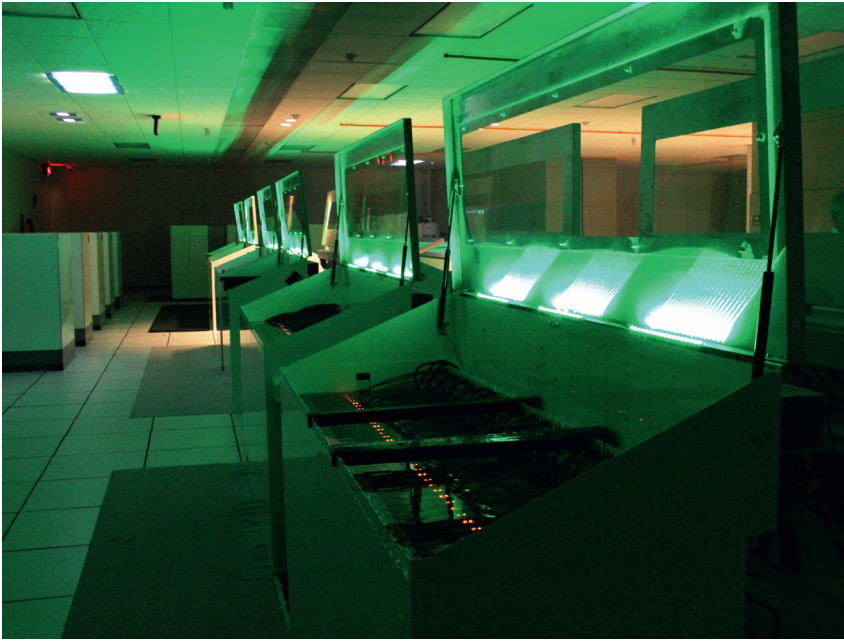


Figure 8 Computer room populated with oil immersion tanks.

- An additive used in some plastics to make them flexible (plastic gets brittle with time)
- Some network optical interfaces (eventually affecting network operations), see Figure 6
- Some electrolytic capacitor sealants (oil mixes with the electrolyte, capacitors swell, etc.)

The most notable operational issues were:

- Oil wicking out along the cables (network, power), even against gravity...
- Micro-codes complaining that fans, etc., were missing
- Tiny pieces of equipment falling off and accumulating at the bottom of the tanks
- Difficulties removing dripping and increasingly heavy servers from the tanks
- Significant testing was needed to find a way to clean parts and components

But all of those situations had fairly practical solutions, and overall, it is possible to operate at scale in oil without significant issues.

With time, we could confirm that oil was not an issue for the systems. Failure rates are similar or better than air systems, and to date we have not encountered any systemic failure due to the systems being in oil over ten years of operations. This has taken place across multiple generations of systems and components, with some systems having been in use for years.

In addition, the thermal environment of the servers is significantly better: temperatures are much more stable and uniform across all components, fluctuations are slow and limited (no thermal shock), even at more than 1.2kW / 1U, and even when the primary cooling loop (see Figure 7) is lost for a limited period. Considering that thermal excursions have a cumulative negative effect on the life expectancy of the components, this is a very positive attribute.

But ultimately, as we now look for an industrial, simple and efficient solution to the very high heat densities we can see

coming, we definitely believe that the solution our company has used for more than a decade (see Figure 8) is going to be the only valid path forward.

There is a learning curve, and we already have a ten-year head start, but we are open to sharing our experience.

Slow industry evolution

Other cooling options are reaching their limitations, both technically and from a cost perspective (which is obviously an important criterion for an industrial player). Air-cooled high-density HPC systems, whether they are classic 19-in racks or OCP 21-in ones, are becoming expensive to cool at scale above 30 kW (per 42U rack), even with some form of confinement (hot/cold aisle being the standard option, that can be made to be quite sophisticated and more efficient, but at a cost). Rear doors are an improvement to the air cooling strategy, as they reduce (quite significantly) the temperature of the hot air but there is still a need to ingest the proper volume of air through the system to evacuate the heat, so the doors increase the heat dissipation limit, but again the cost of a cold door increases rapidly with their cooling capacity. Direct liquid cooling is an interesting recent trend, even if it has existed for quite a few years, but it has two major limitations compared to oil immersion: only dedicated components featuring a specific heatsink are liquid-cooled, and the unitary cost is not competitive compared to immersion cooling, especially in the very dense systems range. And the last option currently used in HPC systems is phase-change cooling. From a thermodynamic perspective only, phase change (in most cases, liquid to gas) has a larger capacity to evacuate heat compared to a flow of liquid. Phase-change cooling systems have been around for several decades (starting with Cray Research systems and SGI supercomputers) but remain very marginally used, they are difficult to operate (due to their physical properties and the nature of their fluid, which is very often toxic), and are fairly expensive. The last two options also require bringing liquid into computer rooms.

From an oil immersion ecosystem perspective, it is now much richer than a decade ago: we now have about a dozen oil immersion cooling system vendors, and many HPC users and hyperscalers have adopted oil immersion. The entire oil immersion cooling commercial offering is proposed by small businesses, as no major IT manufacturer has started to offer such a cooling option for its computer equipment (although some of them have been studying this technology for several years). It is difficult to evaluate the effective scale of this market, as several major HPC users and hyperscalers are not keen to communicate, but it is certain that oil immersion is currently used on the equivalent of several tens of thousands of racks. Some HPC sites have publicly communicated on these deployments, defence agencies of various countries 'might have adopted' this approach and it is likely that some of the internet giants use oil immersion cooling at various scales.

The ecosystem of server vendors, a key driver behind a broad adoption of oil immersion cooling, is currently growing. As explained above, we managed to convince one, then several of our vendors to support oil immersion for some of their models. Even if it has taken several years, we are now able to choose from a variety of architectures that are ready for oil immersion. Not only have those vendors agreed to maintain a warranty for oil-immersed servers, but also some design modifications have been made (BIOS modifications, the positioning of some components, frame geometry alterations). The next step would be for server manufacturers to identify a big enough market for them to release and manufacture architectures specifically designed for oil immersion. CGG is currently actively collaborating on oil immersion projects, not only with server vendors, but also with some key component providers, to prepare for the release of even higher-density HPC systems.

Conclusions

CGG has gained unique operational experience in oil immersion cooling. Experience is key with this technology, as some of the pitfalls encountered in its deployment can only be identified at significant scale and over time (for instance, it is quite difficult to simulate accurately enough the ageing of an electronic component dipped in mineral oil).

After a long maturation, the adoption of this technology is accelerating: we already discussed the increased density of HPC systems and the severe increase in the heat dissipation of the latest generations of components, and the fact that other available cooling options are reaching their limits, especially for industrial deployment, with cost and reliability constraints. Another aspect that currently argues in favor of oil immersion deployment is its energy efficiency, which is of course a very positive impact from a cost standpoint, but it also positions this technology at the top of green IT initiatives in HPC (energy efficiency being the main adjustment parameter in this respect).

Oil immersion seems increasingly required to enable the coming generation of ultra-dense and hot HPC systems. All the major CPU manufacturers are planning to release their new generation of chips with ultra-fast memory and a large core count (with expanded vector calculation capabilities) and every new generation of GPU is significantly hotter than the previous one (especially in the last few years, when GPU manufacturing lithography has reached state-of-the-art technology, and the recent transition from GDDR to HBM memory has further accelerated the growth in memory bandwidth). It looks increasingly like the HPC roadmap in coming years will require the ability to cool systems with such a high heat density. Oil immersion cooling could offer that capability.