



HAL
open science

Datacenters bas carbone : le recours au Liquid Cooling inévitable à terme ?

Jean-Christophe Léonard

► **To cite this version:**

Jean-Christophe Léonard. Datacenters bas carbone : le recours au Liquid Cooling inévitable à terme ?. CVC la revue des climaticiens, 2023. hal-04380456

HAL Id: hal-04380456

<https://edf.hal.science/hal-04380456v1>

Submitted on 8 Jan 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Datacenters bas carbone : le recours au Liquid Cooling inévitable à terme ?

Auteur : Jean-Christophe Léonard. Ingénieur. EDF Lab Les Renardières - Technologies et Recherches en Efficacité Énergétique

Le numérique nous rend d'innombrables services, mais son empreinte carbone, actuelle et future, interroge nos sociétés. Cœur de cette industrie, les datacenters sont donc sommés de viser une efficacité énergétique élevée. Ils y sont d'ailleurs contraints par la hausse du prix de l'électricité. Ce secteur a su relever les défis : alors qu'en 2007, le Power Usage Effectiveness* (PUE) du parc Monde était de 2.5, sept ans plus tard, en 2014, il était de 1.65, soit une baisse de 30%. Cependant, en 2021, il stagne à 1.57. En France, la surface du parc avoisine les 900.000 m², avec un PUE de 1.69 (source ADEME & ACERP 2022). La consommation totale du parc est de 11.6 TWh/an : c'est l'équivalent de la production annuelle d'un EPR. Les émissions directes des datacenters sont voisines de 750 ktonnes éqCO₂/an. Alors, certes, les bonnes pratiques doivent continuer à se diffuser (hot air extraction, free-cooling, etc.). Mais, pour s'approcher du Graal -PUE de 1-, l'utilisation de technologies de refroidissement en rupture sera nécessaire.

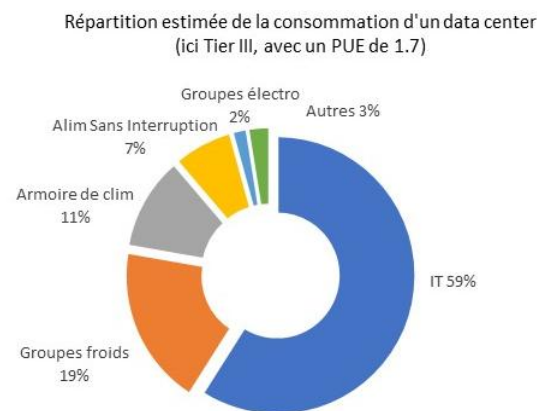


Figure 1 : exemple décomposition des consommations d'un datacenter

* : PUE, consommation totale du datacenter divisée par la consommation des serveurs (IT sur le graphique)

Une performance exigée par tous : Europe, nation, territoires, sociétés et citoyens

Le lobby des datacenters Européens a pris une initiative d'autorégulation : le Climate Neutral Data Center Pact, afin que le législateur s'en inspire pour le Green Deal. Il comporte quatre volets. En 2030, afficher un PUE inférieur à 1.3 ou 1.4, selon la localisation. Ensuite, disposer d'un approvisionnement 100% ENR (via les garanties d'origine). Troisièmement, favoriser la réutilisation des serveurs. Enfin, valoriser l'énergie dissipée, si c'est « économiquement raisonnable ». En France, le secteur est soumis au décret tertiaire qui imposera des PUE maximums (attendus courant 2023). Les contraintes sont aussi parfois locales, à l'exemple d'Amsterdam qui conditionne l'obtention du permis de construire d'un datacenter à la mise en place d'une récupération d'énergie. Le signal est encore (très) faible, mais des entreprises mentionnent dans leur rapport RSE, le recours à des datacenter à faible PUE. Cas plus rare, parfois des collectifs s'opposent même à l'implantation d'un datacenter (exemple à Wissous).

Quand le refroidissement à air atteint ses limites...

En parallèle des exigences citées précédemment, cette industrie fait aussi face à des contraintes endogènes. La performance des serveurs croît, année après année. Cependant, depuis 2005, un ralentissement est observé. Alors que la loi de Moore est toujours vérifiée (jusqu'à quand ?), la contrainte vient de la limitation en fréquence des microprocesseurs. Des limites physiques, complexes à développer ici, ont été atteintes (Dennard scaling). Face à cela, les industriels ont su trouver une réponse : les micro-processeurs multicœurs sont apparus, avec une augmentation « raisonnable » de la puissance dissipée. Malgré cela, la performance des serveurs a continué à se tasser...Innovateurs infatigables, les industriels ont alors depuis exploré une nouvelle piste : la spécialisation des puces. Comprendons que les CPU (Central Processing Unit) sont capables de couvrir, certes, de nombreuses applications (bureautique, vidéos, calculs, etc.), mais avec une performance moyenne ; d'où l'idée de développer des produits très performants, mais sur un champ applicatif plus restreint. Les GPU (G : Graphics), particulièrement adaptés au traitement des flux vidéo, en sont un exemple. Cette spécialisation touche maintenant le

domaine du calcul intensif : l'High Performance Computing (HPC). Au sein du bestiaire HPC, on trouve par exemple, les TPU (T : Tensor) dédiés à l'utilisation de TensorFlow pour l'intelligence artificielle, les DPU (D : Data), les MPPA (Multi Parallel Processor Array), etc. Hyper performants, mais sur un domaine très ciblé, ces objets ont deux points en commun : ils se diffusent rapidement eu égard aux services qu'ils rendent, et ils sont surtout de plus en plus difficiles à refroidir... L'ASHRAE, Intel, Dell et Equinix alertent justement sur le fait que nous rentrons dans l'ère du « Power War Trend » (fig.1). Intel annonce d'ailleurs dans sa roadmap qu'au-delà de 2025, ses produits les plus performants ne seront plus compatibles avec un refroidissement à air (fig.2).

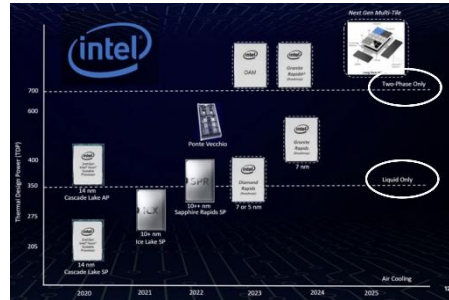
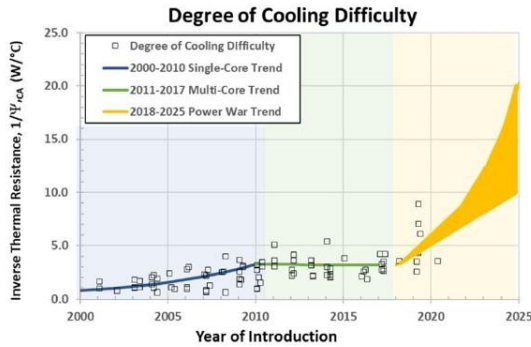


Figure 2 : l'ère du « Power War Trend ». Figure 3 : le refroidissement à l'air des processeurs atteint ses limites

C'est là qu'apparaît le concept de Liquid Cooling : multiplier d'un facteur 100 à 1000, le coefficient d'échange convectif à la surface des puces. Ce terme générique recouvre deux technologies (fig.3) : celle dite Direct (to Chip) Liquid Cooling (DLC) et la seconde, plus en rupture, l'Immersion Cooling (IC). Dans le DCL simple phase (DCL1P), un échangeur de chaleur est disposé sur le processeur pour assurer une conduction parfaite (collé et vissé). Le caloporteur (généralement de l'eau) extrait l'énergie des processeurs vers des dry-coolers (température voisine de 55°C). Il n'y a donc plus besoin de produire, transporter le froid correspondant à cet apport. Il reste, cependant environ 15 % de la puissance à évacuer (effet Joule sur les autres composants). Le serveur conserve donc son ventilateur, tout comme le datacenter ses armoires de climatisation et une production frigorifique pour traiter ces apports résiduels. En DCL1P, la baisse de consommations est voisine de 30%. Une variante existe avec un caloporteur qui s'évapore au contact du processeur (DLC2P).

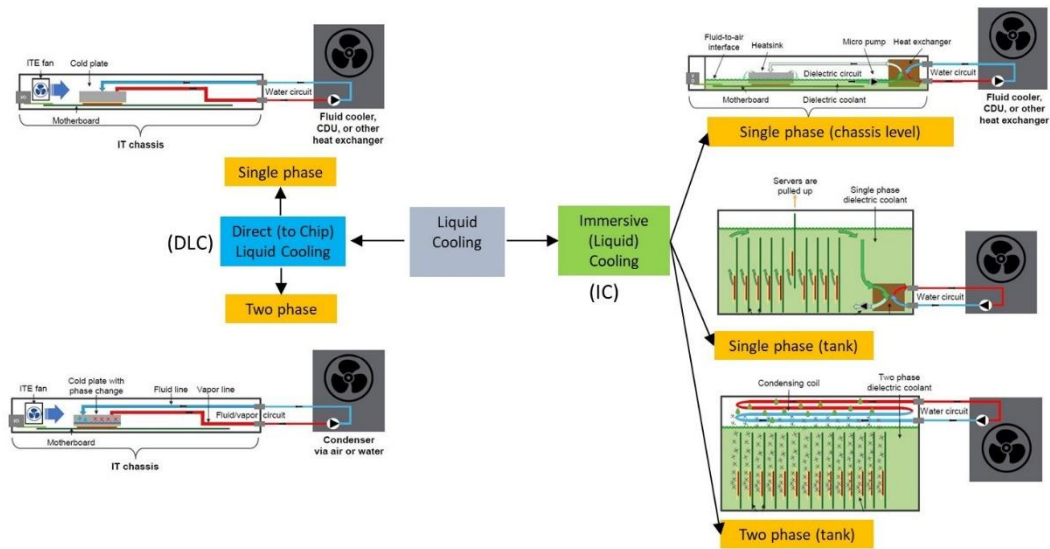


Figure 4 : Liquid Cooling : un terme générique qui recouvre différentes technologies

Avec l'immersion Cooling une phase (IC1P), les serveurs sont immergés verticalement dans un fluide diélectrique. Ils fonctionnent sans ventilateurs (10% de la consommation d'un serveur). Le datacenter n'a plus ici besoin de disposer d'armoires de climatisation, d'eau glacée, etc. La chaleur est envoyée via des pompes vers les dry-coolers. Certains constructeurs recourent à des pompes à l'intérieur de la cuve pour la circulation du fluide. D'autres, réussissent à travailler en convection naturelle. Dans une variante dite chassis level, les serveurs sont en position horizontale et baignent dans une couche de diélectrique. En IC1P, la baisse de consommations est voisine de 35%. Enfin, pour la très forte densité (250 kW/bac), le recours à l'immersion double phase (IC2P) est nécessaire. L'évacuation de chaleur en dehors du bac s'effectue via un condenseur.

Une compatibilité des serveurs à l'immersion une phase aujourd'hui acquise

Les fluides diélectriques utilisés dans l'IC1P peuvent être soit des huiles (minérales ou de synthèse) soit des produits fluorés : PFC (PerFluroCarbures), HFE (HydroFluorEther) ou plus rarement des FK (Fluorokétone). Si l'on compare ces produits, notamment en termes de GWP (Global Warming Potential), de risques sanitaires pour les travailleurs de l'IT liés à une exposition prolongée aux vapeurs, les produits fluorés seront raisonnablement éliminés. In fine, les huiles de synthèse semblent les meilleures candidates, notamment au regard de leurs qualités diélectriques.

La filière a mené de nombreux tests (mécaniques, électriques, analyse microscopique, etc.), les dernières années, sur la compatibilité de ses produits avec les composants des serveurs. Lors des premiers essais sur l'IC, de nombreux défauts ont été observés : décomposition des câbles, des colles, des encres, gonflement des condensateurs électrolytiques, etc. Mais, des solutions ont été trouvées. En 2022, on voit d'ailleurs apparaître les premiers serveurs garantis pour l'IC1P. L'IC2P rencontre encore cependant des problèmes de mise au point. Microsoft Azure, qui y consacre d'importants travaux R&D, le reconnaît d'ailleurs, non sans humour : « we're not ready to dive in ».

Des PUE inférieurs à 1.1

Actuellement, les entreprises testent à petite échelle ces technologies. Au-delà de la seule question de l'efficacité énergétique, il s'agit de comprendre, en situation, les implications opérationnelles de ces technologies. Citons par exemple, l'utilisation de potences pour retirer en toute sécurité un serveur d'un bac, les procédures pour détremper un serveur avant intervention, l'efficacité des connexions rapides sans gouttes (dripless quick connect), etc.

Pour l'instant, les retours d'expérience chiffrés sont peu nombreux. Nous ne citerons pas d'entreprises ici, mais que ce soit en Chine, au US et en Europe, la plupart des PUE annuels obtenus sont inférieurs à 1.1 (1.04, 1.07, etc.). C'est remarquable. Notons qu'il s'agit des PUE calculés au périmètre IT + auxiliaires d'évacuation de la chaleur. On parle alors de partial PUE (pPUE). Les pertes amont (alimentation sans interruption, batteries, fioul pour les tests des groupes électrogènes, etc.) n'étant pas intégrées (rajouter de 5 à 10%). A titre d'exemple, la figure 4 compare les puissances mises en jeu selon le mode de refroidissement sur une expérimentation menée par un opérateur de commerce en ligne. On notera la disparition de la consommation des ventilateurs des serveurs dans le cas immergé. La baisse de consommations est ici de 36%. Des résultats similaires ont été obtenus par une grande banque française avec la même technologie, avec un maintien des performances lors des récentes canicules.

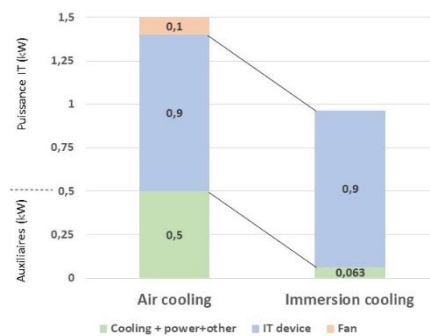


Figure 5 : résultats obtenus par une immersion une phase

Une récupération de chaleur facilitée, donc à développer

Le Liquid Cooling est aussi un facilitateur de solutions de récupération d'énergie. Nous l'avons dit, ces technologies rejettent de l'eau aux alentours de 55°C. Certes, les datacenters sont généralement éloignés des points de valorisation possibles. Mais, dans une configuration de petits datacenters (ordre de grandeur < 100 kW_{IT}) territoriaux (Edge Computing), ceux-ci peuvent être installés justement au plus proche d'un usage pour valoriser leur énergie fatale. Les exemples se multiplient en France et en Europe. Citons par exemple, Qalway qui préchauffe l'ECS en résidentiel collectif à La Chapelle sur Erdre (45 logements), qui alimente aussi un réseau de chaleur basse température à Kankaanpää en Finlande. Autre exemple, Neutral IT qui équipe la piscine de la Butte-aux-Cailles à Paris. Dans chacun de ces cas, l'énergie récupérée vient substituer des consommations fossiles, contribuant ainsi à nos efforts nécessaires de décarbonation. Dans le cas de la piscine citée, la baisse des émissions directes de CO₂ sur le périmètre [datacenter + fraction des besoins de la piscine couverts par la récupération] avoisine les 75%.



Figure 6 : de gauche à droite : préchauffage d'ECS en résidentiel, en piscine, ou sur des réseaux de chaleur basse température

Une filière française

Point remarquable : une filière française peut se structurer autour du Liquid Cooling : des puces MPPA fabriquées par le français Kalray, aux serveurs compatibles à l'immersion de la société Strasbourgeoise 2CRSI, en passant par les cuves particulièrement adaptées à la récupération d'énergie du groupe Lavallois Numains, tout comme les solutions Qalway ou Neutral IT, à des hébergeurs comme OVH, nous disposons du tissu industriel pour développer cette filière. L'idée ici n'est pas d'invoquer le concept à la mode de la souveraineté numérique. Soyons réalistes, en pratique la souveraineté consiste généralement à choisir...ses dépendances. Certes, nous avons déjà la chance de disposer, en France, d'une production d'électricité décarbonée. Mais soyons conscients, que sur ce cas d'usage précis – le datacenter bas carbone - nous avons aussi collectivement, grands groupes, entreprises et collectivités, la responsabilité, sur des appels d'offre, à aller solliciter ces entreprises. La réindustrialisation et la décarbonation souhaitées de notre pays passent par cette démarche.