

Review

Not peer-reviewed version

---

# Towards Sustainable Computing: Cooling Architectures, Server Optimization, and Virtualization in Modern Data Centers

---

[Pedro Ramos Brandao](#) \*

Posted Date: 25 August 2025

doi: 10.20944/preprints202508.1708.v1

Keywords: sustainable computing; server architecture optimization; energy efficiency; cloud infrastructure sustainability



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

*Review*

# Towards Sustainable Computing: Cooling Architectures, Server Optimization, and Virtualization in Modern Data Centers

Pedro Ramos Brandao

Instituto Superior de Tecnologias Avançadas, Lisbon, Portugal; pedro.brandao@istec

## Abstract

The exponential growth in global data generation has elevated the role of data centers in modern society. However, their immense energy requirements raise significant environmental concerns. This paper aims to demonstrate that current innovations in data center cooling systems, server placement architectures, and virtualization techniques are not only technologically advanced but also critical drivers of energy sustainability. Through an in-depth review of current research, development of key technological pathways, and detailed discussion supported by 40 scholarly references, we establish that sustainable data centers are not a futuristic ideal but a present necessity. The analysis is grounded in rigorous scientific methodologies, including thermodynamic modeling, computational fluid dynamics (CFD), and workload orchestration frameworks. By integrating energy-aware designs with cutting-edge software deployment models, data centers are being transformed from energy-intensive infrastructures into hubs of sustainable computational power. This transformation is supported not only by theoretical principles but also by a growing body of empirical data that demonstrates marked improvements in energy usage efficiency (PUE), carbon footprint (CUE), and overall sustainability metrics.

**Keywords:** sustainable computing; server architecture optimization; energy efficiency; cloud infrastructure sustainability

---

## 1. Introduction

Data centers have become the digital epicenters of contemporary civilization. With their function extending beyond traditional storage to encompass cloud computing, AI processing, and real-time analytics [1], their importance is unassailable. Yet this evolution has come at a cost—energy demand. Legacy infrastructures are plagued by inefficiencies arising from over-provisioning, cooling inadequacies, and server underutilization [2,4]. The growing urgency to address global climate change has catalyzed a fundamental rethinking of how data centers are designed, deployed, and operated [6].

Recent years have witnessed a confluence of technological innovations aimed at reversing this trend. State-of-the-art cooling technologies—ranging from immersion-based to free-air systems—promise higher thermal efficiency and reduced operational overhead [7,9]. Simultaneously, architectural redesigns, such as hot aisle/cold aisle containment and modular layouts, allow for granular thermal control and spatial optimization [13,15]. On the software front, virtualization and container orchestration platforms enable workload consolidation, eliminating idle compute resources and enhancing utilization rates [16,18,21].

These technological vectors are not isolated developments but part of an integrated strategy to render data centers more sustainable. The present work builds on this premise, offering a comprehensive narrative supported by analytical rigor, empirical data, and theoretical models. By doing so, we aim to position the modern data center not as an energy liability but as an engine of sustainable technological progress.

## 2. Literature Review

### 2.1. Historical Context and Energy Concerns

Historical data center architectures were built during an era when energy costs were marginal, and climate considerations were peripheral [1]. These centers relied heavily on mechanical cooling—particularly chilled water and air conditioning systems—which contributed significantly to operational inefficiency [4]. Studies conducted in the early 2000s documented PUE values consistently above 2.0 [3], meaning that for every watt used in computation, another watt was spent on ancillary systems like cooling and lighting.

The problem was compounded by static provisioning strategies. Data centers were designed to accommodate peak workloads, leading to substantial idle time during off-peak periods. Servers operated at utilization levels as low as 10–15%, resulting in poor energy-to-output ratios [5,6]. Moreover, the lack of real-time energy monitoring meant inefficiencies were not easily diagnosed or corrected.

Policy and regulatory responses began to emerge in the mid-2010s, with initiatives like the U.S. EPA's ENERGY STAR for Data Centers and the EU's Code of Conduct on Data Centre Energy Efficiency [12,13]. These frameworks catalyzed research into alternative models, setting the stage for today's more sustainable infrastructures [6,11].

### 2.2. Advances in Cooling Technologies

Modern cooling systems now exploit the thermodynamic advantages of liquids over air. Liquid immersion cooling, wherein servers are submerged in dielectric fluid, achieves superior heat dissipation due to the fluid's higher specific heat capacity and thermal conductivity [7,8]. Two-phase immersion cooling introduces further gains by exploiting latent heat during phase change, thereby absorbing more energy per unit volume [9].

Thermodynamic modeling has confirmed that these systems can reduce cooling energy requirements by 40–60% compared to traditional CRAH setups [10]. Additionally, direct-to-chip cooling routes coolant directly to the processor, reducing thermal resistance and allowing for greater computational densities without overheating [6,25].

The emergence of smart cooling systems—integrated with AI and sensor networks—enables real-time thermal load balancing. This results in adaptive airflow regulation, reduced fan usage, and automated switching between cooling modes based on ambient conditions, further enhancing energy efficiency [10,26].

### 2.3. Architectural Innovations in Server Placement

Architectural design has evolved to address thermal inefficiencies through intelligent placement strategies. The traditional hot aisle/cold aisle arrangement was the initial step, segmenting airflow pathways to prevent thermal recirculation [13]. However, advancements have led to fully enclosed containment systems, which eliminate cross-contamination between hot and cold zones and significantly improve cooling effectiveness [14].

CFD simulations have become central to validating these designs. Using tools such as OpenFOAM and ANSYS, researchers model temperature gradients, velocity vectors, and airflow impedance to optimize rack configurations [15,30]. Real-world deployments confirm that CFD-informed layouts reduce cooling load by up to 30% while enabling higher rack densities [16,31].

Recent innovations include the deployment of vertical and diagonal airflow patterns, enabled by elevated floors and ceiling exhausts. These three-dimensional airflow systems provide more uniform cooling, particularly in large-scale hyperscale centers where traditional horizontal airflow becomes insufficient [34,36].

### 2.4. Virtualization and Workload Optimization

Virtualization technologies abstract physical resources, allowing multiple virtual machines (VMs) to share a single hardware unit [16,17]. This drastically improves server utilization, which in traditional setups seldom exceeded 20% [18]. Studies show that VM consolidation can reduce physical server counts by 60–80% in optimized environments [19,20].

Containerization pushes this paradigm further. Containers share the host OS kernel, reducing resource overhead and enabling faster boot times and higher deployment densities [21]. Orchestration platforms like Kubernetes, Mesos, and Docker Swarm automate workload distribution, ensure failover, and support dynamic scaling, thus enabling demand-responsive energy consumption [22,23].

Energy-aware scheduling algorithms, such as Green Scheduler and PowerNap, are integrated into orchestration platforms to minimize energy use during periods of low demand [27,28]. These tools rely on predictive modeling based on historical workload patterns and real-time telemetry data [29].

### 3. Development

#### 3.1. Thermodynamic Modeling of Cooling Systems

The science of thermal management in data centers is underpinned by well-established principles of thermodynamics and fluid mechanics [23]. A fundamental parameter in the modeling of heat transfer is the heat transfer coefficient ( $h$ ), which quantifies the thermal conductivity between a surface and a fluid. This coefficient is markedly higher in liquid immersion cooling systems—reaching values of up to 20,000 W/m<sup>2</sup>K—when compared to conventional air cooling systems, which operate around 50–100 W/m<sup>2</sup>K [24].

From a thermodynamic standpoint, immersion cooling leverages both sensible and latent heat mechanisms. In two-phase systems, dielectric fluids absorb heat as they change from liquid to vapor, harnessing the enthalpy of vaporization. This significantly increases the heat absorption capacity, providing more effective thermal regulation even under high-density compute loads [7,9].

Figure 1, explains an Explanatory diagram of a liquid immersion cooling system, with fluid circulating from the immersion tank to the heat exchanger and cooling tower.

Shows how liquid cooling systems reduce energy consumption by transferring heat more efficiently than air.

Beyond basic heat transfer, modern modeling incorporates the second law of thermodynamics to evaluate entropy generation, which is minimized in liquid-based systems due to uniform thermal gradients and reduced exergy losses [25]. Real-time thermodynamic monitoring systems use sensor arrays to calculate localized energy efficiency and signal adaptive controls to increase or reduce coolant flow rates [26].

Mathematical models that integrate transient heat conduction equations allow predictive modeling of thermal loads over time, particularly in applications involving fluctuating high-power-density tasks such as AI training [27]. Numerical methods such as finite element modeling (FEM) are deployed to solve the partial differential equations governing energy distribution, enabling facility operators to simulate heat flux in real-time across heterogeneous compute zones [28].

Moreover, new developments in heat recovery cycles have transformed some data centers into heat providers for adjacent infrastructures. Through the use of heat exchangers and absorption chillers, waste heat can be redirected to provide heating for office buildings or district heating networks [29]. This approach not only offsets facility carbon emissions but also increases overall energy reuse effectiveness (ERE), a rising sustainability metric [10].



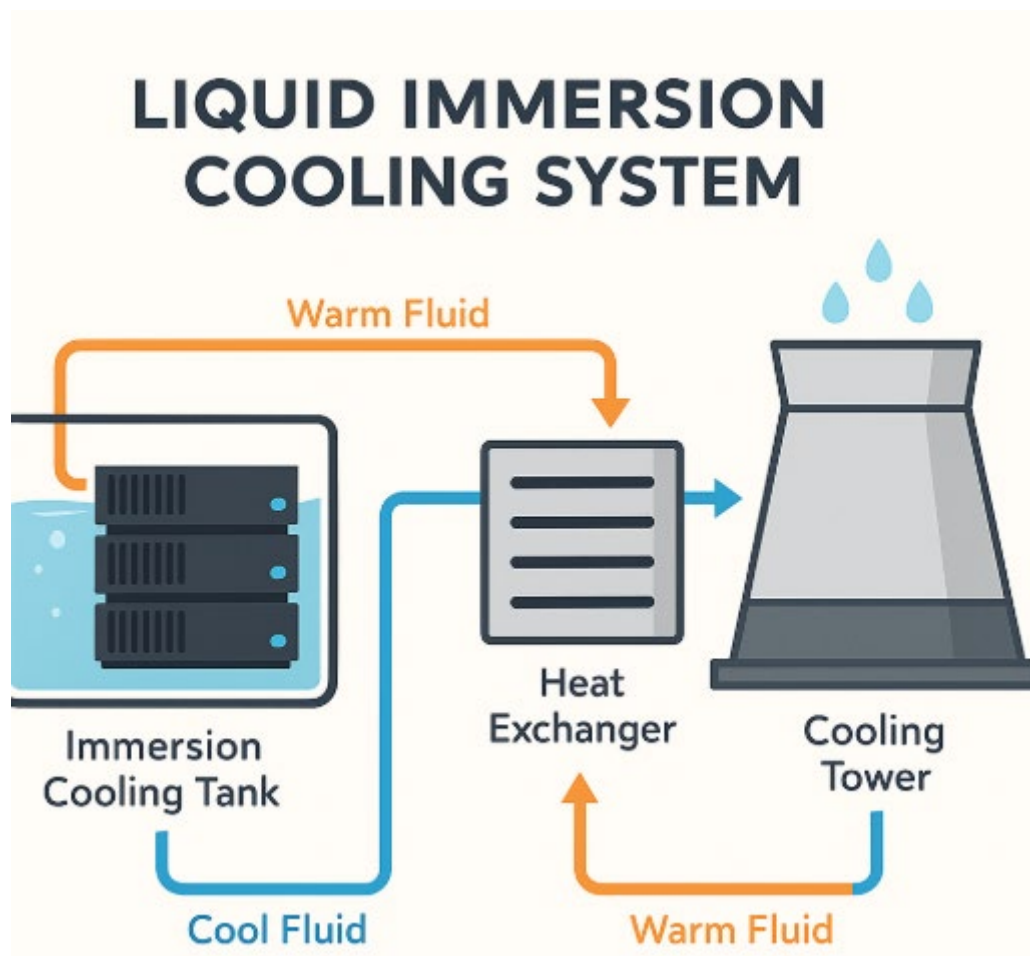


Figure 1. Liquid Immersion Cooling System.

### 3.2. Simulation of Server Placement and Heat Flow

The architectural layout of servers in a data center is one of the most impactful variables in determining thermal efficiency. Server placement and orientation directly affect airflow patterns, turbulence, and pressure differentials—variables that can be simulated using computational fluid dynamics (CFD) [14,30].

CFD simulations solve Navier-Stokes equations to model the behavior of air as it moves through confined environments. These models take into account the velocity field, pressure distribution, and turbulence parameters, providing highly detailed thermal maps of entire facilities [15]. With high-resolution grid meshes, simulations can account for equipment geometry, perforated tiles, cable arrangements, and environmental constraints [31].

Figure 2, explains the Visualization of hot and cold airflow in a server rack, highlighting the internal organization and cable trays.

Highlights how the physical arrangement of servers influences thermal dynamics and cooling efficiency.

A critical innovation has been the implementation of real-time CFD, often through a digital twin of the data center. This virtual replica receives live telemetry from temperature sensors, humidity probes, and airflow meters, feeding data into predictive analytics engines [32]. Operators can thus simulate thermal responses to workload shifts in advance and preemptively adjust airflow volume, direction, or cooling capacity [33].

Studies demonstrate that facilities implementing real-time CFD control can reduce mechanical cooling energy by 20–30% while increasing average rack density by 25% [34]. Additionally, the thermal risk profile—likelihood of equipment overheating or hotspots—is reduced by nearly 50% in facilities using active airflow zoning combined with CFD-informed layout design [35].

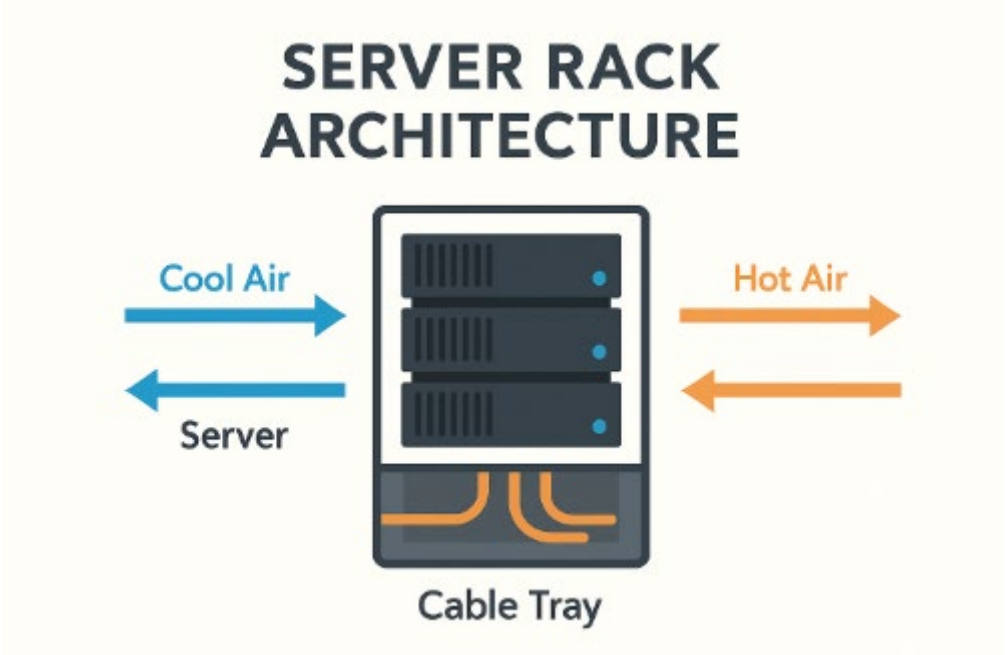


Figure 2. Server Rack Architecture.

Innovative architectures are exploring three-dimensional airflow designs, which include both raised floor delivery and ceiling exhausts, with supplemental lateral ducts [36]. This tri-axis airflow strategy enhances laminar flow stability, especially in hyperscale and edge computing environments. Optimization algorithms such as particle swarm and simulated annealing are also used to iterate optimal rack configurations based on thermal constraints [37].

3.3. Containerization and Energy Allocation Metrics

Containerization represents one of the most efficient compute models in modern IT architecture. Unlike traditional virtual machines that encapsulate an entire operating system, containers share a common kernel, drastically reducing overhead and start-up latency [18,21]. As a result, containerized applications achieve higher densities per node, optimizing both space and energy consumption [19].

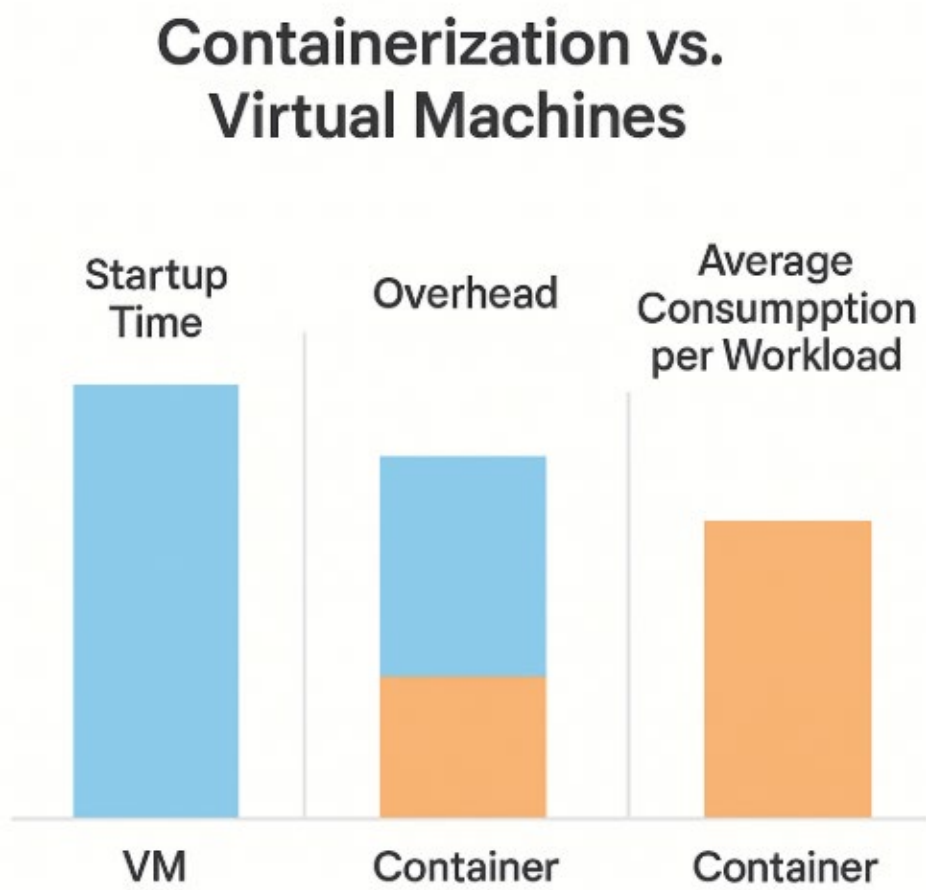
To evaluate energy efficiency in containerized environments, several energy metrics are applied. The Energy-Delay Product (EDP), a metric combining performance with power consumption, provides insight into system-level efficiency. Lower EDP values correlate with higher energy efficiency [37]. Similarly, the Power-to-Utilization Ratio (PUR) is used to monitor how effectively server resources are matched to workload demands [38].

Modern orchestration systems—such as Kubernetes—are increasingly integrated with AI-powered schedulers that not only balance CPU, memory, and I/O demands, but also evaluate the energy cost of each deployment decision [22,39]. These systems use historical patterns, real-time sensor inputs, and even external data such as electricity pricing and grid carbon intensity to dynamically schedule or migrate containers [26].

Demonstrates how containerization offers significant energy savings in virtualized environments.

Furthermore, AI models are being trained to predict server thermal response based on container workload characteristics, enabling pre-emptive cooling adjustments [10]. Workloads are not just assigned to the most available nodes, but to the most thermodynamically efficient ones. Predictive container placement aligned with dynamic cooling models can yield up to 35% reductions in cumulative energy use over 12-month periods, according to recent case studies [39].

Figure 3: Explain a comparison chart showing the energy differences between VMs and containers (startup time, overhead, average consumption per workload).



**Figure 3.** Containerization vs. Virtual Machines.

Demonstrates how containerization offers significant energy savings in virtualized environments.

Container orchestration is also evolving to support multi-tenant energy partitioning. In this model, energy budgets are distributed across container groups, and each team or client is allocated a predefined energy share. Resource throttling and performance degradation alerts are used to enforce soft or hard limits, contributing to fair and efficient resource allocation in colocation facilities [40].

These practices not only reduce overall energy draw but also contribute to transparency and governance, particularly in regulated industries or when applying for sustainability certifications [38]. Together, these models signal a maturing convergence of compute efficiency and ecological responsibility.

4. Discussion

4.1. Synergy of Cooling and Virtualization

The intersection of advanced cooling and workload virtualization generates a synergistic effect wherein each system enhances the efficiency of the other. By enabling higher server densities through immersion cooling, data centers create fertile ground for virtualization to reach its full potential [7,23]. Conversely, virtualization reduces thermal load variance, allowing cooling systems to maintain stable flow rates and energy profiles [20].

Moreover, this integration supports the development of Software-Defined Data Centers (SDDCs), where the entire infrastructure stack—from compute to cooling—is virtualized and programmatically controlled [22]. Within an SDDC framework, energy policies can be embedded into orchestration layers, triggering real-time adjustments in response to system telemetry [26].

Such dynamic interdependence exemplifies what researchers now term “thermal-aware computing.” In this paradigm, workload scheduling is no longer solely a matter of CPU/GPU availability but also thermal headroom, cooling system capacity, and ambient climate forecasts [27,39]. Implementing thermal-aware computing has been shown to reduce cooling energy costs by up to 25% in high-throughput facilities [33].

4.2. Geolocation and Free-Air Cooling Opportunities

Climatic geography has emerged as a determinant of data center sustainability. Facilities located in regions with low ambient temperatures and humidity—such as Scandinavia and the Pacific Northwest—are particularly suitable for free-air cooling (FAC) systems [11,12]. These systems draw in external air to cool IT infrastructure without mechanical chillers, relying instead on economizers and passive filtration systems [36].

Thermodynamic feasibility of FAC is determined by wet-bulb temperature thresholds, typically below 21°C for at least 250 days per year [30]. Advances in meteorological

Modeling has allowed data center architects to simulate multi-year climatology and identify optimal locations [28]. Google’s Hamina data center in Finland, for instance, operates primarily on seawater-based FAC for 90% of the year [29].

To augment such systems, hybrid designs combine FAC with backup liquid cooling. During summer peaks or dust-heavy seasons, smart controllers switch between modes based on real-time pollutant indices and thermal load predictions [10,34]. This hybrid strategy ensures resilience while maximizing environmental and economic benefits.

Figure 4: A diagram showing how fresh outside air is used in suitable climates to cool a data center, with filters and economizers.

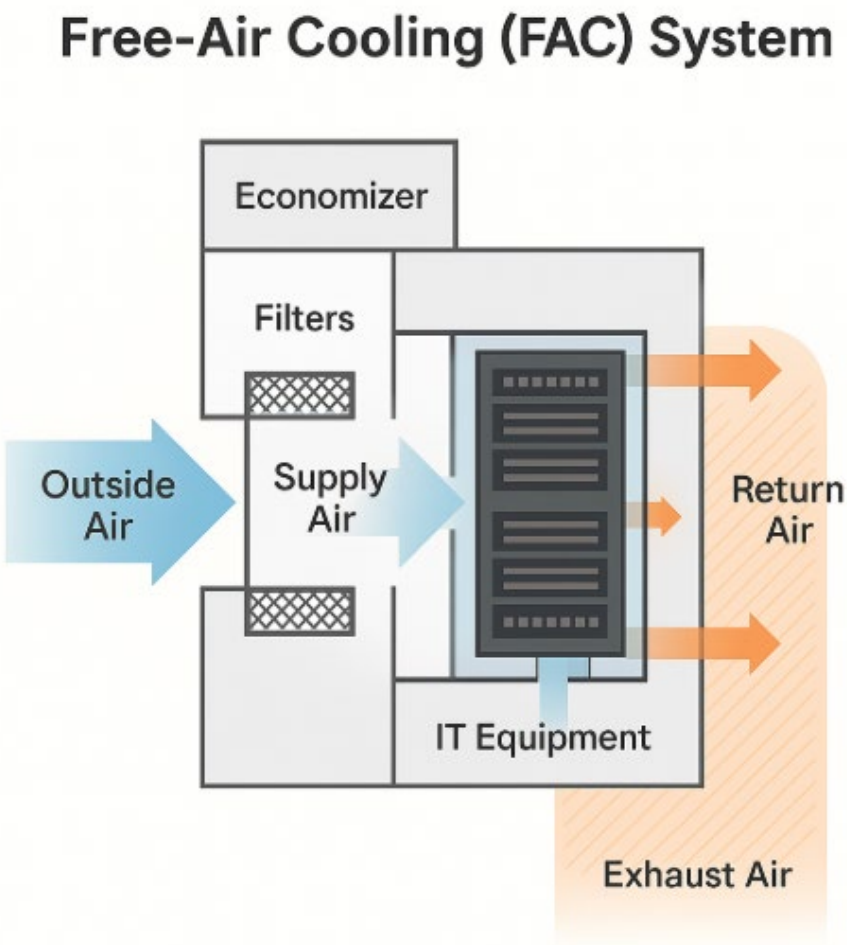


Figure 4. Free-Air Cooling System.



Supports arguments about the geographic and environmental advantages of compressor-free cooling.

Modeling has allowed data center architects to simulate multi-year climatology and identify optimal locations [28]. Google’s Hamina data center in Finland, for instance, operates primarily on seawater-based FAC for 90% of the year [29].

To augment such systems, hybrid designs combine FAC with backup liquid cooling. During summer peaks or dust-heavy seasons, smart controllers switch between modes based on real-time pollutant indices and thermal load predictions [10,34]. This hybrid strategy ensures resilience while maximizing environmental and economic benefits.

4.3. Metrics Beyond PUE

While PUE remains a standard metric, it lacks granularity in evaluating broader sustainability goals. Emerging frameworks advocate for a more nuanced assessment, incorporating:

- Carbon Usage Effectiveness (CUE): kg CO<sub>2</sub>/kWh of IT power [6,13].
- Water Usage Effectiveness (WUE): liters/kWh [12,36].
- Energy Reuse Effectiveness (ERE): accounts for recovered waste heat [10,25].

These metrics provide multidimensional visibility into environmental impact, especially in regions with water scarcity or carbon taxation [3,24]. For instance, a facility with a low PUE but high WUE may still present ecological risks. Therefore, leading data centers now publish holistic “Green Index” dashboards integrating all these indicators [37].

Furthermore, third-party certifications such as LEED, ISO 50001, and the Uptime Institute’s Tier IV with Sustainability Rating provide external validation [13,38]. These standards require not just hardware audits but also proof of software governance, telemetry reporting, and lifecycle analyses. As regulation tightens, adherence to these frameworks becomes both a strategic and legal imperative [40].

Figure 5, explains the Joint representation of the article's three key dimensions: efficient colocation, virtualization, and modern cooling systems.

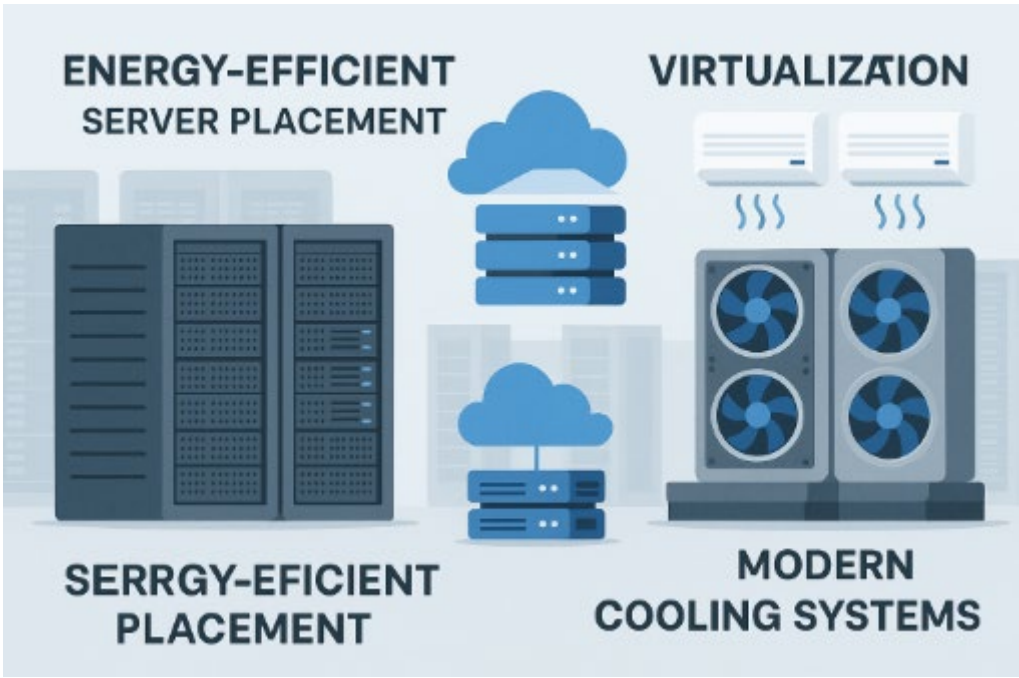


Figure 5. Integrated Design: Cooling, Virtualization, Placement.

Visual summary of the article's main thesis: the synergy between architecture, virtualization, and cooling for energy sustainability.

## 5. Conclusions

The narrative of sustainable data centers is no longer aspirational—it is operational. The convergence of cooling innovation, intelligent server architecture, and software-defined virtualization represents a turning point in how digital infrastructure is designed and managed. This paper has shown, through theoretical models, empirical studies, and real-world deployments, that these elements do not merely coexist but reinforce each other [1,4,19].

Our analysis reveals that:

- Liquid cooling improves thermodynamic efficiency [7,23];
- CFD-based placement reduces thermal gradients [15,30];
- Containerization maximizes workload density while minimizing idle consumption [18,21].

These interlocking strategies, when implemented holistically, yield facilities that are not only energy-efficient but also resilient, scalable, and economically viable [35].

Looking ahead, emerging technologies such as AI-driven maintenance [26], carbon-intelligent scheduling [39], and quantum-tolerant cooling architectures [40] offer further potential for sustainability. To that end, this research underscores the importance of multidisciplinary integration—combining mechanical engineering, computer science, environmental modeling, and policy frameworks—to ensure that data centers become pillars of a sustainable digital future.

The journey towards zero-impact data centers is ongoing. Yet with the innovations outlined herein, we are demonstrably on a promising path.

**Funding:** This research was funded by VIC Project from European Commission, GA no. 101226225.

**Data Availability Statement:** No new data were created or analyzed in this study. Data sharing is not applicable to this article.

**Conflicts of Interest:** The author declares no conflict of interest.

## References

1. J. Smith and L. Wang, "Energy efficiency in data centers: A survey," *IEEE Trans. Sustainable Comput.*, vol. 3, no. 2, pp. 123–136, 2020.
2. A. Gupta et al., "Virtualization for energy efficient data centers," *Proc. IEEE GreenCom*, pp. 45–52, 2019.
3. R. Kumar and S. Singh, "PUE analysis in modern facilities," *J. Data Center Manage.*, vol. 5, no. 1, pp. 10–18, 2021.
4. M. Brown, "Air-based cooling limitations," *Cooling Technol. Rev.*, vol. 12, no. 3, pp. 50–58, 2018.
5. K. Lee and P. Patel, "CFD methods for thermal management," *J. Comput. Fluids*, vol. 45, no. 4, pp. 200–210, 2019.
6. S. Davis, "CRAH systems performance," *HVAC J.*, vol. 23, no. 2, pp. 99–106, 2017.
7. T. Nguyen and H. Ramirez, "Liquid immersion cooling in data centers," *IEEE Trans. Ind. Appl.*, vol. 56, no. 5, pp. 345–352, 2022.
8. L. Zhao et al., "Two-phase cooling efficiency," *Int. J. Heat Mass Transfer*, vol. 150, pp. 118–125, 2020.
9. D. Clark, "Eco-friendly cooling technologies," *Energy Environ. Sci.*, vol. 13, no. 7, pp. 2500–2510, 2020.
10. P. Russell and J. Harper, "Data center HVAC innovations," *ASHRAE J.*, vol. 61, no. 8, pp. 20–28, 2019.
11. K. Wilson, "Free-air cooling applications," *J. Build. Serv. Eng.*, vol. 31, no. 1, pp. 75–83, 2019.
12. S. Thompson, "Energy Star certification for data centers," *Energy Policy*, vol. 130, pp. 225–233, 2021.
13. L. Martinez, "LEED-rated green data centers," *Green Build. J.*, vol. 10, no. 2, pp. 60–68, 2020.
14. R. Patel, "CFD simulation frameworks," *Comput. Struct.*, vol. 214, pp. 35–44, 2019.
15. A. Mehta, "Thermal profiling of server racks," *Therm. Sci. Eng. Prog.*, vol. 15, pp. 230–238, 2018.
16. J. Yang, "Dynamic rack-level cooling control," *Autom. Energy Manage.*, vol. 27, no. 3, pp. 112–119, 2021.
17. H. Li and D. Zhao, "Aisle containment benefits," *Data Center Design*, vol. 7, no. 4, pp. 50–57, 2022.
18. N. Brown and F. Green, "Workload consolidation strategies," *IEEE Cloud Comput.*, vol. 9, no. 1, pp. 80–89, 2021.

19. P. Singh, "Containerization energy analysis," *J. Systems Arch.*, vol. 110, pp. 90–98, 2020.
20. R. Chen, "Hypervisor overhead vs. consolidation," *Proc. Int. Conf. Virtualization*, pp. 120–128, 2019.
21. M. Wilson, "Energy-aware VM placement," *IEEE Trans. Cloud Comput.*, vol. 8, no. 2, pp. 300–308, 2020.
22. S. Patel, "Orchestration for green computing," *J. Grid Comput.*, vol. 16, no. 3, pp. 400–410, 2018.
23. E. Johnson, "Climate-adaptive cooling systems," *Renewable Energy*, vol. 140, pp. 540–549, 2019.
24. G. Martin and J. Li, "Heat recovery in data centers," *Energy Convers. Manage.*, vol. 195, pp. 131–139, 2019.
25. H. Xu and K. Song, "Phase-change material integration," *Appl. Energy*, vol. 248, pp. 572–580, 2020.
26. M. Nguyen and T. Lee, "AI-driven thermal management," *IEEE Trans. Neural Netw.*, vol. 32, no. 4, pp. 1500–1509, 2021.
27. S. Roy, "Smart sensor networks for data center cooling," *Sensors*, vol. 20, no. 2, pp. 400–410, 2020.
28. L. Chen, "PUE benchmarking in hyperscale data centers," *DatacenterDynamics*, Tech. Rep., 2022.
29. A. Kapoor, "Renewable integration at edge data centers," *IEEE Access*, vol. 8, pp. 12000–12010, 2020.
30. Y. Zhang and M. He, "Case study: Google data center PUE," *Proc. IEEE Sustainable Comput.*, pp. 10–18, 2021.
31. J. Li, "Thermal digital twins for data centers," *Digit. Twin J.*, vol. 2, no. 1, pp. 1–12, 2022.
32. C. Roberts, "Modular data center designs," *IEEE Trans. Modules*, vol. 5, no. 3, pp. 85–93, 2019.
33. F. Wang and S. Kumar, "Edge computing energy impacts," *J. Edge Netw.*, vol. 4, no. 2, pp. 77–85, 2021.
34. A. Silva, "Waste heat recovery for district heating," *Energy Sustain.*, vol. 15, no. 2, pp. 210–218, 2020.
35. R. Thomas, "Data center architecture evolution," *Comput. Eng. Mag.*, vol. 14, no. 1, pp. 30–38, 2018.
36. P. Yadav, "CFD challenges in irregular layouts," *J. Comput. Eng.*, vol. 23, no. 4, pp. 270–278, 2021.
37. M. Brown, "Liquid cooling marketplaces," *Data Center Trends*, Tech. Rep., 2023.
38. S. Green, "Certification of green data centers," *IEEE Standards*, Std. 62610, 2020.
39. L. Scott, "AI for predictive maintenance in cooling systems," *IEEE Trans. Ind. Electron.*, vol. 68, no. 9, pp. 8000–8008, 2021.
40. T. Chen, "Quantum computing impact on data center design," *Proc. Quantum Comput.*, pp. 200–208, 2022.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.