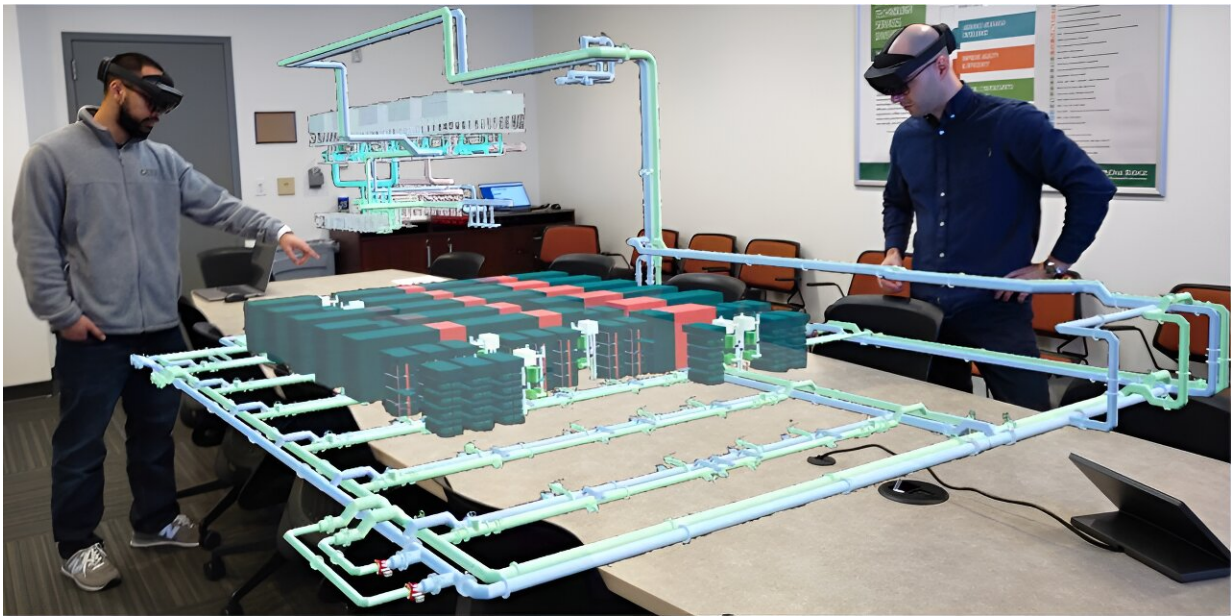


Computer engineers pioneer approaches to energy-efficient supercomputing

September 11 2024, by Coury Z Turczyn



From left, Sedrick Bouknight and Matthias Maiterth of ORNL's Analytics and AI Methods at Scale group demonstrate the VR capabilities of the Frontier digital twin project's ExaDIGIT framework. Using VR allows Frontier's operators to exam the system's telemetry in a more interactive and intuitive way. Credit: Wes Brewer/ORNL, U.S. Dept. of Energy

As high-tech companies ramp up construction of massive data centers to meet the business boom in artificial intelligence, one component is becoming an increasingly rare commodity: electricity.

Commercial demand for electricity has been growing sharply in recent years and is projected to increase by 3% in 2024 alone, according to the U.S. Energy Information Administration. But that [growth has been driven by just a few states](#)—the ones that are rapidly becoming hubs for large-scale computing facilities, such as Virginia and Texas.

The inventory of North American data centers grew 24.4% year over year in the first quarter of 2024, as the real-estate services firm CBRE reports in its "Global Data Center Trends 2024" study. These new centers are being built with capacities of 100 to 1,000 megawatts, or about the same loads that can power from 80,000 to 800,000 homes, notes the Electric Power Research Institute in a [2024 white paper](#).

In this paper, EPRI analyzes AI and data-center energy consumption and predicts that if a projected high growth rate of 10% per year continues, data centers will annually consume up to 6.8% of total U.S. electricity generation by 2030—versus an estimated 4% today.

To satisfy that soaring demand, Goldman Sachs Research estimates that U.S. utilities will need to invest around \$50 billion in new electrical generation capacity. Meanwhile, community opposition to data center construction in some areas is also growing, as grassroots groups protest the potential local impacts of more and more data centers and their increasing demands for electricity for AI and water for cooling.

Whether the nation's private enterprises can pull off the daunting challenge of powering an AI "revolution" may depend less on money and more on ingenuity. That CBRE study concludes with a helpful, or perhaps hopeful, recommendation: "High-performance computing [or HPC] will require rapid innovation in data center design and technology to manage rising [power density](#) needs."

At the Oak Ridge Leadership Computing Facility, a Department of

Energy Office of Science user facility located at Oak Ridge National Laboratory, investigating new approaches to energy-efficient supercomputing has always been part of its mission.

Since its formation in 2004, the OLCF has fielded five generations of world-class supercomputing systems that have produced a nearly 2,000 times increase in energy efficiency per floating point operation per second, or flops. Frontier, the OLCF's latest supercomputer, currently [ranks first in the TOP500 list](#) of the world's most [powerful computers](#), and in 2022, it debuted at the [top of the Green500 list](#) of the world's most energy-efficient computers.

Keeping the electricity bill affordable goes hand in hand with being a government-funded facility. But constructing and maintaining leadership supercomputers are no longer just the domain of government. Major tech companies have entered HPC in a big way but are only just now starting to worry about how much power these mega systems consume.

"Our machines were always the biggest ones on the planet, but that is no longer true. Private companies are now deploying machines that are several times larger than Frontier. Today, they essentially have unlimited deep pockets, so it's easy for them to stand up a data center without concern for efficiency," said Scott Atchley, chief technology officer of the National Center for Computational Sciences, or NCCS, at ORNL. "That will change once they become more power constrained, and they will want to get the most bang for their buck."

With decades of experience in making HPC more energy efficient, the OLCF may serve as a resource for best "bang for the buck" practices in a suddenly burgeoning industry.

"We are uniquely positioned to influence the full energy-efficiency ecosystem of HPC, from the applications to the hardware to the

facilities. And you need efficiency gains in all three of those areas to attack the problem," said Ashley Barker, OLCF program director.

"Striving for improvements in energy efficiency comes into play in every aspect of our facility. What is the most energy-efficient hardware we can buy? What is the most energy-efficient way we can run that hardware? And what are the most energy-efficient ways that we can tweak the applications that run on the hardware?"

As the OLCF plans its successor to Frontier—called Discovery—those questions are asked daily as different teams work together to deliver a new supercomputer by 2028 that will also demonstrate next-generation energy efficiencies in HPC.

System hardware

One of the most significant computational efficiency advancements of the past 30 years originated from an unlikely source: video games.

More specifically, the innovation came from chip makers competing to fulfill the video game industry's need for increasingly sophisticated in-game graphics. To achieve the realistic visuals that drew in gamers, personal computers and game consoles required dedicated chips—also known as the graphics processing unit, or GPU—to render detailed moving images.

Today, GPUs are an indispensable part of most supercomputers, especially ones used for training artificial intelligence models. In 2012, when the OLCF pioneered the use of GPUs in leadership-scale HPC with its Titan supercomputer, the design was considered a bold departure from traditional systems that rely only on central processing units, or CPUs.

It required computational scientists to adapt their codes to fully exploit the GPU's ability to churn through simple calculations and speed up the time to solution. The less time it takes a computer to solve a particular problem, the more problems it can solve in a given time frame.

"A GPU is, by design, more energy efficient than a CPU. Why is it more efficient? If you're going to run electricity into a computer and you want it to do calculations very efficiently, then you want almost all the electricity powering floating point operations. You want as much silicon area to just be floating point units, not all the other stuff that's on every CPU chip.

"A GPU is almost pure floating point units. When you run electricity into a machine with GPUs, it takes roughly about a tenth the amount of energy as a machine that just has CPUs," said ORNL's Al Geist, director of the Frontier project.

The OLCF's gamble on GPUs in 2012 paid off over the next decade with progressively more energy-efficient systems as each generation of OLCF supercomputer increased its number of speedier GPUs. This evolution culminated in the architecture of Frontier, launched in 2022 as the world's first exascale supercomputer, capable of more than 1 quintillion calculations per second and consisting of 9,408 compute nodes.

However, when exascale discussions began in 2008, the Exascale Study Group issued a report outlining its [four biggest challenges](#), foremost of which was power consumption. It foresaw an electric bill of potentially \$500 million a year. Even accounting for the projected technological advances of 2015, the report predicted that a stripped-down 1-exaflop system would use 150 megawatts of electricity.

"DOE said, 'That's a non-starter.' Well, we asked, what would be

acceptable? And the answer that came back was, 'We don't want you to spend more money on electricity than the cost of the machine,'" Geist said. "In the 2009 time frame, supercomputers cost about \$100 million. They have a lifetime of about five years.

"What you end up with is about \$20 million per year that we could spend on electricity. How many megawatts can I get out of \$20 million? It turns out that 1 megawatt here in East Tennessee is \$1 million a year, roughly. So that was the number we set as our target: a 20-megawatt per exaflop system."

There wasn't a clear path to achieving that energy consumption goal. So, in 2012, the DOE Office of Science launched the FastForward and DesignForward programs to work with vendors to advance new technologies.

FastForward initially focused on the processor, memory and storage vendors to address performance, power-consumption and resiliency issues. It later moved its focus to node design (i.e., the individual compute server). DesignForward initially focused on scaling networks to the anticipated system sizes and later focused on whole-system packaging, integration and engineering.

As a result of the FastForward investment, semiconductor chip vendor AMD developed a faster, more powerful compute node for Frontier—consisting of a 64-core 3rd Gen EPYC CPU and four Instinct MI250X GPUs—and figured out a way to make the GPUs more efficient by turning off sections of the chips that are not being used and then turning them back on when needed in just a few milliseconds.

"In the old days, the entire system would light up and sit there idle, still burning electricity. Now we can turn off everything that's not being used—and not just a whole GPU. On Frontier, about 50 different areas

on each GPU can be turned off individually if they're not being used. Now, not only is the silicon area mostly devoted to floating point operations, but in fact I'm not going to waste any energy on anything I'm not using," Geist said.

However, with the next generation of supercomputers, simply continuing to add more GPUs to achieve more calculations per watt may have reached its point of diminishing returns, even with newer and more advanced architectures.

"The processor vendors will really have to reach into their bag of tricks to come up with techniques that will give them just small, incremental improvements. And that's not only true for energy efficiency, but it's also true for performance. They're getting about as much performance out of the silicon as they can," Atchley said.

"We've been benefiting from Moore's Law: transistors got smaller, they got cheaper and they got faster. Our applications ran faster, and the price point was the same or less. That world is over. There are some possible technologies out there that might give us some jumps, but the biggest thing that will help us is a more integrated, holistic approach to energy efficiency."

System operations

Feiyi Wang—leader of the OLCF's Analytics and AI Methods at Scale, or AAIMS, group—has been spending much of his time pondering an elusive goal: how to operate a supercomputer so that it uses less energy. Tackling this problem first required the assembly of a massive amount of HPC operational data.

Long before Frontier was built, he and the AAIMS group collected over one year's worth of power profiling data from Summit, the OLCF's

200-petaflop supercomputer launched in 2018. Summit's 4,608 nodes each have over 100 sensors that report metrics at 1 hertz, meaning that for every second, the system reports over 460,000 metrics.

Using this 10-terabyte dataset, Wang's team analyzed Summit's entire system from end to end, including its central energy plant, which contains all its cooling machinery. They overlaid the system's job allocation history on the telemetry data to construct per-job, fine-grained power-consumption profiles for over 840,000 jobs. This work earned them the [Best Paper Award](#) at the 2021 International Conference for High Performance Computing, Networking, Storage, and Analysis, or SC21.

The effort also led Wang to come up with a few ideas about how such data can be used to make informed operational decisions for better energy efficiency.

Using the energy-profile datasets from Summit, Wang and his team kicked off the Smart Facility for Science project to provide ongoing production insight into HPC systems and give system operators "data-driven operational intelligence," as Wang puts it.

"I want to take this continuous monitoring one step further to 'continuous integration,' meaning that we want to take the computer's ongoing metrics and integrate them into a system so that the user can observe how their energy usage is going to be for their particular job application. Taking this further, we also want to implement 'continuous optimization,' going from just monitoring and integration to actually optimizing the work on the fly," Wang said.

Another one of Wang's ideas may assist in that goal. At SC23, Wang and lead author Wes Brewer, a senior research scientist in the AAIMS group, delivered a presentation, "Toward the Development of a Comprehensive

Digital Twin of an Exascale Supercomputer." They proposed a framework called ExaDIGIT that uses augmented reality, or AR, and virtual reality, or VR, to provide holistic insights into how a facility operates to improve its overall energy efficiency.

Now, ExaDIGIT has evolved into a collaborative project of 10 international and industry partners, and Brewer will present the team's [newest paper](#) at [SC24](#) in Atlanta, Georgia.

At ORNL, the AAIMS group launched the Digital Twin for Frontier project to construct a simulation of the Frontier supercomputer. This virtual Frontier will enable operators to experiment with "What if we tried this?" energy-saving scenarios before attempting them on the real Frontier machine. What if you raised the incoming water temperature of Frontier's cooling system—would that increase its efficiency? Or will you put it at risk of not cooling the system enough, thereby driving up its failure rate?

"Frontier is a system so valuable that you can't just say, 'Let's try it out. Let's experiment on the system,' because the consequences may be destructive if you get it wrong," Wang said. "But with this digital twin idea, we can take all that telemetry data into a system where, if we have enough fidelity modeled for the power and cooling aspects of the system, we can experiment. What if I change this setting—does it have a positive effect on the system or not?"

Frontier's digital twin can be run on a desktop computer, and using VR and AR allows operators to examine the system telemetry in a more interactive and intuitive way as they adjust parameters. The AAIMS group also created a virtual scheduling system to examine the digital twin's power consumption and how it progresses over time as it runs jobs.

Although the virtual Frontier is still being developed, it is already yielding insights into how workloads can affect its cooling system and what happens with the power losses that occur during rectification, which is the process of converting alternating current to direct current. The system is also being used to predict the future power and cooling needs of Discovery.

"We can and will tailor our development as well as the system to address any current and future pressing challenges faced by the OLCF," Wang said.

Facility infrastructure

Powering a supercomputer doesn't just mean turning it on—it also means powering the entire facility that supports it. Most critical is the cooling system that must remove the heat generated by all the computer's cabinets in its data center.

"From a 10,000-foot viewpoint, a supercomputer is really just a giant heater—I take electricity from the grid, I run it into this big box, and it gets hot because it's using electricity. Now I have to run more electricity into an air conditioner to cool it back off again so that I can keep it running and it doesn't melt," Geist said.

"Inside the data center there is a lot of work that goes into cooling these big machines more efficiently. From 2009 to 2022, we have reduced the energy needed for cooling by 10 times, and our team will continue to make cooling optimizations going forward."

Much of the planning for those cooling optimizations is led by David Grant, the lead HPC mechanical engineer in ORNL's Laboratory Modernization Division. Grant oversees the design and construction of new mechanical facilities and is primarily responsible for ensuring that

every new supercomputer system installed at the OLCF has the cooling it requires to reliably operate 24-7.

He started at ORNL in 2009 and worked on operations for the Jaguar supercomputer. Then, he became involved in its transition into Titan in 2012, led Summit's infrastructure design for its launch in 2018, and most recently oversaw all the engineering to support Frontier.

In that span of time, the OLCF's cooling systems have substantially evolved alongside the chip technology, going from loud fans and chiller-based air-conditioning in Jaguar to fan-free liquid cooling in Frontier.

Furthermore, the water temperatures required to cool down the compute nodes have risen from 42°F for Titan to Frontier's 90°F—a target set by the FastForward program. That extra warmth spurs huge energy savings because the circulating water no longer needs to be refrigerated and can be sufficiently cooled by evaporative towers instead.

"We are trying to get the warmest water possible back from the cabinets while serving them the warmest water-supply temperatures—the higher the supply temperatures, the better," Grant said.

"Warmer water coming back to us allows us to minimize the flow that we have to circulate on the facility side of the system, which saves pumping energy. And then the warmer temperatures allow us to be more efficient with our cooling towers to be able to reject that heat to our environment."

Frontier's power usage effectiveness, or PUE—the ratio of the total power used by a computer data-center facility versus the power delivered to computing equipment—is delivering 1.03 at peak usage. This essentially means that for every 1,000 watts of heat, it takes just 30 watts of additional electrical power to maintain the system's appropriate

thermal envelope.

The global, industry-wide average for data centers is around 1.47 PUE, [according to the Uptime Institute](#).

Making further reductions in power usage for a faster system such as Discovery will require even more innovative approaches, which Grant is investigating.

First, the concept of recovering (or using) some of Discovery's excess heat may hold some promise. The facility is well situated to reuse waste heat if it can be moved from the cooling system to the heating system. But this task is challenging because of the elevated temperatures of the heating system, the low-grade heat from the cooling system and the highly dynamic nature of the heat being generated by the HPC systems.

Second, the incoming Discovery system will share Frontier's cooling system. Additional operational efficiencies are expected from this combined-use configuration.

"Right now, Frontier gets to sit on its own cooling system, and we've optimized it for that type of operation. But if you have Frontier demanding up to 30 megawatts and then another system demanding perhaps that much again, what does that do to our cooling system?"

"It is designed to be able to do that, but we're going to be operating at a different place in its operational envelope that we haven't seen before. So, there'll be new opportunities that present themselves once we get there," Grant said.

Third, Grant is examining how construction and equipment choices may benefit the facility's overall energy efficiency. For example, Frontier's cooling system has 20 individual cooling towers that require a process

called pacification to help protect their internal metal surfaces, and this process involves a lot of pumping over time. That step could be eliminated with newer towers that no longer require the pacification process.

Fourth, idle time on a supercomputer can use up a great deal of electricity —Frontier's idle loads are 7 to 8 megawatts. What if that idle load could be greatly reduced or eliminated?

"When we interact with the customers who have influence on the software side, we try to communicate to them how their decisions will translate through the cooling system and to the facility energy use," Grant said.

"I think there's a lot of potential on the software side to try to reduce the idle load requirement and make their models run as efficiently as possible and increase the utilization of the system. In return, they will get higher production on their side for the data that they're trying to produce."

Applications

Optimizing science applications to run more efficiently on the OLCF's supercomputers is the domain of Tom Beck, head of the NCCS's Science Engagement section, and Trey White, a distinguished research scientist in the NCCS's Algorithms and Performance Analysis group. Getting codes to return their results faster is not exactly a new concept, but the goal is now shifting away from just pure speed.

"For a long time, people have wanted to make their codes run faster, and that's what we've concentrated on—that singular goal of running faster applications, which also happened to reduce energy use," White said.

"Hardware is still increasing in speed, just not as fast as it used to, and so now we must look at applications in terms of both time and energy efficiency. For the most part, running faster means less energy, but it's not a perfect correlation. So, we are now starting to look at trade-offs between the two."

One area the team is investigating is how the operating frequency of the GPUs can impact their energy consumption. The maximum frequency for a GPU to achieve its fastest throughput may not necessarily be the most energy-efficient frequency.

"But if you start at the maximum frequency and pull back by 5% to 10%, there are some indications you might get 20% or 25% energy savings. So, then it's an arbitrage of, are you willing to give up a little bit of your performance to get big energy savings?" Beck said.

"Previously, what maximum clock frequency the computer uses was set for all projects to a single number, in general. But now we're looking at adapting that per application and maybe even within a single run," White said. "That 'frequency knob' is one example of something where there's a trade-off between time and energy efficiency, and we're investigating how to give users that choice."

Another area the team is exploring is the use of mixed-precision arithmetic. Historically, full-precision floating point arithmetic at 64 bits was considered the standard for computational accuracy in science applications. Increasingly more powerful supercomputers since the early 2000s made full precision nearly as fast to use as single-precision arithmetic at 32 bits.

Now, with the rise of the AI market, low-precision arithmetic—16 bits or fewer—has demonstrated that it is accurate enough for training neural networks and other data-science applications. Driven by GPUs, low-

precision calculations can offer substantial speedups and energy savings.

"Using lower precision is a scary landscape to users because everybody's used to assuming full precision's 64 bits and partly just because it's already there and accessible," Beck said.

"And if you start deviating from 64 bits, it could impact things in nonlinear ways throughout your code, where it's really hard to track down what's going on. So that's part of our research strategy—to do a broad study of the impacts of going to mixed-precision arithmetic in some applications."

Another area that may reap increases in energy efficiency is data transfer—the less movement of data, the less electricity required. This work could be accomplished by constructing software algorithms that reduce data movement. Beck would like to offer users pie charts that show the percentages of power used by each different operation of an algorithm, thereby allowing them to target potential reductions.

"Without a radical hardware change or revolution in the architecture, applications are really the place that people are looking now for increasing [energy efficiency](#)," Beck said. "Most likely, this is not a game of getting a 300% improvement through coding."

"There are definitely places where we can make improvements, but it's probably going to be a more incremental process of 3% here, 5% there. But if you can accumulate that over a bunch of changes and get to 20%, that's a big accomplishment."

Provided by Oak Ridge National Laboratory

Citation: Computer engineers pioneer approaches to energy-efficient supercomputing (2024,

September 11) retrieved 11 September 2024 from
<https://techxplore.com/news/2024-09-approaches-energy-efficient-supercomputing.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.