

Beyond PUE: Flexible Datacenters Empowering the Cloud to Decarbonize

Andrew A. Chien
*University of Chicago and
Argonne National Laboratory*

Chaojie Zhang, Liuzixuan Lin, Varsha Rao
University of Chicago

Abstract

Traditional datacenter design and optimization for TCO and PUE is based on static views of power grids as well as computational loads. Power grids exhibit increasingly variable price and carbon-emissions, becoming more so as government initiatives drive further decarbonization. The resulting opportunities require dynamic, temporal metrics (eg. not simple averages), flexible systems and intelligent adaptive control.

Two research areas represent new opportunities to reduce both carbon and cost in this world of variable power, carbon, and price. First, the design and optimization of flexible datacenters. Second, cloud resource, power, and application management for variable-capacity datacenters. For each, we describe the challenges and potential benefits.

1 Introduction

The creation of internet-scale services has driven the creation of large-scale datacenters which have achieved dramatic compute, cost, and power efficiency advances over the past two decades. From the launch of Amazon’s EC2 in 2007, now these cloud companies operate networks of datacenters that underpin global cloud computing [12, 23, 32].

Providing new capabilities such as scalability, on-demand access, pay-as-you-go accounting, and an array of new cloud services, the growing commercial success of internet-scale applications and cloud computing continues to drive the scale and reach of computing infrastructure. Measured by revenue growth or renewable purchases, this annual growth rate is as high as 30% for some large cloud providers such as Microsoft and Google [10, 27] over the past 5 years. Behind this rapid growth are newly built or expanded datacenters, with maximum power capacity as large as 200 MW to 1 GW [12, 13, 23, 30, 32].

The design of these datacenters gave rise to the phrase the "Datacenter is the Computer" [1, 2], reflecting a new focus in research and engineering around how to holistically design and optimize them. Key tenets of this technical activity is

the optimization of total-cost-of-ownership (TCO) that takes a holistic view of capital (eg. buy the equipment and facilities) and operating cost (eg. power, water, staffing, etc.). To optimize the rapidly growing power consumption of both individual and networks of datacenters, a range of research grew up around minimizing the use of power, particularly the overheads of power in a datacenter that are not directly delivering computing services [22, 26]. Power-use efficiency (PUE) has become a standard metric for capturing the energy efficiency of a datacenter. This potent pair – TCO and PUE has been the technical foundation for the design, construction, and operation of datacenters and the cloud for two decades.

In recent years, the grid power context for cloud datacenters has been changing. To respond to climate change, power grids in the United States have begun large-scale adoption of wind and solar renewable generation [14, 25], and over the past 20 years, these renewables at fractions of 10%, 20%, and in some grids as high as 40% have changed grid dynamics. These renewables are intermittent generators, producing low-cost, carbon-free energy – but only when the environment enables them to, not necessarily when the power grid needs it. To integrate such intermittent, distributed generation, power grids have adopted multi-market price-based dispatch [5, 16] that produces power grid with large variations in both carbon-intensity and price of power. The variation can be both fast and dramatic, and as discussed in Section 2 has significant implications for large power consumers such as datacenters, directly affecting their power cost and carbon footprint.

There is little doubt that hyperscalers have the technical capability to build more and larger datacenters, but hyperscalers face increasing pressure to reduce their carbon impact (despite growing power consumption) and to reduce their power cost (as Dennard scaling’s end has slowed electronics energy-efficiency improvement). For example, to put this in context, a gigawatt-datacenter could incur \$652M (Virginia)–\$1.3B (New England) annual electricity cost and 1.8M (California)–7.9M (Wyoming) metric-tons of CO₂ (mTCO₂) associated annual carbon emissions. The transformation of the power grid with renewable generation provides opportunity to re-

duce both of these costs, if datacenter design and software can flexibly manage variable capacity to match grid opportunity.

In this paper, we propose a new framework for datacenter design that reflects these new opportunities, based on dynamic metrics and optimization for variation – not the average PUE traditionally used. This dynamic framework exposes a set of challenging new systems research opportunities in datacenter design, workload management, and more. Specifically:

- How to design datacenters that can flexibly exploit the dynamic properties of the renewable power grid? (computer, datacenter architecture)
- What are the new pillars analogous to PUE and TCO that can frame comparison and measurement of effective design? (metrics, measurement)
- How to reinvent resource management software to reflect and efficiently exploit the dynamically varying compute capacity? (resource management)
- How to coordinate optimization of carbon, SLOs, resource efficiency, and application compute efficiency? (resource & sustainability mgmt, applications, ...)
- How to manage the distribution, sharing, and use of power within future variable power (and therefore capacity) datacenter? (power mgmt software, algorithms)

In rest of the paper, we describe varying dynamics in decarbonizing power grids in Section 2 and review traditional datacenter optimization in Section 3. Section 4 presents opportunities in datacenter design and resource management techniques to harness varying grid dynamics for cost and environmental benefits. We summarize the paper and future directions in Section 5.

2 Today’s Renewable-based Power Grid and Price/Carbon Variation

Traditionally, power grids dispatch conventional power plants (eg. nuclear, coal, natural gas) to match power supply and changing demand. These resources form a relatively stable supply curve, and some of them (eg. coal and nuclear) are running nearly all the time to support base demand, producing stable carbon content of generation. In the market clearing process, as the power price is determined by the marginal generation cost of the last power plant dispatched to meet the demand, which is usually gas-fired power plant, the power price is also relatively stable and varies with the demand.

As power grids incorporate increasing renewable generation like wind and solar, their variation and intermittence in power supply have transformed power markets, causing rapid fluctuations in power price [31] and carbon-content more dynamic as shown in Figure 1.

When wind and solar are plentiful, they drive down the price and carbon intensity of power (carbon emissions per

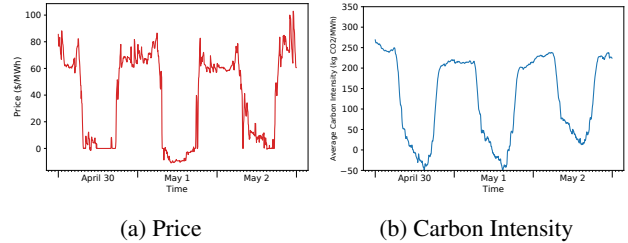


Figure 1: Large variations in Price (\$0–100/MWh) and Carbon Emissions (0–280 kg CO₂/MWh) of power is a feature of decarbonizing power grids (April 30–May 2, 2022, Data Source: CAISO).

MWh). The resulting range is broad, swinging from \$0 to 100/MWh with occasional spikes to \$1000/MWh [4]. Negative prices are also frequent [7]. Carbon intensity swings can be wide, ranging from 0 to 1000 kg CO₂/MWh. The example in Figure 1 shows California Independent System Operator’s (CAISO) time-varying average power price and carbon intensity from April 30 to May 2, 2022. On April 30, CAISO momentarily achieved 100% renewable generation for the first time [18]! As a result, the carbon-emission rate dropped to zero (actually slightly negative at -50 kg CO₂/MWh due to power export), and the power price fell to near zero. These trends toward greater variation and volatility are increasing.



Figure 2: States and Countries have ambitious renewable portfolio goals, adding renewable generation that increases variation in future power grids.

Many countries or regions have proposed ambitious renewable portfolio standard (RPS, renewable fraction in power generation) mandates/goals for the near future, many growing to more than 50 or even 100% (Figure 2); the addition of renewables to meet these goals will drive variation in power price and carbon intensity far beyond that in Figure 1.

This growing divergence of a modern renewable grid when compared with the past stable, unvarying power grid creates new opportunities. Not only can dynamic properties of the power grid be exploited to reduce carbon, but also to reduce cost. Current datacenter design and resource management metrics and approaches do not reflect this opportunity.

3 Traditional Datacenter Optimization

Total-cost-of-ownership (TCO) models provide a framework to think about datacenter optimization, and several variants [1, 3, 15] are widely recognized and used in research and commercial practice. These models share several key elements.

First, traditional cost-optimization of datacenters uses the TCO model, consisting of capital expenses (capex) and operational expenses (opex). Capex includes costs like datacenter construction, land purchase, server purchase, and others, but annualizing them. Opex refers to operational costs such as power, staffing, etc. The model uses an annual average. Thus, TCO optimization is framed as annualized cost minimization.

Further, these costs are based on fixed maximum power consumption (eg. 40MW, 200MW, etc.) with a high availability target (eg. six "9's" or a less than a minute per year outage). A classic example is shown in Figure 3, where capex (server amortization, server interest, DC amortization, and DC interest) are combined with opex (dc opex, pue overhead, server power, and server opex) in a pie chart for TCO. These annualized averages do not capture temporal variation of costs on an hourly or daily basis that can occur in a renewable-based grid. And a fixed target for availability fails to capture the continua of partial availability, temporal extent, and more.

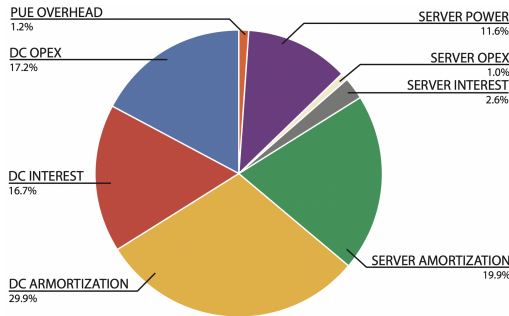
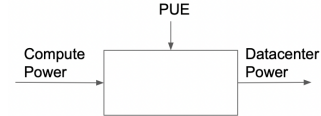


Figure 3: Annualized Capex and Opex models are combined for overall TCO (from [1] Case B)

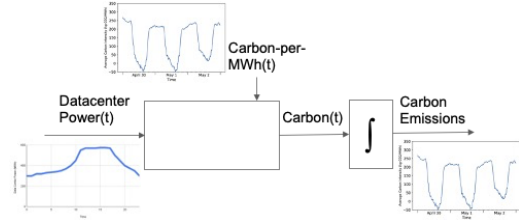
Second, a critical focus on energy efficiency of the entire system, including power distribution and cooling infrastructure has driven use of the power usage effectiveness (PUE) measure. PUE captures the power overhead of cooling, power distribution and more relative to the power supplied to the computing equipment. The average PUE multiplied by the computing equipment power reflects the datacenter power.

$$PUE = \frac{\text{Facility Power}}{\text{IT Equipment Power}} = \frac{\text{Datacenter Power}}{\text{Compute Power}}$$

Some analyses [1] treat PUE as if it were a static property of a datacenter, facility PUE as shown in Figure 4a. But as pointed out in [1], a more accurate measure is "instantaneous PUE", which captures PUE's variation as a function of datacenter occupancy, load, and external weather.



(a) Facility PUE Model



(b) Dynamic Power and Carbon Emissions

Figure 4: Shifting to Carbon metrics requires time-dependent metrics in a variable Carbon Emissions grid.

Further, to compute carbon emissions, even a nearly constant datacenter power consumption can translate into rapidly varying carbon emissions due to power grid variation in generation mix – swinging from renewables to fossil-fuels and back again. In short, average PUE does not capture this dynamism. This idea is illustrated in Figure 4b.

New measures are needed that capture power use efficiency at an operating power, and capture carbon-emissions implied as a function of time. This is critical because future datacenters will vary their power significantly over time to exploit both cheap power and lower-carbon power.

4 Opportunity: Flexing to Reduce Carbon

The variation of the modern renewable-based power grid means that the carbon-content of power (and thereby computing) depends on the time of consumption. For datacenters the opportunity is to flex their demand and properties, so that they consume (and do) more when grid conditions are favorable (low carbon, low price). This reduces their "weighted average" of carbon emissions per unit compute. However, flexing requires datacenter capacity flexibility; for example, as in Zero-carbon Cloud [36].

To exploit varying datacenter hardware capacity, we need flexible workloads and resource management. Ideally, they would have service-level objectives (SLOs) that allow computation at varying speed or perhaps delayed to a later time. By "shifting workload" to exploit power grid opportunity, we can reduce carbon emissions for a workload. We explore hardware and software challenges – both infrastructure and applications – to both create and exploit variable capacity.¹

¹One can think about this as pure upside – increased intermittent capacity – or partially controllable downside – outage. Both framings appear useful.

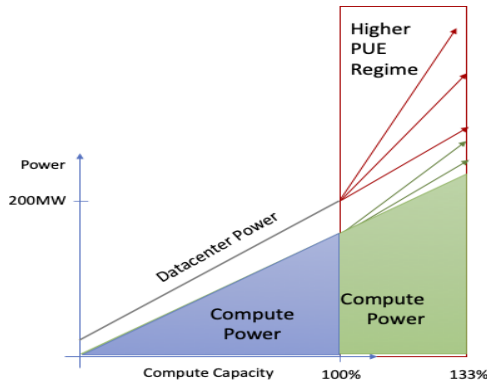


Figure 5: Compute Power, PUE, and Resulting Datacenter Power Consumption

4.1 Opportunity 1: Flexing Datacenters for Cost and Carbon

As introduced in Section 3, we consider the instantaneous PUE of a datacenter at different power levels as in Figure 5. On the left, the blue region indicates power consumed by compute hardware, the line above the total datacenter power, and the ratio between the two, the instantaneous PUE. A traditional datacenter corresponds to this left region.

Consider a new kind of flexible datacenter that has an over-powered mode – a bit like turbo or overclocking – where the datacenter could run its cooling equipment, additional fans, higher coolant flow rates, etc. to remove even more heat, at a cost in power-efficiency, reflecting a higher PUE, and a bigger gap between the colored region and the steeper red lines. Such a datacenter operates in "blue" mode when power is high-carbon, and in "blue+green" mode when power is low-carbon. This variable capacity allows the datacenter to weighted-average reduce power or carbon cost.

Applying this idea to a variable-capacity datacenter can reduce power cost or carbon at fixed capacity or boost capacity under the same budget as shown in Figure 6, using data from the German power grid [11]. And because in many power markets, power price generally covaries with carbon-emissions content, both carbon and cost benefits can be achieved with this same technique. Hence, a flexible datacenter lowers of Opex as well as carbon cost (per unit compute).

The challenge for flexible datacenters is the additional Capex. An operator’s point of view is – "I paid for it already, so use it always". To make flexing cost-effective, its important to minimize the capital cost of the greater compute capacity with innovative design. Ideally, the goal is to increase compute/cooling/power distribution capacity with a sublinear increase in Capex. Here are some potential approaches:

1. Turbo cooling with higher flow rates (or immersion), more pump power per watt removed [17]
2. Run chilled water at lower temperature, decreasing effi-

- ciency but increasing capacity (or run blowers faster) [9]
3. Turbo mode compute [17, 34]
4. Increase power distribution capacity, but with lower redundancy, build for dynamic range [39]
5. Reduce capacity of power, cooling, systems, and don’t use flexibility on hottest days [21]

The first three approaches reflect ways to increase capacity, but the increase is achieved at increased PUE (lower efficiency) in the first two, and lower compute power efficiency in the third. This matches well with the ability to exploit cheap or low-carbon power. The latter two provide increased flexible capacity, but at lower reliability acceptable because this is "bonus" capacity. This lower reliability reduces cost per unit compute. All reverse conventional wisdom for optimizing PUE. Finally, all five approaches reduce the Capex/unit delivered compute. Interesting research problems include:

- What slopes of PUE are possible? What ranges of flexibility are possible? at what Capex/unit cost?
- What is the right granularity to build flexibility? (10, 40, 200 MW units?)
- What datacenter designs (flexibility, headroom) are cost-appropriate for different power grid settings? (variation opportunity) for regional workloads? (flexibility).
- Are there other approaches that can increase flexibility? (heat/cold storage, precooling, etc.)

More generally, design for the "flexible datacenter" challenge involves re-examining the full breadth of computing hardware design, network, cooling, building, etc. [1], in a dynamic (time-varying), region-custom (grid power) setting.

4.2 Opportunity 2: Flexing to Exploit Variable Capacity

Once beyond the notion of a static datacenter capacity, the challenge of resource management changes [38]. The objective is now both filling a variable capacity envelope (high utilization), and to push the envelope to increase carbon or price efficiency (see Figure 6, variation from German power grid data). We adopt the metric of goodput, which represents the system throughput of useful work, capturing both utilization and increased capacity under budget.

Resource managers today deal with uncertainty in load, but generally assume static or slowly changing capacity. In such models, the current capacity will continue, so resource managers can commit resources far into the future and simply strive to optimize system utilization. However, in a renewable-powered grid, both power price and carbon intensity can change violently in response to changes in grid load and renewable generation. In this case, intelligent use of budget can increase available capacity, and hence, significantly improve

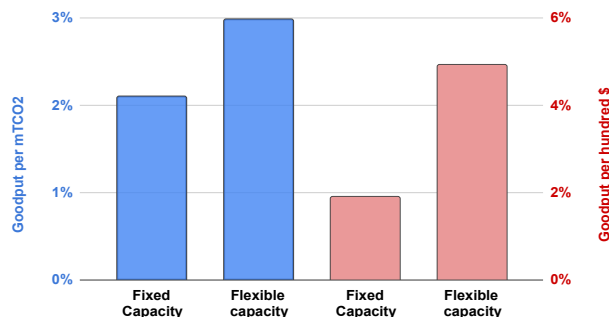
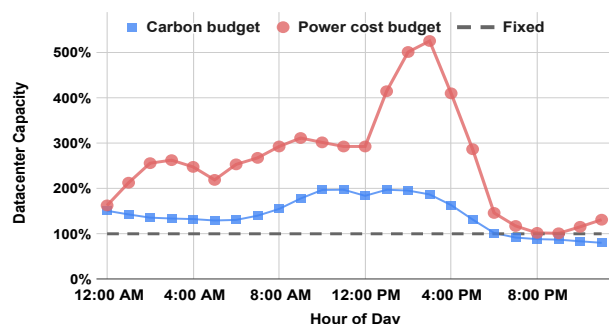


Figure 6: Capacity determined by Carbon or Power Cost Budget (left) and How Flexibility Increases Capacity/carbon or /\$ (right)

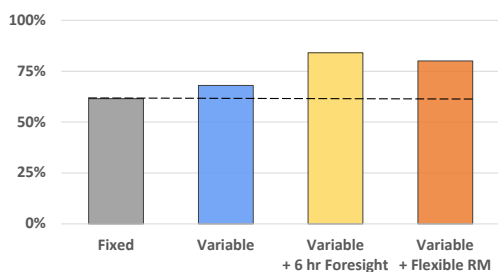


Figure 7: Flexible Resource Managers combined with Flexible Capacity Datacenters can increase goodput significantly

compute efficiency per unit of carbon emission or monetary cost. Figure 6 demonstrates such benefits by simulations of historical workload traces on Mira Supercomputer [28].

With a carbon budget, the capacity available to a resource manager can change fast – relative to the length of resource commitments – and unpredictably [38]. Without improved resource management, long-running jobs could be subject to termination or preemption, violating the service-level objective (SLO). If resource capacity increases, the scheduler may have underestimated and missed an opportunity to achieve more goodput under budget. Therefore, one of the resource management challenges for flexible datacenters is the mismatch between rapid change in capacity and resource managers’ assumption into the future.

One way to reduce the mismatch is to provide schedulers with information about the future (e.g., prediction). For example, given foresight of capacity variation, resource managers can reshape workloads to achieve high utilization with variable capacity under the same budget (Figure 7, Mira Supercomputer). Ideally, to maximize system performance with variable capacity, resource managers must close the gap between future capacity estimates and actual capacity changes. This empowers them to minimize negative SLO impact. Potential approaches include:

1. Accurate prediction of capacity (foresight) to enable better resource management [24, 29]

2. Increase workload flexibility with geo-distributed load shifting [19]
3. Capture workload resource requirements over time (eg. advertised max runtime, variation in memory or cpu use) [6, 8]
4. Intelligent intra-datacenter power allocation and management amongst cells, racks, and servers [35]
5. New service models capturing temporal application flexibility (eg. SLO, performance) [20, 33]
6. Invent new techniques for workload reliability through failures [37]

The first four approaches empower resource management ideas to optimize system utilization at the right times. The latter two focus on workloads providing more information to enable flexibility to resource managers as they tolerate variations (subject to SLOs) or even failures. All enable resource managers to better exploit variable capacity to reduce carbon or power cost. Interesting Research Problems:

- How to determine the right power level? Should the power grid have input?
- How to distribute power changes across applications? hardware? What is fair? How does this limit the ability to adapt to variation?
- What new service or workload models are best? For different variation? Reliability? User experience?

5 Summary

Climate change and the radical transformation of the power grid have forced dynamic management onto cloud datacenters. We highlight several exciting new research directions, and we suspect, the change has only begun. It is too early to tell how these vectors of flexibility and dynamism will evolve, but we are excited to pursue these growing opportunities.

Acknowledgments

Work supported in part by NSF Grants CMMI-1832230, OAC-2019506, and the VMware University Research Fund Thanks Zero-carbon Cloud group!

References

- [1] Luiz André Barroso, Urs Hölzle, and Parthasarathy Ranganathan. The datacenter as a computer: Designing warehouse-scale machines. *Synthesis Lectures on Computer Architecture*, 13(3):i–189, 2018.
- [2] Luiz André Barroso. A brief history of warehouse-scale computing. *IEEE Micro*, 41(2):78–83, 2021.
- [3] Josep L Berral, Ínigo Goiri, Thu D Nguyen, Ricard Gavaldà, Jordi Torres, and Ricardo Bianchini. Building green cloud services at low cost. In *2014 IEEE 34th International Conference on Distributed Computing Systems*, pages 449–460. IEEE, 2014.
- [4] Joshua W Busby, Kyri Baker, Morgan D Bazilian, Alex Q Gilbert, Emily Grubert, Varun Rai, Joshua D Rhodes, Sarang Shidore, Caitlin A Smith, and Michael E Webber. Cascading risks: Understanding the 2021 winter blackout in texas. *Energy Research & Social Science*, 77:102106, 2021.
- [5] CAISO. Market processes and products, 2022. Retrieved from <http://www.caiso.com/market/Pages/MarketProcesses.aspx>.
- [6] Abhishek Chandra, Weibo Gong, and Prashant Shenoy. Dynamic resource allocation for shared data centers using online measurements. In *International Workshop on Quality of Service*, pages 381–398. Springer, 2003.
- [7] Andrew A. Chien, Fan Yang, and Chaojie Zhang. Characterizing curtailed and uneconomic renewable power in the mid-continent independent system operator. *AIMS Energy*, 6(2):376–401, December 2018.
- [8] Eli Cortez, Anand Bonde, Alexandre Muzio, Mark Russinovich, Marcus Fontoura, and Ricardo Bianchini. Resource central: Understanding and predicting workloads for improved resource management in large cloud platforms. In *Proceedings of the 26th Symposium on Operating Systems Principles*, pages 153–167, 2017.
- [9] Schneider Electric. Schneider high performance uni-flair chillers, 2022. Describes how warmer water temperatures increase energy efficiency. Retrieved from <https://www.se.com/us/en/>.
- [10] Mary Jo Foley. Cloud revenues power microsoft’s \$51.7 billion q2 in fiscal year 2022, January 2022. Retrieved from <https://www.zdnet.com/article/microsoft-cloud-revenues-power-microsofts-51-7-billion-second-fy22-quarter/>.
- [11] Electricity market data: Generations, prices, power, 2020. Retrieved from <https://www.smard.de>.
- [12] Google. About google datacenters, 2022. Retrieved from <https://www.google.com/about/datacenters/>.
- [13] Greenpeace. Clicking Clean Virginia: The Dirty Energy Powering Data Center Alley, February 2019. Retrieved from <https://www.greenpeace.org/usa/reports/click-clean-virginia/>.
- [14] GWEC. Global wind report: Annual market update. Technical report, Global Wind Energy Council, 2016. Documents curtailment around the world.
- [15] J. Hamilton. Overall data center costs, (2010). Retrieved May 12, 2022 from <https://perspectives.mvdirona.com/2010/09/overall-data-center-costs/>.
- [16] M Huneault and FD Galiana. A survey of the optimal power flow literature. *IEEE transactions on Power Systems*, 6(2):762–770, 1991.
- [17] Majid Jalili, Ioannis Manousakis, Ínigo Goiri, Pulkit A. Misra, Ashish Raniwala, Husam Alissa, Bharath Ramakrishnan, Phillip Tuma, Christian Belady, Marcus Fontoura, and Ricardo Bianchini. *Cost-Efficient Overclocking in Immersion-Cooled Datacenters*, page 623–636. IEEE Press, 2021.
- [18] Ryan Kennedy. For the first time in history, california’s demand was 100% matched by renewable energy generation, (2022). Retrieved May 12, 2022 from <https://pv-magazine-usa.com/2022/05/02/for-the-first-time-in-history-california-was-100-powered-by-renewable-energy/>.
- [19] Julia Lindberg, Yasmine Abdennadher, Jiaqi Chen, Bernard C Lesieutre, and Line Roald. A guide to reducing carbon emissions through data center geographical load shifting. In *Proceedings of the Twelfth ACM International Conference on Future Energy Systems*, pages 430–436, 2021.
- [20] Michel J Litzkow, Miron Livny, and Matt W Mutka. Condor—a hunter of idle workstations. Technical report, University of Wisconsin-Madison Department of Computer Sciences, 1987.
- [21] Lancium LLC. Lancium compute pricing and qos structure, 2022. Includes availability levels, ranging from 90 pct down to 25 pct. Retrieved from <https://portal.lancium.com/about/pricing>.
- [22] Ioannis Manousakis, Ínigo Goiri, Sriram Sankar, Thu D Nguyen, and Ricardo Bianchini. Coolprovision: Underprovisioning datacenter cooling. In *Proceedings of the Sixth ACM Symposium on Cloud Computing*, pages 356–367, 2015.

- [23] Microsoft Azure. Azure global datacenters, 2022. Retrieved from <https://azure.microsoft.com/en-us/global-infrastructure/>.
- [24] Jakub Nowotarski and Rafał Weron. Recent advances in electricity price forecasting: A review of probabilistic forecasting. *Renewable and Sustainable Energy Reviews*, 81:1548–1568, 2018.
- [25] U.S. Department of Energy. Solar futures study. Technical report, September 2021. Retrieved from <https://www.energy.gov/eere/solar/solar-futures-study>.
- [26] Ehsan Pakbaznia and Massoud Pedram. Minimizing data center cooling and server power costs. In *Proceedings of the 2009 ACM/IEEE international symposium on Low power electronics and design*, pages 145–150, 2009.
- [27] Sundar Pichai. Our biggest renewable energy purchase ever, 2019. Retrieved from <https://blog.google/outreach-initiatives/sustainability/our-biggest-renewable-energy-purchase-ever/>.
- [28] Argonne Leadership Computing Facility. Mira supercomputer, 2022. Retrieved from <https://www.alcf.anl.gov/mira>.
- [29] Ana Radovanovic, Ross Koningstein, Ian Schneider, Bokan Chen, Alexandre Duarte, Binz Roy, Diyue Xiao, Maya Haridasan, Patrick Hung, Nick Care, et al. Carbon-aware computing for datacenters. *arXiv preprint arXiv:2106.11750*, 2021.
- [30] John Roach. Microsoft’s virtual datacenter grounds ‘the cloud’ in reality, April 2021. Microsoft to build 50 to 100 datacenters per year. Retrieved from <https://news.microsoft.com/innovation-stories/microsofts-virtual-datacenter-grounds-the-cloud-in-reality/>.
- [31] Joachim Seel, Andrew D Mills, and Ryan H Wiser. Impacts of high variable renewable energy futures on wholesale electricity prices, and on electric-sector decision making. 2018.
- [32] Amazon Web Services. Amazon global datacenters, 2022. Retrieved from <https://aws.amazon.com/about-aws/global-infrastructure/>.
- [33] Amazon Web Services. Amazon spot instances. Retrieved from <https://aws.amazon.com/ec2/spot/>, 2022.
- [34] Wikipedia. Overclocking — Wikipedia, the free encyclopedia. (2022). Retrieved May 12, 2022 from <https://en.wikipedia.org/wiki/Overclocking>.
- [35] Qiang Wu, Qingyuan Deng, Lakshmi Ganesh, Chang-Hong Hsu, Yun Jin, Sanjeev Kumar, Bin Li, Justin Meza, and Yee Jiun Song. Dynamo: Facebook’s data center-wide power management system. *ACM SIGARCH Computer Architecture News*, 44(3):469–480, 2016.
- [36] Fan Yang and Andrew A. Chien. Large-scale and extreme-scale computing with stranded green power: Opportunities and costs. *IEEE Transactions on Parallel and Distributed Systems*, 29(5):1103–1116, 2018.
- [37] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauly, Michael J Franklin, Scott Shenker, and Ion Stoica. Resilient distributed datasets: A {Fault-Tolerant} abstraction for {In-Memory} cluster computing. In *9th USENIX Symposium on Networked Systems Design and Implementation (NSDI 12)*, pages 15–28, 2012.
- [38] Chaojie Zhang and Andrew A. Chien. Scheduling challenges for variable capacity resources. In *Job Scheduling Strategies for Parallel Processing: 24th International Workshop, JSSPP 2021, Virtual Event, May 21, 2021, Revised Selected Papers*, page 190–209, Berlin, Heidelberg, 2021. Springer-Verlag.
- [39] Chaojie Zhang, Alok Gautam Kumbhare, Ioannis Manousakis, Deli Zhang, Pulkit A. Misra, Rod Assis, Kyle Woolcock, Nithish Mahalingam, Brijesh Warrior, David Gauthier, Lalu Kunnath, Steve Solomon, Osvaldo Morales, Marcus Fontoura, and Ricardo Bianchini. *Flex: High-Availability Datacenters with Zero Reserved Power*, page 319–332. IEEE Press, 2021.