

Grid Integration of AI Data Centers: A Brief Review of Energy Storage Systems Solutions

Sina Mohammadi, Marcus Chen I Wada,
Rouzbeh Haghighi, Ali Hassan, Hualong Liu,
Wencong Su
University of Michigan-Dearborn
Dearborn, United States
sinamo, wadamc, haghighi, alihssn, wencong@umich.edu

Wayne Wang, Archit Bhatnagar, Ang Chen
University of Michigan-Ann Arbor
Ann Arbor, United States
wswang, architb, chenang@umich.edu

ABSTRACT

The proliferation of artificial intelligence (AI) workloads has necessitated AI-specific data centers (DCs), which have stochastic power profiles and new reliability constraints. Unlike traditional DCs, AI DCs exhibit high-frequency power fluctuations and significant ramping events, because of the tight coupling of high-density and magnitude compute, thermal management, and power electronics. This paper studies the role of energy storage systems (ESS) and provides an analysis of multi-layer ESS architectures that can help mitigate grid-integration challenges. We take a hierarchical view that spans from grid-scale BESS and grid-interactive UPS down to rack-level units and server/GPU-level energy buffers. The analysis reveals the distinct operating timescales and control coordination required for heterogeneous storage integration. Furthermore, we evaluate grid-support functionalities, such as frequency regulation, while addressing the techno-economic trade-offs of degradation and reliability. Our review of ESS capabilities provides roadmaps for designing resilient, sustainable, and grid-compatible AI DCs.

KEYWORDS

AI Data Center, Battery Backup Unit (BBU), Battery Energy Storage Systems (BESS), GPU, UPS

1 INTRODUCTION

Artificial Intelligence (AI) has fundamentally altered society's relationship with technology; spanning applications from simple search tasks to new avenues of scientific exploration. Large language models (LLMs) have been at the forefront of AI systems, with models such as GPT and Claude increasingly embedded in everyday use. These large-scale models are trained on trillions of tokens collected from online sources and are accessible from web platforms and APIs. The surge in demand has prompted AI industry stakeholders to design specialized data centers (DCs) capable of supporting the high-computational tasks required for training, fine-tuning, and inference of these complex models. Energy requirements for AI queries are now nearly an order of magnitude greater than those of Google search, while AI training workloads double approximately every 3.4 months [7]. Collectively, these trends highlight the necessity for rethinking data center architectures to sustainably support the rapid scaling of AI workloads.

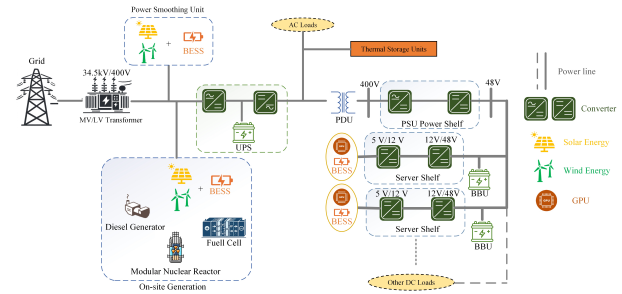


Figure 1: An overview of AI DC structure including hierarchal ESSs

DCs are dedicated facilities that are designed to house compute servers, network devices, and data storage infrastructure. Energy demand within a DC is unevenly distributed: IT equipment and cooling systems account for the majority of electricity consumption, while auxiliary loads such as lighting and security represent a comparatively small fraction [1]. Traditional data centers (TDCs) have historically operated at total power levels typically below 30 MW [1], with individual electrical distribution branches commonly limited to a few tens of kilowatts. These facilities are predominantly connected to power distribution networks in densely populated areas and operate under an architectural paradigm that emphasizes redundancy and exceptionally high uptime, making uninterrupted power delivery from the utility grid essential. In contrast, AI DCs (especially for training) comprise of high-density computing units that usually operate in the hundreds of megawatt power range. Under these conditions, AI DCs are better suited for connection to transmission networks to enable reliable grid integration and reduced grid impact. Inference AI DCs operate at comparatively lower power levels with fewer restrictions regarding connecting them to transmission networks.

Beyond the differences in scale and infrastructure, AI DCs exhibit two distinct load profiles (training and inference) that differ fundamentally from those of TDCs. AI training workloads are characterized by rapid power fluctuations arising from idle periods, peak utilization events, communication pauses, and checkpointing operations, with significant load variations occurring on sub-second timescales, producing high-frequency, jitter-like power spikes. In contrast, inference workloads generally lack such rapid transients and exhibit comparatively smoother demand profiles. Traditional data center loads, by comparison, behave more like conventional grid loads and are typically more predictable using historical data and standard forecasting methods.

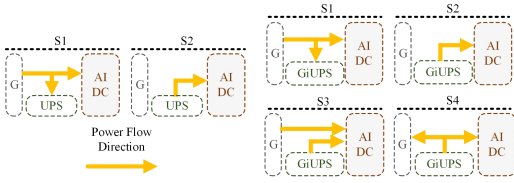


Figure 2: UPS and GiUPS operation modes

The highly dynamic and irregular nature of AI DC load profiles has important implications for grid operation. Rapid demand changes can challenge the maintenance of power quality and system stability by inducing voltage fluctuations, increasing harmonic distortion through grid-connected power electronic interfaces, and imposing additional stress on supporting electrical infrastructure, particularly at high penetration levels. As a result, load smoothing becomes a critical requirement for AI DC operation, both to improve load predictability and to mitigate adverse power-quality impacts on the power system.

The combination of high power density, rapid load transitions, and limited predictability challenges the maintenance of grid reliability, stability, and power quality using traditional infrastructure alone. In this context, energy storage systems (ESSs) emerge as a critical enabling resource for buffering fast transients and enhancing operational flexibility. Traditionally, ESSs have existed in various forms in DCs mainly for the purpose of supplying power for the DC during power outages, to ensure the continuing operation of computing devices and sustain a high up-time. Stored energy, if well leveraged, can address many of the aforementioned challenges for AI DCs. However, existing studies largely focus on isolated applications such as backup power or peak shaving, lacking a comprehensive framework that addresses ESS technologies, architectural placement, and coordinated grid-operational roles specific to AI DCs.

In this paper, we address this gap and investigate the role of ESSs in AI DCs, not only as backup units but also as on-site cooling support via thermal ESSs; battery backup units (BBUs) inside IT racks; GPU-level capacitors for chip-level load smoothing; battery energy storage systems (BESSs) for data-center-level load smoothing; and even the interaction of ESSs inside the UPS with the utility grid for supporting services such as frequency regulation or peak shaving. In contrast to [8, 24], this paper focuses on the role of ESSs in enabling reliable grid integration of AI DCs, rather than only characterizing the grid impacts of AI DCs. Moreover, our work does not focus on ESS chemistry technologies [43, 45], but instead attempts to study architectures for the optimal application of ESSs in AI DC grid integration. This paper also provides insights into the control and coordination of ESSs across different levels of AI DC architecture to address existing challenges. Therefore, we hope that our work will serve as a useful source of information for both academic and industry partners that study the role of ESSs in future AI DCs.

2 THE ROLE OF UPS SYSTEMS

2.1 Traditional UPS Operation and Limitations

UPS systems are fundamental components in AI DCs, ensuring continuous operation of critical loads during grid disturbances or

outages. Conventional UPS systems [22, 29, 53] typically operate in standby or online modes. In standby mode, the UPS remains bypassed during normal operation and transfers to battery-supported inverter operation only when a disturbance is detected. In online mode, the rectifier and inverter continuously process power, providing isolation from grid disturbances and maintaining tight voltage and frequency regulation at the critical bus (see Fig. 2). In both modes, the UPS interacts with the grid only through battery charging and synchronization. The UPS is therefore treated as a passive protective device.

2.2 GiUPS systems in AI DCs

Traditional UPS architectures were designed primarily for load protection with minimal interaction with the utility grid. This paradigm is increasingly insufficient for AI DCs. First, AI workloads introduce rapid load variations that propagate upstream and cause voltage flicker, transformer thermal stress, and feeder congestion [38]. These effects are amplified when multiple AI racks operate in synchrony, which is common in distributed training. Second, the growing penetration of renewable energy sources (RESs) reduces system inertia and increases the rate of change of grid frequency. Conventional UPS systems do not contribute to frequency support beyond simple ride-through, which leaves a large amount of inverter capacity unused during normal operation [13]. Third, the energy stored in UPS batteries remains idle for more than 99 percent of their lifetime, representing a significant opportunity cost. These limitations motivate the development of UPS architectures capable of dynamic grid interaction.

GiUPS systems enhance traditional UPS designs by incorporating bidirectional power converters, advanced digital controllers, and communication interfaces that allow coordination with the utility grid or market operator. These systems can modulate active power in response to grid frequency deviations or dispatch signals, enabling participation in services such as primary frequency regulation, fast frequency response, synthetic inertia, and peak load management [16, 39]. Unlike conventional UPS systems, which operate independently of grid conditions except during outages, GiUPS systems continuously monitor grid conditions and adjust their output accordingly.

The control architecture of GiUPS typically includes a multi-layer structure. The inner control loops regulate current and voltage at sub-millisecond timescales. The outer loops implement grid-following or grid-forming behavior depending on the operating mode. A supervisory controller coordinates battery SoC, grid service participation, and critical load protection. This layered structure allows the UPS to respond to frequency deviations within tens of milliseconds, which is significantly faster than most synchronous generators and comparable to dedicated BESSs [49].

GiUPS can provide fast frequency response within 0.5 to 10 seconds after a grid event. It can operate in a dynamic frequency response mode. In this mode, the injected battery power is adjusted in real time based on frequency regulation needs. The UPS can also operate in a static frequency response mode. In this case, a fixed amount of power is injected after a grid event such as a fault or large load switching [17].

2.2.1 Operational Modes of GiUPS. GiUPS can operate in several modes depending on the grid conditions and external control signals, including *standard operation*, *discharge modes* (full and partial disconnection), *recharge mode*, and *energy export (bi-directional) mode* [16, 20]. When a frequency variation is detected by an external controller, the UPS adjusts its operation to provide both positive and negative regulation by charging or discharging its batteries within operational limits.

Under normal conditions, the UPS receives 100% of the input power from the utility grid through the rectifier, delivering it to the load (see S1 in Fig. 2). In this mode, the UPS functions as a standard double-conversion system, with the batteries maintained in standby and not actively supplying power. This operation is referred to as the *standard UPS operation*.

Discharge modes are employed to meet external control requests, supplying energy to the load or the grid. In the *full disconnection mode*, the UPS is completely disconnected from the utility, and 100% of the load power is supplied by the batteries (see S2 in Fig. 2). In the *partial disconnection mode*, the input power from the grid is reduced according to external commands, with the remaining load power supplied by the batteries (see S3 in Fig. 2). For example, the batteries may provide 25% of the total load power while the grid supplies the remainder.

When the battery SOC is below 100% (e.g., 80%) and an over-frequency event is detected, the UPS draws power from the grid to recharge the batteries. This operation is referred to as *recharge mode*. The maximum recharge power is typically limited to 20–25% of the UPS nominal capacity, constrained by battery recharge characteristics and the maximum input current.

The UPS can also operate as a bi-directional power converter to inject energy back into the grid (see S4 in Fig. 2). This *energy export mode* is subject to local regulatory and grid requirements and allows the UPS to discharge batteries upstream.

2.2.2 GiUPS Topologies. For GiUPS applications, topologies that support bidirectional power flow and high-speed digital control are preferred. Double-conversion and delta conversion architectures with bidirectional converters can rapidly transition between load-following, grid-support, and islanded operation. Modular architectures are particularly attractive in AI DCs because they allow some modules to participate in grid services while others remain dedicated to critical load protection. This partitioning reduces operational risk and improves economic performance [16]. Emerging DC-coupled architectures further reduce conversion stages and improve round-trip efficiency. These architectures also enable direct coupling with on-site co-generation, which enhances the ability of the AI DCs to operate as controllable grid resource. Modular UPS design presents significant research opportunities not only for frequency regulation and peak shaving scenarios, but also for enabling demand response programs and synthetic inertia generation. However, detailed studies are required on voltage ride-through capabilities, black-start mechanisms, and battery degradation analysis.

3 THE ROLE OF BESS SYSTEMS

BESS integrated with DCs represent a practical and scalable solution to address the challenges imposed by large and highly dynamic

loads on the power system [9]. When appropriately designed and controlled, BESS can play a critical role in mitigating large load fluctuations, enhancing local power quality, and supporting overall grid stability [36]. To establish a clear understanding of the underlying mechanisms and achievable benefits, this section reviews the key grid-support functions enabled by AI DC-connected BESS, including co-generation, power smoothing and frequency regulation. Finally, the cost implications of BESS deployment are examined.

3.1 Grid Support Roles of BESS in AI DCs

3.1.1 On-Site Clean Power Integration. BESS play a critical role in enabling RESs integration at AI DCs. Given the energy-intensive nature of AI DCs, on-site RESs generation and off-site procurement mechanisms such as power purchase agreements are increasingly adopted to reduce operating costs and meet sustainability objectives. However, the inherent intermittency of RESs complicates reliable power supply. BESS mitigate these challenges by storing excess energy during periods of high generation and supplying it during low availability, while also providing backup power, and improved RESs utilization [3]. AI DC-integrated BESS can be architected either as large, multi-megawatt grid-forming resources capable of replacing conventional backup generators, or as large-capacity batteries integrated within static UPS systems [54]. In this context, modern UPS systems, traditionally designed for short-duration ride-through, are increasingly complemented by dedicated BESS that are optimized for extended energy management and grid-interactive operation rather than transient backup alone [19]. At the converter-control level, BESS can be configured for grid forming (GFM), grid following (GFL), or hybrid GFM and GFL operation to balance fast dynamic response, voltage/frequency support, and grid-service capability, although with added control and sensing complexity.

3.1.2 Power Smoothing. BESS can inject or absorb power at the facility interface to attenuate rapid AI DC demand variations and present a smoother power profile to the grid. A BESS co-located with an AI DC can provide grid-level power smoothing while supplying sufficient energy capacity to support sustained mitigation actions. This architecture improves operational reliability by enabling load shaping, demand flexibility, and extended ride-through for AI workloads.

As shown in Fig. 3, the GFM BESS effectively smooths the power fluctuations of the AI DC resulting from training and inference jobs. An alternative approach is a hybrid E-STATCOM and BESS configuration, which combines the fast, high-speed response of supercapacitors with the higher energy capacity of BESS to achieve effective smoothing and sustained load support. This hybrid architecture can enhance resilience by improving demand flexibility, reducing flicker, and enabling broader grid-service capability. The primary trade-off is increased cost and footprint compared with standalone solutions. In addition, standalone supercapacitor-based E-STATCOM solutions can mitigate high-frequency AI load fluctuations and improve power quality for sensitive IT equipment by reducing flicker, although they provide limited long-duration energy support [42].

3.1.3 UPS-Integrated BESS. UPS-integrated BESS are motivated by a simple premise: AI DCs already deploy UPS systems with fast

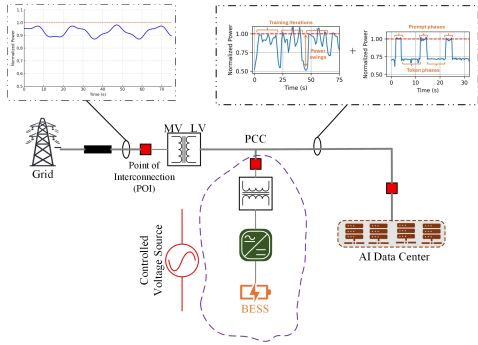


Figure 3: GFM BESS as a power smoothing unit [42]

power-electronic interfaces, and augmenting this mandatory infrastructure with BESS enables grid-interactive capability with limited additional integration complexity. In practice, grid-interactive deployments also show that battery technology selection depends on the targeted services and operating constraints. For example, UPS-integrated BESS have been used to provide fast frequency response (FFR) by leveraging the rapid controllability of UPS power converters [50, 54].

The UPS-integrated BESS, labeled as *power smoothing units* in Fig. 1, extend this concept by supporting long-duration outage operation and enabling participation in demand response programs, while the UPS maintains continuity of critical loads when neither the grid nor the BESS can fully supply demand. Although this architecture increases installation cost, it can significantly improve overall system reliability and may provide additional benefits, including reduced resource consumption, lower carbon emissions, and enhanced grid resilience [39].

As summarized in Fig. 4, a GiUPS-integrated BESS can operate in (i) standard mode, (ii) discharge mode under full or partial disconnection, (iii) grid-services mode, and (iv) recharge mode for SoC recovery [16]. This architecture can reduce the need for bidirectional power flow across the MV/LV distribution transformer by localizing fast power exchanges behind the meter, allowing the transformer to operate primarily under its normal unidirectional loading conditions [34]. These operating modes within the UPS-BESS subsystem enable long-duration BESS to enhance DC resilience through extended backup, reduce or eliminate reliance on diesel generation, participate in ancillary and reserve markets, manage demand charges and time-of-use tariffs, and increase RES utilization [19].

S1) Normal operation: The AI DC load is supplied entirely from the utility grid, and the system delivers 100% of the input power to the AI DC.

S2) Full disconnection with sufficient BESS reserve: The site is fully isolated from the utility, and upon command from an external controller, the BESS supplies 100% of the AI DC load power. In this case, long-duration, large-scale BESS becomes essential; salt-cavern redox flow batteries are a promising due to their high safety, large storage capacity, stable temperature, and low cost [41].

During the transition from Scenario S1 to S2, the UPS initially supports the load because its response time is on the order of milliseconds, while the BESS typically requires a few seconds to assume full load supply. Once engaged, the BESS can sustain operation for

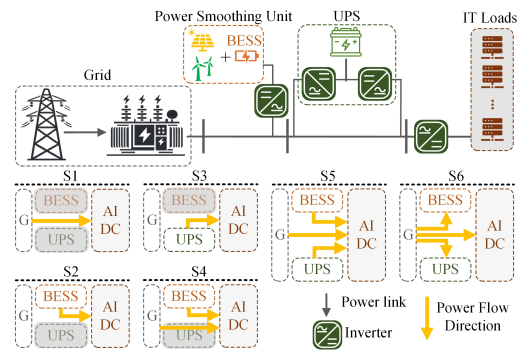


Figure 4: UPS-Integrated BESS operation modes

extended durations (e.g., on the order of 1–4 hours, depending on sizing and operating conditions) [19].

S3) Full disconnection with insufficient BESS reserve: If the outage duration is long and the BESS does not have sufficient available capacity to supply the full AI DC load, the UPS assumes responsibility for supporting the critical load.

S4) Partial disconnection: The external controller reduces the grid input power, and the remaining portion of the AI DC demand is supplied by the BESS. This operating mode enables the BESS to support grid frequency regulation and to present a smoother power profile at the grid interface.

S5) Coordinated UPS-BESS power smoothing: The grid, BESS, and UPS jointly supply the AI DC load, with the UPS helping the BESS by providing fast support during critical transient conditions to assist power smoothing and frequency response. This case can be interpreted as an extension of S4 in which the UPS contributes rapid, short-duration buffering when required.

S6) Recharge mode: When the SOC drops below full charge, the utility initiates battery recharging (e.g., following an over-frequency detection). Under this condition, the UPS is prioritized for recharging first, after which the BESS can be recharged from the grid or from on-site RESs generation.

4 THE ROLE OF RACK-, SERVER-LEVEL ESS

Rack- and server-level ESSs exist to handle two constraints that centralized UPS/BESS alone cannot satisfy: (i) *local buffering* (GPU/server-scale transients that require local energy buffering), and (ii) *fault domain* (limiting the blast radius of power disturbances and energy storage failures). In modern DC-architecture AI facilities, upstream layers (SST/MVDC front-ends and facility UPS/BESS) primarily set facility envelopes and ride-through at longer timescales, while downstream layers move buffering closer to the load to smoothen fast dynamics and localize contingencies under faults [27, 28]. This section focuses on rack-level BBUs as distributed ride-through and backup resources, and on server/GPU-level capacitive storage as power profile smoothening for the highest-frequency behaviors [18, 47, 48].

4.1 Rack-Level ESS: Role, Benefits, and Constraints

The rack represents a natural operational and electrical boundary for AI infrastructure, with dense DC-architecture racks alone reaching MW-scale power demands [27]. Deploying energy storage at

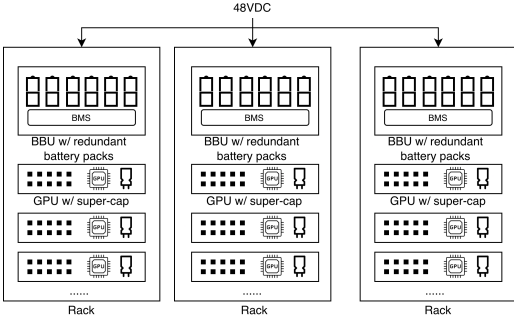


Figure 5: Rack-level BBUs and server-level capacitors

the rack level offers several advantages. First, the storage is electrically close to the load, enabling local buffering before propagating upstream to the facility distribution system. Second, failures in battery, converters, or wiring can be contained to individual racks, limiting the blast radius of faults. Third, storage capacity budgets and ride-through priorities can be configured on a per-rack basis, allowing workload-specific policies that reflect criticality differences across the data center.

Despite these benefits, rack-level deployment introduces operational challenges. Coordination becomes necessary, as many independent storage units must be managed for state-of-charge tracking, recharge scheduling, and availability monitoring. Safety and compliance requirements increase, since colocation of lithium-ion energy storage with IT equipment raises fire and thermal-runaway mitigation concerns. Finally, maintenance overhead grows substantially, as consistent monitoring, battery replacement, and preventive maintenance must be performed across racks. Fig. 5 shows the implementation of rack-level BBUs and server-level capacitors in the IT equipment.

4.2 Battery Backup Units (BBUs) as Distributed DC UPS

In OCP Open Rack designs, BBUs are integrated into the rack as distributed DC ride-through modules that supply the rack DC bus during upstream AC disturbances [48]. In Open Rack V3 (ORV3) architecture, a BBU shelf typically hosts multiple BBU modules with 5+1 redundancy and supports both *charge mode* and *discharge mode* operations with monitoring of SoC/SoH and maintenance tests [47]. Public reference designs summarize the intended operating point (e.g., per-module backup power on the order of kW for minutes, and lower-power charging over hours), reflecting a design goal of short ride-through and fast recovery rather than continuous power profile smoothing [5]. Meta reports ORV3 BBU shelves designed for minutes of backup, and notes that paired shelves can be used for higher rack power configurations [6].

4.2.1 BBU vs centralized UPS. BBUs and centralized UPS systems can both respond quickly in principle to supply the gap between grid power disturbance or outage and longer-term energy storage systems coming online (e.g., gensets / BESS), but they differ in *where* energy is buffered and *what* the buffer is optimized for. Centralized UPS concentrates energy and power conversion at room/facility scale, simplifying management and enabling grid-interactive use cases when permitted [20, 44]. Rack BBUs push

storage behind fewer conversion stages in DC racks and reduce some double-conversion penalties compared to traditional centralized UPS architectures [46]. Practically, BBUs are best viewed as a distributed energy availability layer for the rack (seconds–minutes), not as the primary solution for sub-second smoothing which is handled more effectively by server/GPU capacitive buffers.

4.2.2 BBUs as power absorbers. BBU shelves inherently include charging control; in normal operation they remain in charge mode and can modulate charging current/power within design limits [48]. At hyperscale, uncontrolled simultaneous recharge can create step increases in aggregate load that trip upstream protection. Production case studies show that coordinated and priority-aware charging of distributed batteries can dramatically reduce recharge power (reported reductions up to ~80%) while meeting recharge constraints [35].

However, using BBUs as a frequent, high-rate smoothing solution can accelerate degradation depending on cycling throughput, C-rate, SoC range, and temperature [14, 15]. Practically, server/GPU capacitors and software solution are better suited to handle high-frequency power fluctuation, and BBUs are better suited for ride-through, controlled recovery, and low-frequency shaping with budgeted cycling.

4.2.3 Safety and deployment considerations. Rack BBUs introduce Li-ion safety considerations near IT equipment. OCP BBU specifications explicitly reference safety and propagation constraints and relevant certification/testing practices [48, 51]. This strongly motivates designs that (i) enforce conservative SoC windows and thermal monitoring, and (ii) reduce unnecessary cycling.

If the design intent is “UPS-like” ride-through for all computing devices, deployment of BBU on every rack is the straightforward approach used in OCP-style architectures [6, 47, 48]. Selective deployment (e.g., BBUs only on the highest-power or highest-priority racks) can reduce capex and coordination overhead, but it changes the availability model: unprotected racks become dependent on workload-level fault tolerance, or upstream ESS (BESS/gensets) without transient protection.

4.3 Server- and GPU-Level Storage

4.3.1 Capacitors instead of batteries. The highest-frequency components of AI load power dynamics are at server-level: rapid power consumption swings on the GPUs. These events are handled by local decoupling and bulk capacitance [9]. Buffering at this layer prevents microsecond-to-millisecond disturbances from propagating to the PSU and rack bus, reducing upstream stress and allowing higher utilization at the rack power infrastructure [32].

Compared with batteries, supercapacitors tolerate extremely frequent charge–discharge cycles with minimal degradation, making them attractive for sub-second smoothing [21, 23]. Industry offerings and proposals include compact supercapacitor banks intended to suppress short spikes that would otherwise require overbuilding upstream power infrastructure. The main constraint is energy density and physical volume: supercapacitors can buffer transients and short bursts, but they cannot provide minutes of ride-through and therefore complement (not replace) rack BBUs and other upstream ESSs [21].

4.3.2 Firmware/software shaping as energy “storage”. A key trend is co-design of electrical buffering with firmware-level power shaping. Vendors explicitly target smoother power draw by controlling ramp rates and limiting transient excursions. For example, NVIDIA describes firmware-controlled ramp-up behavior and a “power burner” mode to manage ramp-down and stabilize facility-level power dynamics during AI job transitions [2]. At the platform level, software interfaces such as NVIDIA’s NVML expose enforced GPU power limits and real-time power telemetry, enabling operators and higher-level controllers to constrain GPU power draw and coordinate device-level consumption with rack- or facility-level power envelopes [37]. It is framed as an industry direction to approach power stabilization as a cross-stack problem that couples hardware energy storage and software control to prevent large excursions and meet grid/facility ramp constraints [9].

4.4 Coordination Across Rack and Server Layers

Coordination should be explicit about objectives and timescales: GPU/server controllers shape fast dynamics under performance constraints; rack BBU controllers manage SoC availability and recovery without creating recharge spikes at facility level; facility controllers enforce facility envelopes and power-quality limits [30, 35]. OCP rack ecosystems already expose the necessary hooks (telemetry, shelf controllers, parallel shelf operation) to implement multi-layer coordination in practice [47, 48]. Under MVDC architectures, server/GPU capacitive buffering can reduce required BBU power bandwidth and cycling, but does not eliminate the seconds–minutes energy role of BBUs for ride-through and controlled recovery [28].

5 CHALLENGES AND FUTURE DIRECTIONS

5.1 Challenges for Integrating ESSs in AI DCs

In the following section the main challenges and gaps are discussed considering multi-level ESSs in AI DCs. As discussed in previous sections, the ESSs not only can support the reliable operation of AI DCs but also can smooth the AI DC grid integration. In addition, it showed that integrating ESSs in different levels in AI DC power infrastructure leads to a grid interactive AI DC that can support the grid in emergency cases like frequency regulation, reactive power compensation or peak shaving scenarios. However, there are still challenges for implementing hierarchical ESSs in AI DCs, which will be discussed in the following sections.

5.1.1 GPU Scheduling. Current GPU-based AI workloads are predominantly executed at the highest supported core frequency, prioritizing peak performance at the expense of energy efficiency. Unlike CPUs, where dynamic voltage and frequency scaling (DVFS) is a mature and widely adopted mechanism, energy-aware frequency regulation for GPUs remains at an early stage of development. [12, 55] Moreover, existing GPU scheduling policies largely optimize for throughput and latency, without explicitly considering the energy consumption characteristics of heterogeneous deep learning tasks. [10, 11] This limitation is further exacerbated by the fact that many AI systems developers possess limited insight into the underlying power system implications. This can accelerate the degradation of ESSs and reveals the lack of control over GPU power profiles, which is essential for power smoothing scenarios.

5.1.2 Load Forecasting in AI DCs. As discussed earlier, the unique characteristics of AI DC power profile is the result of training and inference jobs. These special power profiles are harder for forecast. The traditional loads in power grids have predictable trajectories that can be predicted considering the proper historical data. However, this is not applicable for AI DC load profile. As a result, it can directly affected the sizing and ESS operation modes in AI DCs. For instance, the rack-level BESS or chip-level batteries which are designed for power smoothing applications need the accurate power consumption to operate effectively. In addition, the grid-scale BESSs are also need the accurate load profile of the DC to smooth the load and provide effective grid support.

5.1.3 Advanced ESSs Degradation Modeling and Life Time Prediction. Due to high power variability in AI DCs, more frequent charging and discharging scenarios for ESSs are inevitable. This require a precise degradation analysis to maximize the life time and optimize the ESS operation. As discussed in section 4, the ESS at chip or rack level experiences a very high frequency in different operating conditions to smooth the power profile. In this regard, the online monitoring of the ESS could help for better life time prediction and proactive maintenance operations. As Lithium-ion batteries are more desirable for future AI DCs, the study of cyclic aging is important in both small scale ESSs or large scale BESSs. In addition, the installed SLBESSs in AI DCs are more vulnerable in case of aging and they need more precise motorization especially for larger scale storage units.

5.1.4 Hierarchical ESSs Coordination. Multi-level ESS implementation in AI DCs requires a centralized management system to monitor the operating modes of different ESSs and ensure coordination among them. The on-site grid-scale BESS should operate properly with the GiUPS to enable accurate power sharing, either to efficiently utilize RESs or to support the grid in emergency cases. Moreover, the GiUPS battery size and its switching time are related to the size of the on-site grid-scale BESS. Reliability analysis and cost-comparison analysis are required to determine the optimal combination of both for supplying the AI DC and supporting the power grid.

5.1.5 Optimal Sizing and Cost Considerations. There is little public research on the optimal sizing and cost analysis of ESSs inside DCs. In AI DCs, due to the presence of multi-layer ESSs, sizing and cost analysis are even more critical. For instance, is a grid-scale BESS more cost-effective, or a GiUPS system with high power-density batteries? Alternatively, is it more cost-effective to consider large-scale supercapacitors at the rack level, or to use a combination of BBUs and supercapacitors in load-smoothing scenarios? To answer these questions, multiple analyses should be considered to enable precise comparisons.

5.2 Possible Future Directions

5.2.1 Power-Aware GPU Scheduling. During the training of deep learning models for predicting responses under different configuration settings, it is essential to jointly learn and optimize both throughput and energy consumption. This requires adaptive strategies for local and global batch size scheduling to balance convergence speed with computational efficiency [25]. In addition, GPU devices could be dynamically scheduled across training tasks by

accounting for time-varying electricity prices and the availability of RESs, enabling energy-aware model training [4, 8, 52]. To support such flexibility, elastic resource allocation is a key requirement, allowing the number of GPUs assigned to a training job to be adjusted dynamically during execution without disrupting the learning process. Beyond this, dynamic batch size selection can be implemented to regulate the utility grid voltage profile [33]. The work [33] demonstrates the potential of chip-level utilization to enhance system-level objectives, in which ESSs implementation can also be included in the optimization problem.

5.2.2 AI DC Load Observability. High-resolution monitoring of AI DC load dynamics can be achieved by employing waveform measurement units capable of capturing load profiles at very short time intervals (below 10 ms) [7]. Such fine-grained measurements enable the observation of fast transient behaviors and rapid load fluctuations that are otherwise invisible to conventional metering infrastructure. Building on this high-fidelity data, machine learning techniques can be leveraged to accurately predict short-term and long-term load patterns, supporting proactive control and optimization strategies. Furthermore, close collaboration between utilities and DC operators is essential to enable secure access to real-time load profiles, ensuring data confidentiality while facilitating coordinated grid-DC operation and enhanced system reliability.

5.2.3 Remaining Useful Life (RUL) Prediction for ESSs. Lithium-ion batteries undergo both calendar degradation and cyclic aging [31]. Calendar degradation is time-dependent and occurs regardless of battery usage, while cyclic degradation results from charge-discharge cycles. Degradation pathways, including solid electrolyte interphase (SEI) growth, lithium plating, and particle fracture [26] lead to capacity fade and increased internal resistance. To avoid the BESSs from entering the non-linear region of the degradation curve, a degradation-aware dispatch strategy needs to be devised. The dispatch strategy must count for hierarchical operation of ESSs inside AI DCs.

5.2.4 FTM grid-scale BESS as A Reliable Reserve. By monitoring and communicating the grid-scale BESS status, located in AI DC sites, with the grid operator, more optimized reserve planning can be achieved [40, 56]. In this case, the total amount of unused energy of grid-scale BESS can be considered a potential reserve that can be injected into the utility grid during contingency events such as the tripping of synchronous generators. The FTM grid-scale BESS can also play the role of an active reserve source in emergency cases by locally feeding the AI DC load. In this case, the grid operator disconnects the AI DC from the grid and allows the AI DC to supply its load using the available BESS co-generation resources. However, this requires precise coordination and control strategies between both the utility and AI DC owners.

5.2.5 AI DC as a Microgrid. AI DC microgrids adopt hybrid power architectures. These architectures integrate BESS, conventional generators, distributed generation units, and utility grid connections (see Fig. 6). The objective is to achieve high reliability and redundancy. The placement of BESS is a key design choice. Proper placement allows BESS to damp fast load fluctuations. It also improves local voltage and frequency stability. BESS enable coordinated operation among different power generation units. They

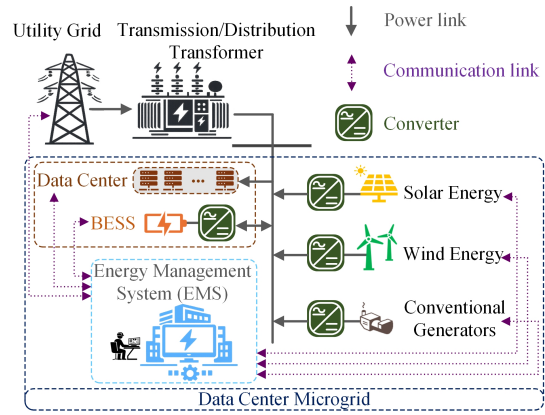


Figure 6: AI DC microgrid.

support grid services such as demand response and peak shaving. They also provide fast reserves. Optimal integration of BESS with other units requires joint power flow and control design. This integration raises the question of BESS versus grid-interactive UPS systems. The boundary between backup and active grid support becomes unclear. In this context, BESS redefine power interface thresholds and reshape power quality metrics. This enables data center microgrids to act as active and flexible grid participants.

6 CONCLUSIONS

This paper provides a brief review of multi-layer ESSs in future AI DCs and their role in supporting both DCs and the utility grid. ESSs ranging from chip-level units to large-scale BESSs are discussed, with particular emphasis on their contribution to grid support and ancillary services. Chip-level ESSs act as the first layer for smoothing heterogeneous GPU power profiles, while large-scale BESSs constitute the final layer in this hierarchical structure. The concept of GiUPS is thoroughly examined, highlighting bidirectional power flow capability as a key enabler of grid-aware UPS design. Furthermore, the main challenges associated with ESS integration in AI DCs are outlined. Finally, as a future direction, the deployment of multi-layer ESS architectures necessitates advanced monitoring systems and EMSs to enable optimal coordination across different layers.

REFERENCES

- [1] 2025. Characteristics and Risks of Emerging Large Loads. <https://www.nerc.com/globalassets/who-we-are/standing-committees/rstc/whitepaper-characteristics-and-risks-of-emerging-large-loads.pdf> White Paper.
- [2] 2025. How New GB300 NVL72 Features Provide Steady Power for AI. <https://developer.nvidia.com/blog/how-new-gb300-nvl72-features-provide-steady-power-for-ai/>
- [3] Rouzbeh Reza Ahrabi, Alireza Mousavi, Ebrahim Mohammadi, Ryan Wu, and Aoxia Kevin Chen. 2025. AI-Driven Data Center Energy Profile, Power Quality, Sustainable Siting, and Energy Management: A Comprehensive Survey. In *2025 IEEE Conference on Technologies for Sustainability (SusTech)*. IEEE, 1–8.
- [4] Dizar Al Kez and Aoife Foley. 2025. Instability Risks from Programmable AI Load Ramping in Low-Inertia Grids. (2025).
- [5] Analog Devices. 2024. ADI OCP ORV3 BBU Reference Design. <https://wiki.analog.com/resources/eval/adi-ocp-orv3-bbu-reference-design>.
- [6] Bjorlin, Alexis. 2022. OCP Summit 2022: Open hardware for AI infrastructure (Grand Teton / ORV3 rack and power). <https://engineering.fb.com/2022/10/18/open-source/ocp-summit-2022-grand-teton/>.
- [7] Babu Chalamala et al. 2025. Data center growth and grid readiness (tr131). *IEEE Power and Energy Society* (2025).

- [8] Xin Chen, Xiaoyang Wang, Ana Colacelli, Matt Lee, and Le Xie. 2025. Electricity demand and grid impacts of AI data centers: Challenges and prospects. *arXiv preprint arXiv:2509.07218* (2025).
- [9] Esha Choukse, Brijesh Warriar, Scot Heath, Luz Belmont, April Zhao, Hassan Ali Khan, Brian Harry, Matthew Kappel, Russell J Hewett, Kushal Datta, et al. 2025. Power stabilization for AI training datacenters. *arXiv preprint arXiv:2508.14318* (2025).
- [10] Jae-Won Chung, Yile Gu, Insu Jang, Luoxi Meng, Nikhil Bansal, and Mosharaf Chowdhury. 2024. Reducing energy bloat in large model training. In *Proceedings of the ACM SIGOPS 30th Symposium on Operating Systems Principles*. 144–159.
- [11] Jae-Won Chung, Jeff J Ma, Ruofan Wu, Jiachen Liu, Oh Jun Kweon, Yuxuan Xia, Zhiyu Wu, and Mosharaf Chowdhury. 2025. The ML ENERGY benchmark: Toward automated inference energy measurement and optimization. *arXiv preprint arXiv:2505.06371* (2025).
- [12] Jae-Won Chung, Ruofan Wu, Jeff J Ma, and Mosharaf Chowdhury. 2026. Where Do the Joules Go? Diagnosing Inference Energy Consumption. *arXiv preprint arXiv:2601.22076* (2026).
- [13] Philip Colangelo, Ayse K Coskun, Jack Megrue, Ciaran Roberts, Shayan Sengupta, Varun Sivaram, Ethan Tiao, Aroon Vijaykar, Chris Williams, Daniel C Wilson, et al. 2025. AI data centres as grid-interactive assets. *Nature Energy* (2025), 1–8.
- [14] Nils Collath, Benedikt Tepe, Stefan Englberger, Andreas Jossen, and Holger Hesse. 2022. Aging aware operation of lithium-ion battery energy storage systems: A review. *Journal of Energy Storage* 55 (2022), 105634. <https://doi.org/10.1016/j.est.2022.105634>
- [15] Alasdair J. Crawford, Qian Huang, Michael C.W. Kintner-Meyer, Ji-Guang Zhang, David M. Reed, Vincent L. Sprenkle, Vilayanur V. Viswanathan, and Daiwon Choi. 2018. Lifecycle comparison of selected Li-ion battery chemistries under grid and electric vehicle duty cycle combinations. *Journal of Power Sources* 380 (2018), 185–193. <https://doi.org/10.1016/j.jpowsour.2018.01.080>
- [16] Arturo Di Filippi and Luca Valentini. 2025. *How to Maximize Revenues from Your Data Center Energy Storage System with Grid Interactive UPS*. Technical Report. Vertiv. Vertiv White Paper.
- [17] Arturo Di Filippi and Luca Valentini. 2025. How to Maximize Revenues from Your Data Center Energy Storage System with Grid Interactive UPS. https://www.vertiv.com/4918e5/globalassets/documents/white-papers/white-paper-maximize-revenues-data-center-energy-storage-grid-ups_329440_2.pdf
- [18] Rouslan Dimitrov, Harry Petty, Neeraj Srivastava, and Mathias Blake. 2025. How New GB300 NVL72 Features Provide Steady Power for AI. <https://developer.nvidia.com/blog/how-new-gb300-nvl72-features-provide-steady-power-for-ai/>.
- [19] Patrick Donovan. 2025. *Understanding BESS: Battery Energy Storage Systems for Data Centers*. Technical Report White Paper 185. Schneider Electric Energy Management Research Center.
- [20] Eaton. [n.d.]. Eaton and Microsoft’s EnergyAware UPS technology pilot project. <https://www.eaton.com/us/en-us/products/backup-power-ups-surge-it-power-distribution/backup-power-ups/dual-purpose-ups-technology.html>.
- [21] Eaton. [n.d.]. Supercapacitors in AI Data Centers. <https://www.eaton.com/us/en-us/products/electronic-components/infographics/supercaps-in-ai-datacenters.html>.
- [22] AL-Hazemi Fawaz, Josip Lorincz, and Alaeldin FY Mohammed. 2019. Minimizing data center uninterruptible power supply overload by server power capping. *IEEE Communications Letters* 23, 8 (2019), 1342–1346.
- [23] Dina Genkina. 2025. Will Supercapacitors Come to AI’s Rescue? Power bursts in large AI workloads can threaten to overwhelm the grid. <https://spectrum.ieee.org/supercapacitor-2671883490>.
- [24] Elinor Ginzburg-Ganz, Pavel Lifshits, Ram Machlev, Juri Belikov, Ziv Krieger, and Yoash Levron. 2025. Technical Challenges of AI Data Center Integration into Power Grids—A Survey. *Energies* 19, 1 (2025), 137.
- [25] D Gu, X Xie, G Huang, X Jin, and X Liu. [n.d.]. Energy-Efficient GPU Clusters Scheduling for Deep Learning. arXiv 2023. *arXiv preprint arXiv:2304.06381* ([n.d.]).
- [26] A. Hassan, S. A. Khan, R. Li, W. Su, X. Zhou, M. Wang, and B. Wang. 2023. Second-Life Batteries: A Review on Power Grid Applications, Degradation Mechanisms, and Power Electronics Interface Architectures. *Batteries* 9, 12 (2023), 571.
- [27] Jared Huntington and Mike Tu. 2025. 800 VDC Architecture for Next-Generation AI Infrastructure. <https://nvdam.nvidia.com/assets/share/asset/zlg5snufeo>.
- [28] IEC. 2025. Medium voltage DC (MVDC) grids for all-electric society. <https://www.iec.ch/basecamp/medium-voltage-dc-mvdc-grids-all-electric-society>.
- [29] A Karpati, Gy Zsigmond, M Vörös, and Marianna Lendvay. 2012. Uninterruptible Power Supplies (UPS) for data center. In *2012 IEEE 10th Jubilee International Symposium on Intelligent Systems and Informatics*. IEEE, 351–355.
- [30] Jagdeep Kaur and Sarbjeet Kaur Bath. 2025. Harmonic distortion in power systems due to electronic control and renewable energy integration: a comprehensive review. *Discover Electronics* 2, 1 (2025), 67. <https://doi.org/10.1007/s44291-025-00111-9>
- [31] G. P. Kostenko. 2024. Accounting Calendar and Cyclic Ageing Factors in Diagnostic and Prognostic Models of Second-Life EV Batteries.
- [32] Yuzhuo Li and Yunwei Li. 2025. AI Load Dynamics—A Power Electronics Perspective. *arXiv preprint arXiv:2502.01647* (2025).
- [33] Zhirui Liang, Jae-Won Chung, Mosharaf Chowdhury, Jiashi Chen, and Vladimir Dvorkin. 2026. GPU-to-Grid: Voltage Regulation via GPU Utilization Control. *arXiv preprint arXiv:2602.05116* (2026).
- [34] Issah Babatunde Majeed and Nnamdi I Nwulu. 2022. Impact of reverse power flow on distributed transformers in a solar-photovoltaic-integrated low-voltage network. *Energies* 15, 23 (2022), 9238.
- [35] Sulav Malla, Qingyuan Deng, Zoh Ebrahimzadeh, Joe Gasperetti, Sajal Jain, Parimala Kondety, Thiara Ortiz, and Debra Vieira. 2020. Coordinated Priority-aware Charging of Distributed Batteries in Oversubscribed Data Centers. In *2020 53rd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*. 839–851. <https://doi.org/10.1109/MICRO50266.2020.00073>
- [36] North American Electric Reliability Corporation (NERC). 2025. *Characteristics and Risks of Emerging Large Loads*. Technical Report. North American Electric Reliability Corporation. Large Loads Task Force White Paper.
- [37] NVIDIA. 2025. NVML API Reference Guide - GPU Deployment and Management Documentation. https://docs.nvidia.com/deploy/nvml-api/group__nvmlDeviceQueries.html#group__nvmlDeviceQueries_1gf754f109beca3a4a8c8c1cd650d7d66c.
- [38] Janne Paananen. 2023. Grid-Interactive Data Centers Enabling Energy Transition: Data center’s hidden potential to provide essential grid services of a future power system. *IEEE Electrification Magazine* 11, 3 (2023), 26–34. <https://doi.org/10.1109/MELE.2023.3291195>
- [39] Janne Paananen and Ehsan Nasr. 2021. Grid-interactive data centers: enabling decarbonization and system stability. *Dublin, Ireland* (2021).
- [40] Nitin Padmanabhan, Mohamed Ahmed, and Kankar Bhattacharya. 2020. Battery Energy Storage Systems in Energy and Reserve Markets. *IEEE Transactions on Power Systems* 35, 1 (2020), 215–226. <https://doi.org/10.1109/TPWRS.2019.2936131>
- [41] Lyuming Pan, Manrong Song, Nimra Muzaffar, Liuping Chen, Chao Ji, Shengxin Yao, Junhui Xu, Weixiong Wu, Yubai Li, Jie Chen, et al. 2025. Salt cavern redox flow battery: The next-generation long-duration, large-scale energy storage system. *Current Opinion in Electrochemistry* 49 (2025), 101604.
- [42] Quanta Tech LLC. 2025. Understanding AI Load Profiles and Their Impact on Power Systems. Online Seminar (Webinar).
- [43] Saifur Rahman and Tafsir Ahmed Khan. 2026. Energy Storage Systems for AI Data Centers: A Review of Technologies, Characteristics, and Applicability. *Energies* 19, 3 (2026), 634.
- [44] Roach, John. 2022. Microsoft datacenter batteries to support growth of renewables on the power grid. <https://news.microsoft.com/source/features/sustainability/ireland-wind-farm-datacenter-ups/>.
- [45] Ashkan Safari, Frede Blaabjerg, and Arman Oshnoei. 2026. A research-industry perspective of battery systems technology for next-generation data centers. *Journal of Energy Storage* 152 (2026), 120386. <https://doi.org/10.1016/j.est.2026.120386>
- [46] Ashkan Safari, Hoda Sorouri, Afshin Rahimi, and Arman Oshnoei. 2025. A Systematic Review of Energy Efficiency Metrics for Optimizing Cloud Data Center Operations and Management. *Electronics* 14, 11 (2025), 2214.
- [47] David Sun, Dmitriy Shapiro, Ben Kim, Jayati Athavale, and Rommel Mercado. 2022. Open Compute Project: Open Rack V3 BBU Shelf (Rev 1.1). <https://www.opencompute.org/documents/open-rack-v3-bbu-shelf-spec-rev1-1-pdf-1>
- [48] David Sun, Dmitriy Shapiro, Ben Kim, Jayati Athavale, and Rommel Mercado. 2023. Open Compute Project: Open Rack V3 48V BBU (Rev 1.4). <https://www.opencompute.org/documents/open-rack-v3-bbu-module-spec-1-4-pdf>
- [49] Mehmet Türker Takci, Meysam Qadrdan, Jon Summers, and Jonas Gustafsson. 2025. Data centres as a source of flexibility for power systems. *Energy Reports* 13 (2025), 3661–3671.
- [50] Xiaojie Tao and Rajit Gadh. 2025. Coordinated Fast Frequency Response from Electric Vehicles, Data Centers, and Battery Energy Storage Systems. *arXiv preprint arXiv:2512.14136* (2025).
- [51] UL Solutions. [n.d.]. UL 9540A Test Method for Battery Energy Storage Systems (BESS). <https://www.ul.com/services/ul-9540a-test-method>.
- [52] Yi Wang, Qinglai Guo, and Min Chen. 2025. Providing load flexibility by reshaping power profiles of large language model workloads. *Advances in Applied Energy* (2025), 100232.
- [53] Zimu Wang, Zhiqiang Yin, Jinyu Yang, and Jiangjiang Wang. 2025. Coordinated optimization of distributed energy system and storage-enhanced uninterruptible power supply in data center: A three-level optimization framework with model predictive control. *Energy Conversion and Management* 342 (2025), 120137.
- [54] Keith Watson. 2025. Data Centers – A Good Grid Citizen. Presentation. Slide on grid support and batteries; Mission Critical Solutions presentation.
- [55] Ruofan Wu, Jae-Won Chung, and Mosharaf Chowdhury. 2026. Kareus: Joint Reduction of Dynamic and Static Energy in Large Model Training. *arXiv preprint arXiv:2601.17654* (2026).
- [56] Bolun Xu, Yishen Wang, Yury Dvorkin, Ricardo Fernández-Blanco, Cesar A. Silva-Monroy, Jean-Paul Watson, and Daniel S. Kirschen. 2017. Scalable Planning for Energy Storage in Energy and Reserve Markets. *IEEE Transactions on Power Systems* 32, 6 (2017), 4515–4527. <https://doi.org/10.1109/TPWRS.2017.2682790>