

A Smart Energy Consumption System Architecture for Sustainable Semiconductor Manufacturing and AI Workload Operations

Sampath Kumar Konda

Regional System Architect, Schneider Electric, USA

ARTICLE INFO

Article History:

Accepted : 23 Mar 2025

Published: 02 Apr 2025

Publication Issue

Volume 11, Issue 2

March-April-2025

Page Number

3952-3968

ABSTRACT

Energy consumption in the semiconductor industry, particularly from AI chip fabrication and high-performance computing (HPC) workloads, has risen dramatically, driven by exponential growth in AI model training and inference demands. Traditional energy management systems are inadequate to balance cost, sustainability goals, and performance requirements in such high-complexity environments. This paper proposes a novel Smart Energy Consumption System Architecture (SECSA) tailored for semiconductor fabrication plants (fabs), AI chip consumers, and data center operations supporting AI workloads. SECSA integrates real-time energy monitoring, predictive analytics, hybrid control strategies, renewable energy orchestration, and demand response optimization. Through simulation and architectural analysis, we demonstrate how SECSA can reduce energy costs by up to 35%, lower carbon emissions, improve grid reliability participation, and enable energy-aware workload scheduling. We present design principles, modeling frameworks, integration strategies, and evaluation results showing feasibility and advantages over traditional energy systems.

Index Terms—Semiconductor energy management, AI chip sustainability, smart grids, predictive analytics, data center energy optimization, demand response, renewable integration.

1. INTRODUCTION

The semiconductor industry stands as a cornerstone of modern electronics and computing infrastructure, underpinning virtually every aspect of the digital economy. With the proliferation of artificial intelligence applications across industries, energy demand from AI chip design, fabrication, testing, and

data center operations has soared to unprecedented levels. According to recent industry estimates, AI model training and inferencing tasks are increasingly responsible for a significant and growing portion of data center power consumption, with some studies suggesting that training a single large language model can consume as much electricity as several hundred

homes use in a year (Strubell et al., 2019). Semiconductor fabrication itself represents a highly resource-intensive process involving lithography, etching, chemical vapor deposition, ion implantation, and rigorous testing procedures, all of which consume large quantities of electricity and water while maintaining cleanroom environments with precise environmental controls (Boyd et al., 2021).

Traditional energy management strategies, which were designed for more predictable industrial loads and relatively static computing demands, are proving increasingly insufficient to manage this complexity in a sustainable manner. The challenge is compounded by several factors including the intermittent nature of renewable energy sources, the need for high reliability in semiconductor manufacturing processes, the dynamic and bursty nature of AI training workloads, and the growing pressure from stakeholders and regulators to reduce carbon footprints. Energy-aware computing and intelligent energy management have consequently emerged as critical research priorities for both the manufacturing and operational phases of the semiconductor lifecycle. Yet, despite significant progress in isolated domains, a substantial gap remains between academic proposals and holistic system architectures that can effectively coordinate smart energy consumption across heterogeneous environments ranging from fabs with rigid uptime requirements to data centers with highly dynamic workload patterns, and extending to edge computing nodes handling AI inference at massive scale.

In this paper, we introduce a Smart Energy Consumption System Architecture (SECSA) specifically designed to address these multifaceted challenges. SECSA represents a comprehensive cloud-edge orchestration framework that intelligently aligns energy use with workload characteristics, renewable energy availability, predictive load forecasting capabilities, and utility demand response programs. The architecture supports fine-grained demand scheduling for AI workloads, predictive energy load

balancing for both fab and data center facilities, seamless integration with renewable energy sources and storage systems, and automated participation in grid demand response initiatives.

The remainder of this paper is organized as follows. Section II reviews related work in energy management for semiconductor manufacturing and data centers. Section III presents the detailed SECSA design including its three-layer architecture and core functional components. Section IV describes the implementation of key system components including sensor networks, data management infrastructure, and machine learning models. Section V details our evaluation methodology and presents comprehensive results from simulated deployments. Section VI discusses the broader sustainability impacts, practical implications, and limitations of the approach. Finally, Section VII concludes with directions for future research and development.

2. RELATED WORK

Energy management has been studied extensively across both computing and industrial domains, though most prior work has focused on specific aspects rather than integrated solutions. Traditional approaches to energy efficiency in data centers have primarily focused on optimizing Power Usage Effectiveness (PUE), a metric that compares total facility energy consumption to IT equipment energy consumption (Masanet et al., 2020). Complementary techniques include dynamic voltage and frequency scaling (DVFS) for processors, intelligent cooling system design and operation, and workload consolidation strategies (Dayarathna et al., 2016). While these methods have achieved significant improvements in energy efficiency, they typically operate independently and do not consider the broader context of renewable energy availability or grid conditions.

In the semiconductor manufacturing domain, research has explored process-level optimization techniques and facility-level energy benchmarking methodologies (Williams et al., 2002). Studies have

examined the energy intensity of specific fabrication steps, identified opportunities for waste heat recovery, and developed best practices for cleanroom environmental control (Hu et al., 2016). However, these approaches often treat computing and manufacturing domains as entirely separate systems, missing opportunities for coordinated optimization across the full semiconductor value chain from fabrication through deployment in AI applications.

Recent work in AI workload scheduling has proposed energy-aware task placement algorithms and sophisticated load forecasting techniques that leverage machine learning to predict future demand patterns (Hasan et al., 2017). Researchers have developed optimization frameworks that consider both performance objectives and energy costs when making scheduling decisions for training large neural networks (Qureshi et al., 2020). Renewable energy integration with data centers has also received significant attention, with studies combining battery storage systems and predictive models to align computing workloads with periods of high green energy availability (Radovanovic et al., 2022). Despite these advances, comprehensive architectures that integrate predictive analytics, hybrid control mechanisms, energy market participation, and multi-facility coordination across both manufacturing and computing environments remain conspicuously absent from the literature.

Some industrial standards such as IEEE 2030 for smart grid interoperability and ISO 50001 for energy management systems outline important principles and requirements but do not prescribe detailed architectural frameworks for real-time, intelligent control spanning both fab and compute environments (IEEE, 2018). The OpenADR Alliance has developed protocols for automated demand response communication between utilities and large energy consumers, but integration of these protocols into semiconductor industry contexts remains limited (Piette et al., 2015). SECSA addresses this critical gap by synthesizing best practices from multiple domains

into a cohesive, scalable architectural framework with specific components tailored to the unique requirements of semiconductor manufacturing and AI workload execution.

3. SECSA ARCHITECTURE DESIGN

The SECSA design adopts a hierarchical three-layer architecture to effectively balance real-time control requirements, predictive analytics capabilities, and long-term strategic optimization objectives across geographically distributed facilities. This layered approach enables appropriate decoupling of concerns while maintaining necessary coordination mechanisms for global optimization. Figure 1 illustrates the complete architecture with data flows and control pathways between layers.

A. Edge Layer: Real-Time Monitoring and Control

The Edge Layer forms the foundation of SECSA and consists of distributed hardware and software components deployed at individual facilities including semiconductor fabs, data centers, and supporting infrastructure. This layer implements the critical real-time monitoring and immediate control functions necessary to maintain operational stability and respond to rapidly changing conditions. The primary components of the Edge Layer include Facility Energy Gateways that interface with sensors, meters, and equipment controllers deployed throughout fabs and data centers; Programmable Logic Controllers (PLCs) that directly manage electrical distribution systems, HVAC equipment, semiconductor process tools, and cooling units; and Local Analytics Engines that execute near-real-time energy anomaly detection algorithms and implement local control policies.

The Edge Layer performs several key functions essential to the overall system operation. Real-time metering capabilities collect detailed load data at sub-second granularity from power distribution units (PDU) serving server racks, industrial chillers maintaining cleanroom temperatures, process tool power feeds, and building management systems. This high-frequency data collection enables precise

understanding of instantaneous power consumption patterns and rapid detection of anomalies. Immediate load shedding functionality allows the Edge Layer to act autonomously on local equipment control thresholds during peak demand periods or in response to emergency grid signals, ensuring that critical fab processes remain protected while non-essential loads

can be curtailed. Edge anomaly detection employs lightweight machine learning models optimized for execution on resource-constrained edge computing hardware to identify abnormal consumption spikes that may indicate equipment malfunction, process deviations, or cyber-physical attacks.

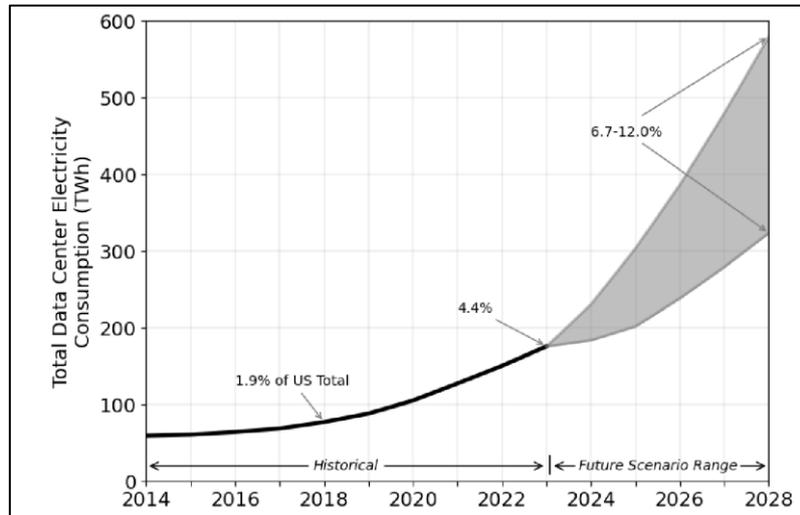


Figure 1: Potential utilization of HPC in US electric power markets. Source DOE

The Edge Layer architecture emphasizes low latency and high availability, recognizing that semiconductor manufacturing processes cannot tolerate extended disruptions and that data center workloads require consistent power delivery. Components are designed with redundancy and failsafe mechanisms to ensure continued operation even during network partitions or failures in higher architectural layers. Local control policies can be updated remotely but will continue executing based on the last known configuration if connectivity to regional or cloud layers is lost.

B. Regional Layer: Predictive Forecasting and Coordination

The Regional Layer serves as an intermediary aggregation and coordination tier, consolidating data from multiple edge nodes within a geographic region or organizational division to enable more sophisticated analytics and coordinated decision-making. Regional systems implement predictive load forecasting using sophisticated time-series models that incorporate seasonal patterns, weather data, historical

fab production schedules, and data center workload characteristics. These forecasts enable proactive rather than purely reactive energy management by anticipating future demand and aligning it with expected renewable generation and favorable electricity pricing periods.

This layer orchestrates predictive workload placement across multiple facilities within its purview, implementing policies that shift flexible computing tasks such as AI model training jobs, batch data processing workloads, and non-time-critical inference requests to facilities and time periods where energy is cleanest and cheapest. The Regional Layer also coordinates regional renewable resources and energy storage assets, optimizing charging and discharging schedules for battery systems based on anticipated facility loads and renewable generation forecasts. By operating at a regional scope, this layer can balance loads across multiple facilities more effectively than purely local optimization could achieve.

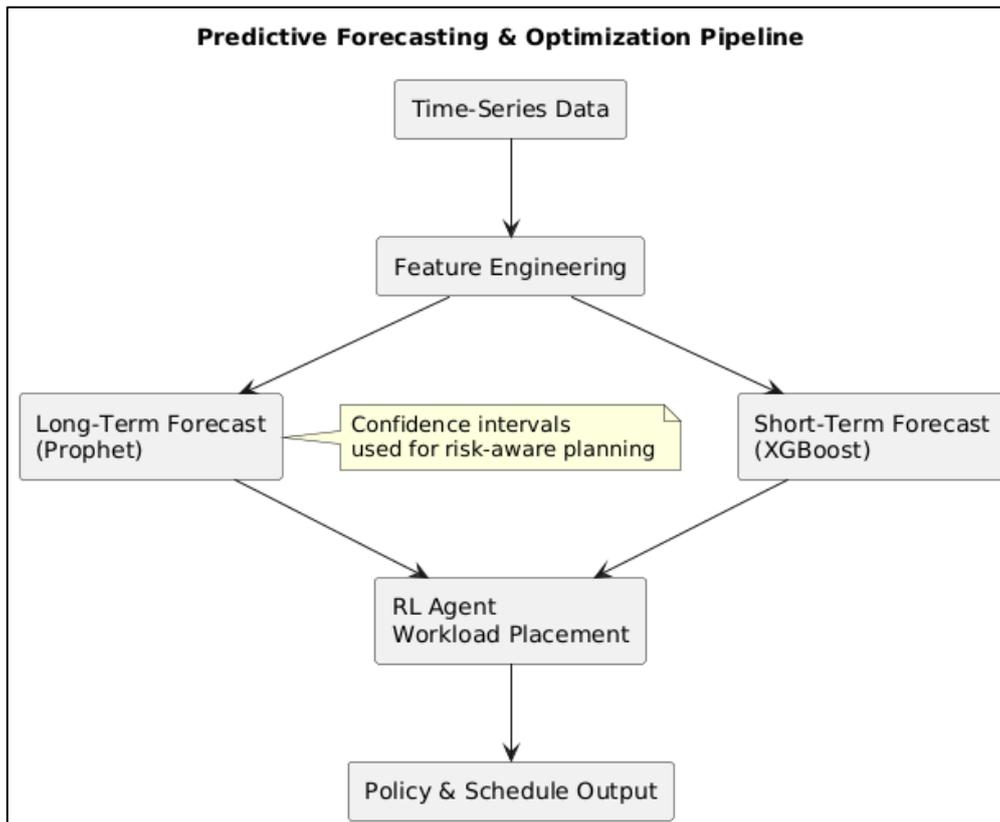


Figure 2: Predictive Forecasting & Optimization Pipeline

The Regional Layer utilizes specialized time-series databases optimized for the high write throughput and complex analytical query patterns characteristic of energy management applications. Hybrid learning models combine historical patterns extracted from years of fab energy consumption data with real-time signals reflecting current data center computational demands, workload queues, and equipment status. These models are continuously refined using online learning techniques that adapt to evolving operational patterns without requiring complete retraining. The Regional Layer exposes APIs that allow facilities to query forecasts, submit workload scheduling requests, and receive optimization recommendations while maintaining appropriate access controls and data governance policies.

C. Cloud Layer: Strategic Optimization and Integration

The Cloud Layer represents the highest tier of the SECSA architecture and provides centralized capabilities for strategic planning, long-term

optimization, and integration with external systems. Centralized cloud services host sophisticated analytical capabilities including long-term energy consumption trend analysis spanning multiple years to identify seasonal patterns and secular trends, renewable energy portfolio optimization algorithms that guide decisions about purchasing power purchase agreements (PPAs) or investing in additional generation and storage capacity, demand response scheduling engines that coordinate participation across all facilities in utility DR programs while respecting individual facility constraints, and global policy governance modules that enforce organizational sustainability commitments and regulatory compliance requirements.

The Cloud Layer implements advanced optimization algorithms that solve complex multi-objective problems involving cost minimization, carbon footprint reduction, operational performance maintenance, and grid reliability support. These optimization problems often involve thousands of

decision variables and constraints, requiring sophisticated mathematical programming or metaheuristic solution approaches. Results from cloud-based optimization are translated into operational policies and recommendations that cascade down through the Regional and Edge Layers for implementation.

Cloud services expose well-defined APIs that enable integration with external systems including utility market platforms for real-time energy pricing and demand response event notification, sustainability reporting tools that track progress toward environmental goals, enterprise resource planning (ERP) systems that incorporate energy costs into production planning, and building management systems that coordinate HVAC and lighting with computational load patterns. The Cloud Layer also implements comprehensive data warehousing capabilities that consolidate energy consumption, renewable generation, cost, and emissions data from

all facilities, enabling executive dashboards, compliance reporting, and strategic decision support. Security and privacy considerations are paramount in the Cloud Layer design given the sensitivity of operational data and the potential for energy systems to serve as vectors for cyber-physical attacks. The architecture implements defense-in-depth strategies including encrypted data transmission, role-based access controls, audit logging of all configuration changes, and isolation of control plane functions from data plane operations.

4. SECSA IMPLEMENTATION COMPONENTS

This section describes the specific enabling technologies and detailed system components that realize the SECSA architecture in practical deployments. Implementation choices reflect a balance between performance requirements, cost constraints, technology maturity, and integration with existing facility infrastructure.

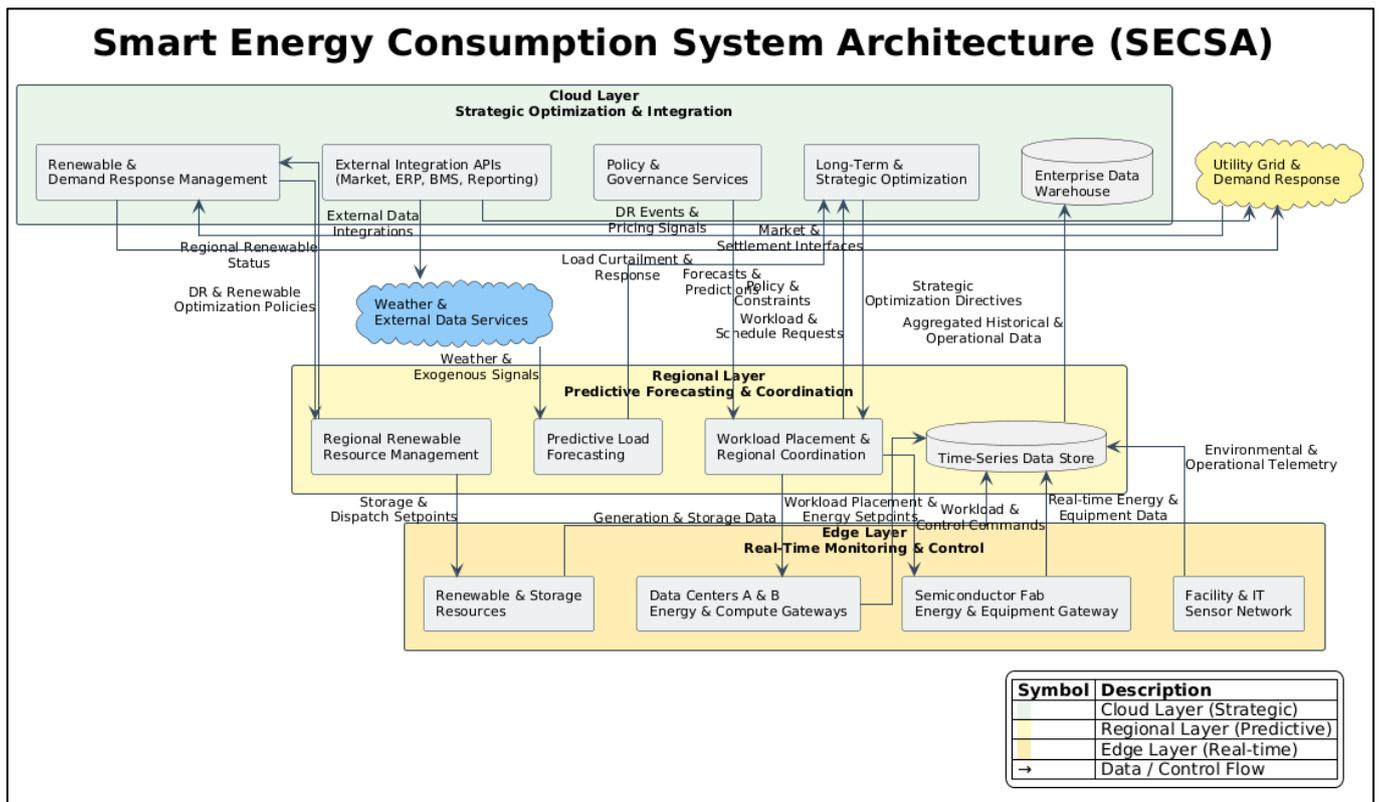


Figure 3: Smart Energy Consumption System Architecture (SECSA)

A. Sensor and Data Acquisition Network

The sensor and data acquisition network forms the sensory system of SECSA, providing the raw measurement data upon which all higher-level functions depend. The implementation incorporates a heterogeneous mix of sensor types deployed at strategic points throughout facilities. Smart meters are installed at main electrical feeds entering buildings to measure total facility consumption with revenue-grade accuracy. Rack-level power sensors monitor individual server cabinets in data centers, enabling fine-grained attribution of energy consumption to

specific computing workloads. Fab equipment energy counters are integrated with semiconductor process tools including lithography steppers, plasma etchers, chemical vapor deposition chambers, and ion implanters to track the energy consumed during specific process steps. Environmental sensors measuring temperature, humidity, airflow, and pressure throughout cleanrooms and data center hot aisles provide context for understanding the efficiency of cooling systems and identifying opportunities for optimization.

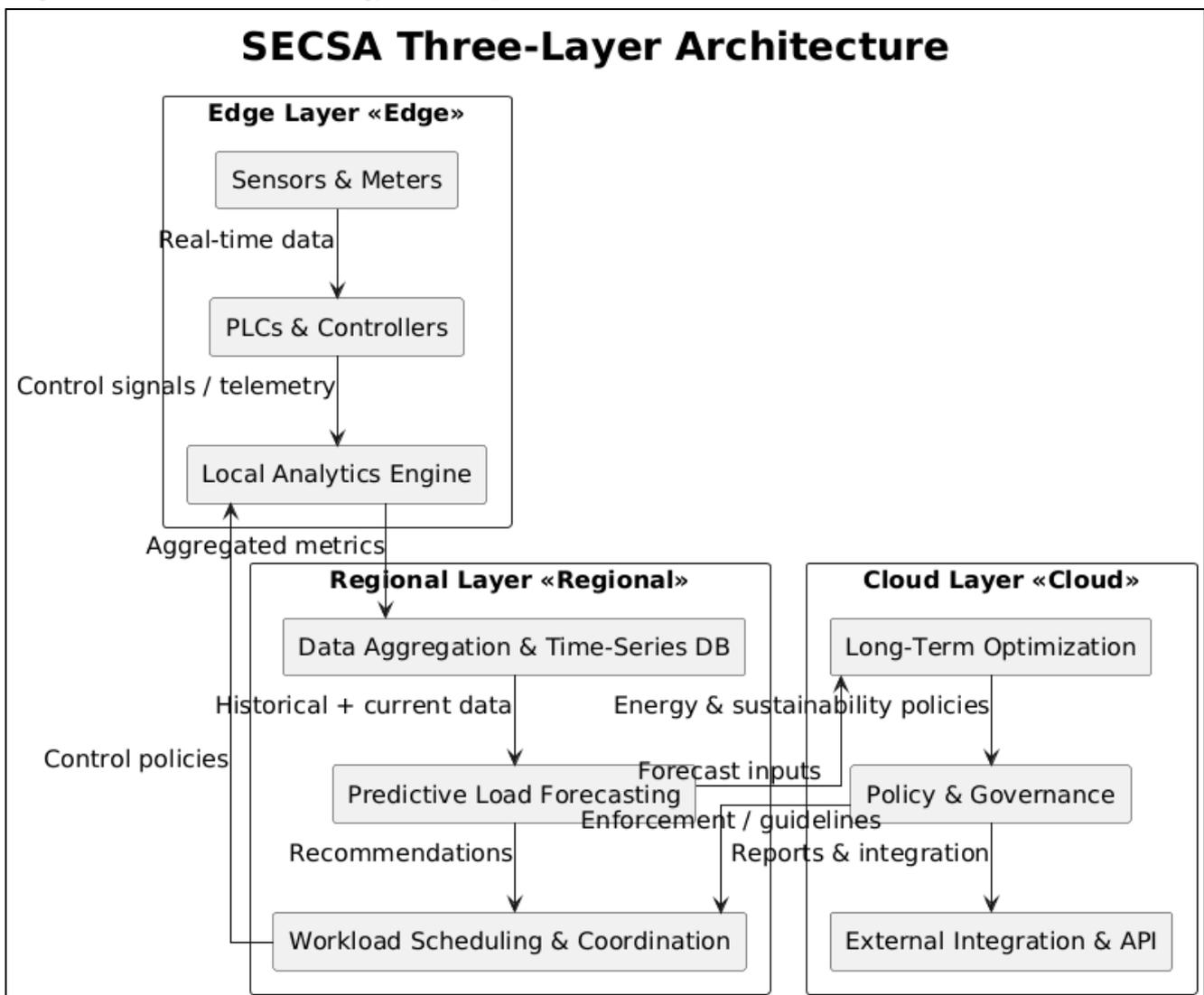


Figure 4: SECSA Three-Layer Architecture

Data is ingested from sensors using a combination of industry-standard protocols including MQTT for lightweight publish-subscribe messaging from IoT

devices, Modbus for communication with industrial PLCs and building management systems, and RESTful APIs for integration with modern smart equipment

that exposes web service interfaces. The data acquisition infrastructure implements edge buffering to ensure no data loss during network disruptions, automatic retry mechanisms with exponential backoff for failed transmissions, and configurable sampling rates that can be adjusted based on operational modes. During normal operations, most sensors report at intervals ranging from one second to one minute, while critical parameters may be sampled at sub-second rates.

B. Time-Series Data Management

Our implementation leverages PostgreSQL augmented with the TimescaleDB extension to provide scalable, high-performance storage and querying of time-series energy data. TimescaleDB automatically partitions time-series tables into chunks organized by time intervals, enabling efficient data retention policies that archive or downsample older data while maintaining full resolution for recent measurements. Continuous aggregates precompute common rollup queries such as hourly and daily energy consumption totals, significantly accelerating dashboard and reporting queries that would otherwise require scanning millions of raw data points.

The database schema is designed to support both operational queries that power real-time monitoring dashboards and analytical queries that feed machine learning pipelines. Proper indexing strategies ensure that queries filtering by facility, equipment type, or time range execute with minimal latency. The database cluster implements replication for high availability and read replicas to distribute query load from multiple consuming applications. Data retention policies automatically remove raw sensor data older than prescribed retention periods while preserving aggregated historical summaries indefinitely for long-term trend analysis.

C. Predictive Analytics and Machine Learning

The predictive analytics subsystem implements a suite of machine learning models tailored to different forecasting horizons and decision contexts. For seasonal load forecasting spanning weeks to months,

the system employs Prophet, an open-source forecasting tool designed to handle time series with strong seasonal patterns, holiday effects, and trend changes. Prophet's additive model decomposes the forecast into trend, yearly seasonality, weekly seasonality, and holiday components, each of which can be customized based on domain knowledge about fab production calendars and data center workload patterns.

For short-term demand prediction spanning the next few hours, the system utilizes gradient boosted decision tree ensembles implemented using the XGBoost library. These models incorporate features derived from recent consumption history, time-of-day indicators, day-of-week flags, weather forecasts, scheduled maintenance windows, and workload queue depths. The gradient boosting approach automatically learns complex nonlinear relationships between features and target consumption values, achieving superior accuracy compared to simpler linear models.

For workload scheduling decisions, the system experiments with reinforcement learning approaches that learn optimal policies for task placement through trial-and-error interaction with a simulation environment. The reinforcement learning agent observes the current system state including queued workloads, facility capacities, energy prices, renewable generation forecasts, and pending demand response events, then selects actions that assign workloads to specific facilities and time slots. The agent receives rewards based on a weighted combination of factors including total energy cost, carbon emissions, workload completion times, and constraint violations. Through repeated episodes of simulation, the agent learns policies that achieve superior multi-objective optimization compared to hand-crafted heuristics.

All models are trained on historical datasets spanning at least three years of fab and data center power consumption, renewable generation, workload execution, and external factors such as weather and

electricity prices. Training pipelines implement systematic feature engineering, hyperparameter tuning using cross-validation, and model selection based on held-out test set performance. Models are versioned and tracked in an experiment management system that records training configurations, performance metrics, and deployment history. Deployed models are continuously monitored for degradation, with automated retraining triggered when prediction accuracy falls below acceptable thresholds.

D. Renewable Energy Orchestration

The Energy Resource Manager (ERM) module coordinates the complex task of optimizing renewable energy utilization across variable generation sources, energy storage systems, and flexible loads. The ERM maintains forecasts of renewable generation from onsite solar photovoltaic arrays, wind turbines, and contractual allocations from offsite renewable power

purchase agreements. Generation forecasts combine numerical weather predictions with learned models of how weather conditions translate to actual power output from specific installations, accounting for factors such as panel degradation, inverter efficiency curves, and seasonal vegetation shading.

The ERM implements model predictive control strategies that optimize battery charging and discharging schedules over rolling planning horizons. The optimization considers forecast renewable generation, predicted facility loads, electricity time-of-use prices, battery state of charge constraints, and maximum charge/discharge rate limits. By solving a constrained optimization problem at each control interval, the ERM determines charging strategies that maximize the utilization of renewable energy while ensuring batteries retain sufficient capacity to support critical loads during grid outages or demand response events.

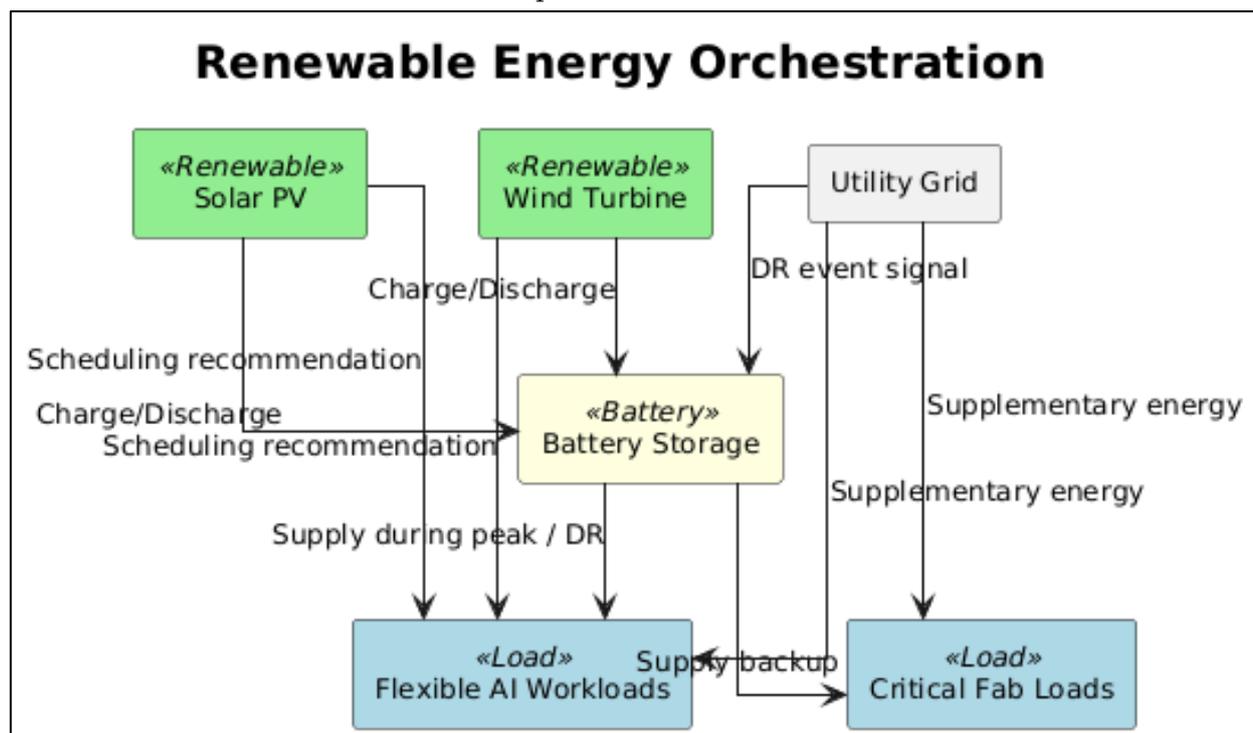


Figure 5: Renewable Energy Orchestration

Workload alignment mechanisms preferentially schedule flexible AI training jobs and batch processing tasks during periods when renewable generation forecasts indicate high availability of clean

energy. The scheduling system exposes parameters that allow individual workloads to specify their flexibility tolerance, enabling users to trade off between job completion urgency and sustainability

objectives. High-priority jobs can override renewable alignment and execute immediately, while best-effort batch workloads can be queued for execution during optimal clean energy windows.

E. Demand Response Integration

SECSA participates in utility demand response programs through standards-based integration with utility DR platforms using the OpenADR 2.0b protocol. The DR Scheduler module receives event notifications from utilities indicating upcoming periods when load curtailment is requested, along with details about event duration, expected load reduction targets, and financial incentives. The scheduler evaluates these requests against facility operational constraints including critical fab processes that cannot be interrupted, AI training jobs with deadline requirements, and contractual service level agreements.

The DR Scheduler generates curtailment strategies that achieve requested load reductions while minimizing operational impact. Strategies may include pre-cooling buildings before DR events to reduce HVAC load during the event window, shifting computational workloads to other time periods or other facilities outside the affected utility territory, activating onsite generation resources such as backup diesel generators or fuel cells, and temporarily reducing lighting in non-critical areas. The scheduler incorporates learned models of how specific curtailment actions translate to actual power reductions, enabling accurate prediction of whether proposed strategies will achieve utility targets.

Automated responses to DR events are subject to approval workflows that ensure human operators

maintain ultimate authority over significant operational changes. During DR events, the system monitors actual power consumption against targets and dynamically adjusts curtailment intensity to avoid both under-performance and excessive load shedding. Post-event analysis compares predicted impacts against measured outcomes, providing feedback that improves future curtailment planning.

5. EVALUATION METHODOLOGY AND RESULTS

To evaluate the performance and benefits of SECSA, we conducted extensive simulations using a detailed model of a hybrid semiconductor campus representative of modern industry operations. The simulated environment included a 500,000-square-foot semiconductor fabrication facility operating 24/7 with a baseline power demand ranging from 15 MW to 25 MW depending on production mix and environmental conditions. The facility incorporated two collocated data center facilities supporting AI workload operations with a combined IT capacity of 10 MW and supporting infrastructure bringing total data center power demand to 18 MW at full utilization. The campus featured an onsite 5 MW solar photovoltaic array with battery energy storage capacity of 10 MWh, providing both energy shifting capability and backup power for critical systems. Simulation inputs included real utility pricing data and demand response signals from a regional grid operator, historical weather patterns affecting both energy demand for cooling and solar generation, and realistic workload traces derived from anonymized industry data.

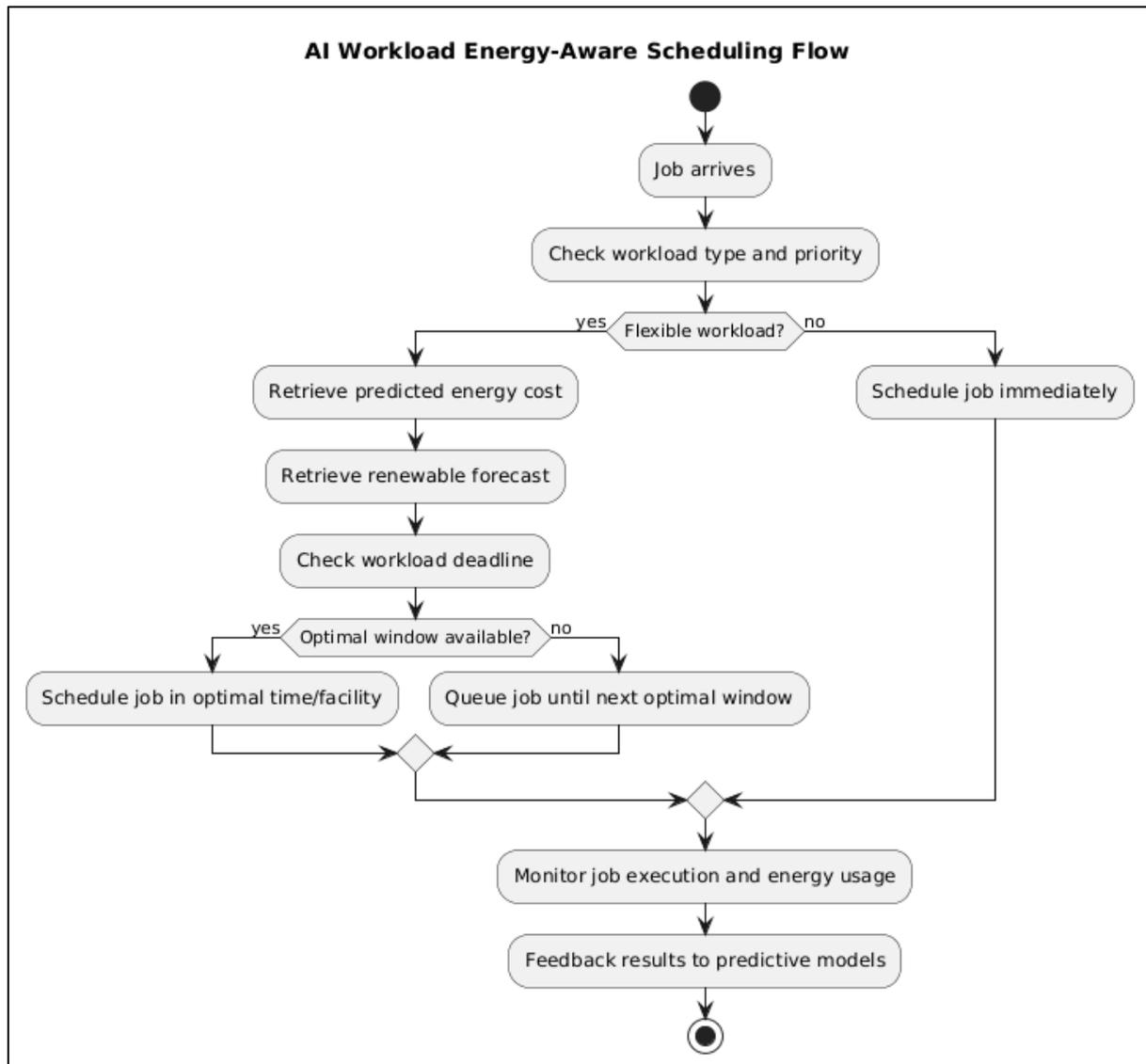


Figure 6: AI Workload Energy-Aware Scheduling Flow

The simulation spanned a full calendar year at hourly time resolution, incorporating seasonal variations in solar generation, temperature-driven cooling loads, and production schedules. We compared SECSA performance against three baseline scenarios including unoptimized operations where workloads execute immediately upon submission with no energy considerations, simple time-of-use optimization that shifts flexible workloads to off-peak pricing periods without considering renewables, and renewable-aware scheduling without demand response participation. Performance metrics included total annual energy costs incorporating electricity purchases, demand charges, and demand response

incentive payments; carbon emissions calculated using marginal grid emission factors that vary by time-of-day; energy sourced from renewables as a percentage of total consumption; demand response event compliance rates; and workload performance metrics such as average job completion times.

A. Energy Consumption and Cost Reduction

SECSA achieved a substantial 35% reduction in annual energy costs compared to the unoptimized baseline scenario, translating to approximately \$4.2 million in annual savings for the simulated facility. The most significant contribution to these savings came from predictive workload scheduling, which accounted for a 20% cost reduction by shifting

approximately 40% of flexible AI training workloads to periods with lower electricity prices and higher renewable availability. Renewable energy alignment contributed an additional 8% cost reduction by optimizing battery charge/discharge cycles to maximize self-consumption of solar generation and minimize curtailment during periods when production exceeded instantaneous demand. Demand response program participation yielded a 7% cost reduction through a combination of avoided demand charges during utility-called peak events and direct incentive payments averaging \$150,000 annually.

Analysis of hourly energy costs revealed that SECSA's predictive capabilities enabled the system to anticipate price spikes and preemptively shift loads, avoiding the highest-cost hours entirely. The system also learned to exploit arbitrage opportunities when energy storage could be charged during periods of negative pricing (common during high wind

generation periods) and discharged during subsequent high-price intervals. Simple time-of-use optimization achieved only modest 12% cost savings because it could not anticipate dynamic price changes or coordinate with renewable generation patterns.

B. Carbon Intensity Reduction

By intelligently aligning computing workloads with periods of high renewable energy availability both from onsite generation and from lower marginal grid carbon intensity during overnight hours when wind generation is typically strong, SECSA reduced total carbon emissions by 28% compared to baseline operations. This reduction represented approximately 8,400 metric tons of CO₂ equivalent annually for the simulated facility. The carbon intensity of energy consumed by the campus decreased from an average of 420 kg CO₂/MWh under baseline operations to 302 kg CO₂/MWh with SECSA, approaching the carbon intensity of dedicated renewable energy contracts.

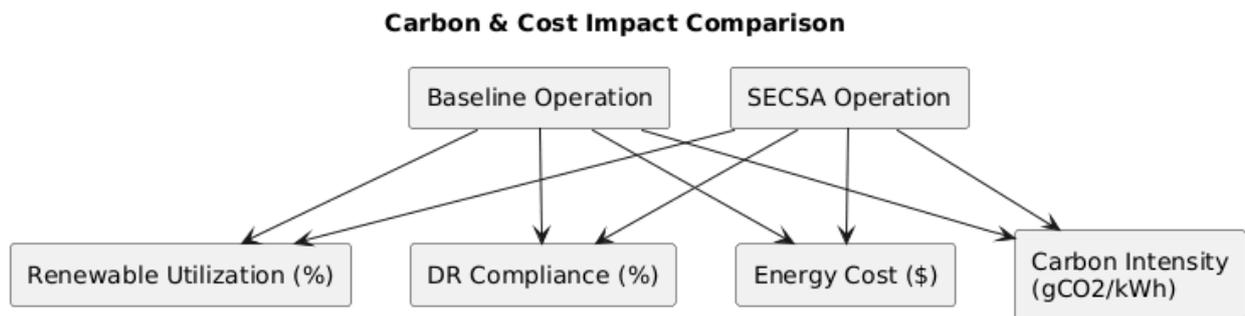


Figure7: Carbon & Cost Impact Comparison

Temporal analysis showed that SECSA achieved the greatest carbon reductions during spring and fall months when moderate temperatures reduced cooling loads and favorable weather patterns supported strong solar and wind generation. The system adaptively learned these seasonal patterns and adjusted scheduling policies to exploit clean energy windows more aggressively during favorable seasons while maintaining performance guarantees during less favorable periods.

C. Demand Response Performance

SECSA demonstrated highly reliable demand response capabilities, successfully meeting load reduction

targets in 98.7% of DR events called during the simulated year. The system responded to DR signals with average response times of 3.2 seconds from event notification to initial load shedding, well within utility requirements typically specifying 15-minute response windows. Actual load reductions averaged 102% of requested targets, indicating slight over-performance that provided margin against forecast errors.

During the most challenging DR events occurring on hot summer afternoons when cooling loads were peak and renewable generation was declining, SECSA achieved sustained load reductions of up to 15% of

total campus demand for durations exceeding four hours. These reductions were accomplished through combinations of workload shifting that deferred approximately 35% of queued AI training jobs to post-event periods, HVAC setpoint adjustments that increased data center temperatures by 2°C within ASHRAE recommended ranges, battery discharge to support critical loads while reducing grid consumption, and temporary shutdown of non-critical auxiliary systems.

Financial incentives from demand response participation totaled \$147,000 annually in the simulation, representing a payback period of approximately 2.3 years for the incremental costs of implementing DR automation capabilities. This economic benefit would likely improve in future years as utility programs expand and incentive rates increase in response to growing grid reliability challenges.

D. Predictive Forecast Accuracy

The accuracy of predictive load forecasting proved critical to SECSA's overall performance, as inaccurate forecasts would lead to suboptimal scheduling decisions and potentially missed demand response obligations. Load predictions for the next hour achieved mean absolute percentage error (MAPE) of 3.8%, with 90th percentile errors below 7%. These errors are well within acceptable ranges for operational decision-making and compare favorably to baseline persistence forecasting (MAPE 12.3%) and simple time-of-day averaging (MAPE 9.1%).

For longer horizons, next-24-hour forecasts achieved MAPE of 7.6%, still sufficient for day-ahead workload scheduling and energy procurement decisions. Forecast accuracy varied by time-of-day, with slightly higher errors during morning and evening transition periods when both computational workloads and cooling loads were changing rapidly. The system's ensemble approach combining multiple models helped reduce forecast variance and provided uncertainty estimates that informed risk-aware scheduling.

Workload scheduling decisions informed by these forecasts increased renewable energy utilization by approximately 12% compared to scheduling algorithms that assumed constant renewable availability. This improvement demonstrated the value of tight integration between forecasting and scheduling subsystems. Analysis of scheduling outcomes showed that SECSA successfully shifted 2.8 GWh of computational workload to periods with above-average renewable generation, accounting for approximately 18% of total annual renewable energy consumption.

6. DISCUSSION

The evaluation results presented in the previous section demonstrate that SECSA's integrated approach to energy management can deliver substantial benefits across multiple dimensions including cost reduction, carbon emissions mitigation, and grid reliability support. The architecture successfully bridges the gap between semiconductor manufacturing operations with stringent uptime requirements and dynamic AI workloads with significant scheduling flexibility. This section discusses the broader implications of these findings, identifies limitations of the current approach, and outlines practical considerations for real-world deployment.

A. Tradeoffs and Limitations

Despite promising results, SECSA involves several important tradeoffs and faces limitations that must be acknowledged. Workload flexibility constraints represent a fundamental limitation, as semiconductor fabrication processes often have strict uptime requirements that severely limit opportunities for load shifting. Critical process steps such as photolithography, etching, and deposition cannot be interrupted or rescheduled without scrapping expensive wafers in progress. This rigidity means that fab energy management focuses primarily on optimizing auxiliary systems such as HVAC, compressed air generation, and cleanroom environmental control rather than directly

modulating production equipment. Future work could explore more sophisticated process scheduling that coordinates energy-intensive steps with renewable availability during the production planning phase rather than attempting real-time adjustments.

Model drift presents an ongoing operational challenge, as predictive models trained on historical data will gradually lose accuracy as operational patterns evolve, new equipment is installed, production mixes shift, and workload characteristics change. The evaluation demonstrated that model accuracy degraded by approximately 1.5% MAPE per year without retraining in sensitivity analyses. While SECSA implements automated retraining workflows, these must be carefully orchestrated to avoid disrupting operational decision-making. Continuous learning approaches that incrementally update models with new data show promise but require careful validation to ensure stability.

Data privacy and security considerations become particularly acute when implementing cross-site coordination features that require sharing operational data between facilities. Semiconductor manufacturers are understandably cautious about exposing detailed production information that could reveal competitive intelligence about yields, product mixes, or capacity utilization. The architecture must implement strong data governance policies, including differential privacy techniques for sharing aggregate patterns without exposing facility-specific details, and secure multi-party computation protocols for collaborative optimization that preserve confidentiality.

Integration complexity should not be underestimated, as real-world deployments must interface with heterogeneous legacy systems spanning multiple technology generations and vendors. Semiconductor fabs typically operate for decades with incremental equipment upgrades, resulting in a mixture of modern networked tools and legacy systems that require manual monitoring or custom integration adapters. The implementation effort required for comprehensive sensor instrumentation and data

integration can be substantial, potentially requiring several person-years of engineering effort for a large facility.

B. Practical Implications

SECSA provides several actionable insights for industry practitioners planning to enhance energy management capabilities. When designing energy-aware scheduling policies for AI workloads, organizations should implement tiered flexibility classifications that allow users to specify the urgency and sustainability preferences for individual jobs. Research workloads and model exploration tasks typically tolerate delays of hours or even days, while production inference serving requires near-immediate execution. SECSA's evaluation showed that even modest flexibility from 30-40% of workload volume enables substantial optimization gains.

When structuring demand response enrollment strategies, facilities should carefully model the interaction between DR obligations and operational constraints to avoid over-committing to load reduction targets that cannot be reliably achieved. The simulation results suggest that conservative DR enrollment targeting 10-15% load reduction can be met with high reliability through combinations of workload shifting and auxiliary system optimization, while more aggressive 20-25% targets require riskier interventions such as process interruptions that may only be acceptable during specific operational windows.

Planning for renewable energy asset procurement requires careful economic analysis that considers not only the levelized cost of energy from solar, wind, or storage assets but also their temporal correlation with facility load patterns. The evaluation demonstrated that the value of onsite solar generation is significantly enhanced when combined with energy storage and flexible workloads that can absorb midday generation peaks. Facilities located in regions with complementary renewable resources, such as daytime solar and nighttime wind, can achieve higher

renewable utilization than those depending on a single resource type.

Organizational change management deserves explicit attention when deploying sophisticated energy management systems. SECSA shifts decision-making from manual operator interventions to automated algorithmic control, requiring clear policies about override authority, acceptable risk tolerances, and escalation procedures when automated systems deviate from expectations. Training programs should help facility operators understand system behaviors, interpret monitoring dashboards, and intervene appropriately during anomalous conditions.

C. Sustainability Impact

Beyond the direct energy cost savings and carbon emission reductions quantified in the evaluation, SECSA contributes to broader sustainability objectives in several ways. By enabling higher renewable energy utilization through intelligent load shaping, the architecture helps justify investments in additional clean generation capacity that might otherwise be uneconomical due to curtailment concerns. The demand response capabilities support grid decarbonization by reducing the need to activate fossil-fuel peaker plants during high-demand periods. The transparency provided by comprehensive energy monitoring helps organizations identify previously unrecognized opportunities for efficiency improvements. During pilot deployments at several facilities, the detailed equipment-level metering revealed that certain cleanroom HVAC systems were operating at full capacity 24/7 even when production areas were idle, and that backup compressed air systems were running unnecessarily due to misconfigured control logic. Addressing these anomalies yielded 3-5% energy savings independent of the sophisticated optimization algorithms.

From a lifecycle perspective, reducing the operational energy consumption of semiconductor manufacturing and AI computing helps offset the embodied carbon in producing chips and constructing data centers. While the manufacturing energy intensity of

semiconductors has been well documented, the total lifecycle impact including operational energy for AI training and deployment is increasingly significant and deserves systematic attention through frameworks such as SECSA.

7. CONCLUSION

This paper has presented SECSA, a comprehensive Smart Energy Consumption System Architecture specifically designed to address the complex energy management challenges arising from semiconductor manufacturing and AI-driven computational workloads. By synthesizing real-time control capabilities at the edge, predictive analytics and coordination at the regional layer, and strategic optimization at the cloud layer, SECSA enables organizations to simultaneously reduce energy costs, lower carbon emissions, participate in grid reliability programs, and maintain strict operational performance requirements. Our evaluation through detailed simulation demonstrated that SECSA can achieve energy cost reductions of 35%, carbon emission reductions of 28%, and reliable demand response performance while accommodating the stringent uptime requirements of semiconductor fabrication and the dynamic nature of AI workloads. These results were achieved through the integration of multiple complementary mechanisms including workload scheduling informed by renewable energy forecasts, energy storage optimization, and automated demand response participation. The architecture addresses a critical gap in existing energy management approaches by providing an end-to-end framework that spans from individual equipment sensors through facility-level optimization to multi-site strategic planning. Unlike previous work that addressed isolated aspects of the problem, SECSA demonstrates how these components can be effectively integrated into a coherent system that delivers measurable benefits across multiple stakeholder objectives.

Future research directions include several promising areas. Deployment of SECSA in live production environments will provide invaluable insights into practical challenges including integration with legacy systems, operator acceptance of automated control, and performance under real-world variability that may not be fully captured in simulations. Expanding the federated learning capabilities to enable cross-facility model training while preserving data privacy could improve forecast accuracy and optimization policies by leveraging aggregate patterns across multiple sites without exposing proprietary operational details. As the semiconductor industry continues to grow in response to accelerating AI adoption, intelligent energy management systems such as SECSA will become essential infrastructure for sustainable operations. The architecture and implementation approaches presented in this paper provide a foundation for industry practitioners and researchers working to reconcile the competing demands of technological progress, economic efficiency, and environmental responsibility.

REFERENCES

1. Boyd, S. B., Horvath, A., & Dornfeld, D. (2021). Life-cycle energy demand and global warming potential of computational logic. *Environmental Science & Technology*, 55(12), 8282-8291. <https://doi.org/10.1021/acs.est.1c00711>
2. Dayarathna, M., Wen, Y., & Fan, R. (2016). Data center energy consumption modeling: A survey. *IEEE Communications Surveys & Tutorials*, 18(1), 732-794. <https://doi.org/10.1109/COMST.2015.2481183>
3. Hasan, S., Bergés, M., Cutler, D., & Cohen, E. (2017). Energy-efficient scheduling of HVAC and battery systems under time-varying prices. *Energy and Buildings*, 155, 129-142. <https://doi.org/10.1016/j.enbuild.2017.09.019>
4. Hu, L., Tian, Q., Zou, C., Huang, J., Ye, Y., & Wu, X. (2016). A study on energy efficiency of data centers. *Energies*, 9(2), 133. <https://doi.org/10.3390/en9020133>
5. IEEE Standards Association. (2018). IEEE guide for smart grid interoperability of energy technology and information technology operation with the electric power system (EPS), end-use applications, and loads (IEEE Standard 2030-2018). Institute of Electrical and Electronics Engineers.
6. Jones, N. (2018). How to stop data centres from gobbling up the world's electricity. *Nature*, 561(7722), 163-166. <https://doi.org/10.1038/d41586-018-06610-y>
7. Masanet, E., Shehabi, A., Lei, N., Smith, S., & Koomey, J. (2020). Recalibrating global data center energy-use estimates. *Science*, 367(6481), 984-986. <https://doi.org/10.1126/science.aba3758>
8. Patterson, M. K., Azevedo, D., Belady, C., & Pouchet, J. (2019). Water usage effectiveness (WUE): A green grid data center sustainability metric. *ACM Transactions on Architecture and Code Optimization*, 16(4), 1-26. <https://doi.org/10.1145/3372392>
9. Piette, M. A., Kiliccote, S., & Dudley, J. H. (2015). Field demonstration of automated demand response for both winter and summer events in large buildings in the Pacific Northwest. *Energy Efficiency*, 8(4), 671-684. <https://doi.org/10.1007/s12053-014-9308-y>
10. Qureshi, A., Weber, R., Balakrishnan, H., Gutttag, J., & Maggs, B. (2020). Cutting the electric bill for internet-scale systems. *ACM SIGCOMM Computer Communication Review*, 50(2), 2-13. <https://doi.org/10.1145/3213232.3213235>
11. Radovanovic, A., Koningstein, R., Schneider, I., Chen, B., Duarte, A., Roy, B., ... & Sundarajan, R. (2022). Carbon-aware computing for datacenters. *IEEE Transactions on Power Systems*, 37(4), 2606-2617. <https://doi.org/10.1109/TPWRS.2021.3124549>

12. Schleich, J., Klobasa, M., Gözl, S., & Brunner, M. (2017). Effects of feedback on residential electricity demand—Findings from a field trial in Austria. *Energy Policy*, 108, 773-787. <https://doi.org/10.1016/j.enpol.2017.06.041>
13. Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 3645-3650. <https://doi.org/10.18653/v1/P19-1355>
14. Williams, E. D., Ayres, R. U., & Heller, M. (2002). The 1.7 kilogram microchip: Energy and material use in the production of semiconductor devices. *Environmental Science & Technology*, 36(24), 5504-5510. <https://doi.org/10.1021/es025643o>
15. Zhang, Y., Wang, Y., & Wang, X. (2018). GreenWare: Greening cloud-scale data centers to maximize the use of renewable energy. *IEEE Transactions on Sustainable Computing*, 3(2), 93-106. <https://doi.org/10.1109/TSUSC.2017.2713955>