

How Can Datacenters Join the Smart Grid to Address the Climate Crisis?

*Using simulation to explore power and cost effects
of direct participation in the energy market*

By

Hongyu He

(VU: 2632195, UvA: 12958794)

hongyu.he@vu.nl

1st supervisor: Prof.Dr.ir. Alexandru Iosup (VU Amsterdam)
daily supervisor: Fabian Mastenbroek (TU Delft)
3rd supervisor: Leon Overweel (Dexter Energy Services B.V.)
2nd reader: Dr.ir. Animesh Trivedi (VU Amsterdam)

Submitted in fulfilment of the requirements for the degree
of Bachelor of Science (Honours)
in Computer Science

VRIJE UNIVERSITEIT AMSTERDAM
Amsterdam, The Netherlands
July 2021

© 2021
Hongyu He
All rights reserved

*To my mom, who always reminds me to shave my beard off,
and to my grandparents.*

Acknowledgements

I would like to express my gratitude to my supervisors, Prof.Dr.ir. Alexandru Iosup, Fabian Mastebroek, Leon Overweel, and Dr.ir. Animesh Trivedi, for their guidance through each stage of this research, and for inspiring my interest in the development of innovative technologies. Also, I appreciate all the training and support received from the AtLarge Research group.

This thesis is the culmination of three-year education and research, marking the first milestone of my scientific contribution to the community and society.

Abstract

Amidst the climate crisis, the massive introduction of renewable energy sources has brought tremendous challenges to both the power grid and its surrounding markets. As datacenters have become ever-larger and more powerful, they play an increasingly significant role in the energy arena. With their unique characteristics, datacenters have been proved to be well-suited for regulating the power grid yet currently provide little, if any, such active response. This problem is due to issues such as unsuitability of the market design, high complexity of the currently proposed solutions, as well as the potential risks thereof. This work aims to provide individual datacenters with insights on the feasibility and profitability of directly participating in the energy market. By modelling the power system of datacenters, and by conducting simulations on real-world datacenter traces, we demonstrate the substantial financial incentive for individual datacenters to directly participate in both the day-ahead and the balancing markets. In turn, we suggest a new short-term, direct scheme of market participation for individual datacenters in place of the current long-term, inactive participation. Furthermore, we develop a novel proactive DVFS scheduling algorithm that can both reduce energy consumption and save energy costs during the market participation of datacenters. Also, in developing this scheduler, we propose an innovative combination of machine learning methods and the DVFS technology that can provide the power grid with indirect demand response (DR). Our experimental results strongly support that individual datacenters can and should directly participate in the energy market both to save their energy costs and to curb their energy consumption, whilst providing the power grid with indirect DR.

Index Terms—Datacenter modelling, cloud simulation, smart grid, energy market, demand response, frequency scaling, DVFS scheduling, machine learning.

Table of Contents

Nomenclature	vi
List of Algorithms	x
List of Figures	xi
List of Tables	xiv
1 Introduction	1
1.1 Problem Statement	4
1.2 Research Questions	6
1.3 Research Methodology	10
1.4 Thesis Contribution	11
1.5 Plagiarism Declaration	13
1.6 Thesis Structure	13
2 Background	14
2.1 Metrics	14

2.2	Energy Proportionality	18
2.3	DVFS as Mechanism to Manage Energy	20
2.4	Power Grid	28
2.5	Energy Modelling for Datacenters	41
3	Energy Modelling & Management	45
3.1	Development Pipeline	45
3.2	Requirement Engineering	46
3.3	System Architecture	57
4	Implementation	64
4.1	Energy Modelling & Management System	65
4.2	Market Extension	72
5	Evaluation	75
5.1	Experiment Setup	75
5.2	Energy Market	82
5.3	DVFS Scheduling	98
6	Conclusion	118
6.1	Answers to Research Questions	118
6.2	Limitations	120
6.3	Future Work	121
6.4	Summary	122
	References	123

Nomenclature

α	Proportional loss coefficient
β	Square-law loss coefficient
η_e	Percentage of PSU energy efficiency
η_l	Percentage of PSU load
λ	Tare loss coefficient
\mathbb{P}	Power function
\mathbb{Q}	Quantile function
$\mathbb{1}$	Indicator function
π	Nameplate loss coefficient
τ	Rated output power of PSU
ζ	Locational marginal pricing
a	Asymptotic Factor
C	Capacitance
c	Number of (logical) cores of a machine
E	Additive Gaussian random variable
E_{PSU}	Accumulative energy consumption of PSU

F	CPU capacity/frequency
f	Instant CPU frequency
f_{\max}	Maximum scaling frequency
f_{\min}	Minimum scaling frequency
G	Financial gain
m	Size of memory unit
N_{PDU}	Number of attached PDUs
N_{server}	Number of active servers
N_{blocked}	Number of tasks blocked by I/O
N_{BRP}	Number of participated BRPs in the balancing market
N_{BSP}	Number of participated BSPs in the balancing market
N_{ISP}	Number of ISPs
N_{queued}	Number of tasks in the run queue
N_{run}	Number of tasks that are being processed
P	Power draw
P^{compute}	Power draw of computing equipment
P^{IT}	Power draw of IT infrastructure
P^{server}	Power draw of server
P^{total}	Power draw of datacenter
$P^{2\text{nd}}$	Power draw of secondary power support
P^{idle}	Idle power
P^{in}	Inlet power

P^{loss}	Total power loss
P^{max}	Maximum power
P^{rated}	Nameplate power
P^{tare}	Tare power loss
p_+^B	Surplus price in the balancing market
p_-^B	Shortage price in the balancing market
p^F	Forecasted price
p^f	Synthetically forecasted price
p^S	Spot market price
p^S	Spot price in the day-ahead market
Q_+	Quantity of energy required by downwards regulations
Q_+	Quantity of surplus energy of a BRP
Q_-	Quantity of shortage energy of a BRP
Q_{\downarrow}	Quantity of energy required by downwards regulations
Q_{\uparrow}	Quantity of energy required by upwards regulations
Q_d	Quantity of energy delivered
Q_s	Quantity of energy scheduled in the spot market
r	Approximation factor
s	Scalar value
t_{CPU}	CPU time
t_{idle}	Accumulative CPU idle time
t_{wall}	Wall time

U Utilization of IT equipment

u CPU usage

u_{os} CPU utilization

V Voltage

List of Algorithms

1	P-state scaling algorithm	66
2	Scaling algorithm in the <code>OnDemandScalingGovernor</code>	67
3	Scaling algorithm in the <code>ConservativeScalingGovernor</code>	69
4	Max-min fair-sharing power distribution algorithm	72
5	DVFS scheduling algorithm	74

List of Figures

1.1	Functions of the smart grid	2
1.2	Demand response in the smart grids	3
1.3	Characteristics of datacenters	5
1.4	Aerial view of the project.	7
2.1	Transitions between the power states defined in the ACPI standard	21
2.2	P-states and corresponding power consumption levels	23
2.3	The ondemand & conservative governors in the Linux kernel .	24
2.4	Challenges faced by energy transmission and trading	28
2.5	Markets around the power grids	29
2.6	Types of energy markets in the EU	30
2.7	Sequence of energy markets	31
2.8	Duck curve showing the net load of the power grid	34
2.9	Solar power generation over three consecutive days	35
2.10	Energy flow in datacenters	38
2.11	Approximate distribution of energy usage in a datacenter	41

2.12	Approximate distribution of energy losses in IT equipments	42
3.1	Development pipeline	46
3.2	Use case diagram of the system	52
3.3	Overview of the architecture of the entire system	56
3.4	Architecture of the modelled power system of datacenters	59
3.5	Architecture of the power support subsystem	61
3.6	Architecture of EEMM	63
4.1	UML diagram of the energy modelling and management system	68
4.2	UML diagram of the power support subsystem	71
5.1	Different power loads	83
5.2	Distributions of day-ahead prices and imbalance prices	84
5.3	Energy Costs of the two power loads in different markets	85
5.4	Comparison of CPU energy efficiency of two machine models	86
5.5	Comparison of power loads of the two machine models	87
5.6	Comparison of energy costs of the two machine models	88
5.7	Simulated load schedules in the day-ahead market	89
5.8	Comparison of energy costs of simulated load schedules	90
5.9	Comparison of energy costs of the two imbalance pricing system	92
5.10	Stacked comparison of energy costs of the imbalance pricing systems	94
5.11	Correlation between imbalance prices	95
5.12	Correlation between day-ahead prices and imbalance prices	96
5.13	Potential benefit of obtaining predictions for imbalance prices early	97

5.14	Relationships between CPU demand and usage of DVFS	98
5.15	PDF and CDF of the CPU usage of different governors	99
5.16	Detailed distributions of CPU usage of different governors	100
5.17	Comparison of power estimation of the two models for governors .	101
5.18	Comparison of instant power draw of different governors	102
5.19	Instant CPU over-commission level of different scaling governors	103
5.20	Over-commission at different CPU usage levels	104
5.21	Instant power draw at different over-commission levels	105
5.22	Comparison of granted work and instant power draw	106
5.23	Effect of the damping factor on energy use and over-commission .	107
5.24	Timeline of ML prediction	108
5.25	Distributions of Gaussian noises with different σ values	110
5.26	AA scores of ML and synthetic predictors	111
5.27	Energy and over-commission of ML and synthetic predictors . . .	112
5.28	Energy costs of ML methods and synthetic predictors	113
5.29	Demonstration of indirect demand response	114
5.30	Comparison of total energy consumption with DVFS scheduling .	115
5.31	Comparison of total CPU over-commission with DVFS scheduling	116
5.32	Comparison of total energy costs for DVFS scheduling	117

List of Tables

2.1	Balancing reserves in the EU	33
2.2	Overview of the nine surveyed datacenter simulators	44
3.1	Potential stakeholders	47
3.2	Stakeholder classification	48
3.3	Experts involved in requirement verification and validation	55
5.1	On-demand energy prices considered in the experiments	76
5.2	P-states consumption levels of the old machine model	77
5.3	Machine models	77
5.4	Host setup for each machine model	78
5.5	Load-to-power data of the new machine model	79
5.6	Coefficient values of the UPS and PDU power model	80
5.7	Symbols used in defining the two balancing systems	91
5.8	Symbols used in defining the AA score	109

Introduction

Taking their toll at a rising speed over the past decades, environmental problems such as global warming have fast become a worldwide focal point [22], so much so that even the COVID-19 pandemic can barely arrest the development of this alarming trend [185]. To combat this exacerbating issue, remedies such as emission limits, carbon taxes, and perhaps most importantly, ambitious renewable energy targets adopted in 2015 [123] and enhanced in 2021 [126], have been introduced. As a result of such a push towards sustainability, the energy market has become increasingly volatile, bringing both challenges and opportunities to the ever-larger energy consumers, datacenters.

Smart Grid. Despite the fact that there is virtually no cost in the production of renewable energy as they are free of charge from nature, the substantial operational costs induced by its huge intermittency and stochasticity, however, greatly impedes the continuous penetration of renewable energy sources into the power grid [103]. Consequently, large numbers of expensive and carbon-intensive system operating reserves, which hinge on more reliable energy sources like petroleum or even diesel, are often required as hot/cold standby reserves to back up renewables in order to maintain the equilibrium of the power grid (§2.4). In addressing the above challenges, the power grid is becoming more and more intelligent (Figure 1.1) — the smart grid [49]. Such recent advances in functions of the smart grid enable real-time, fluent interactions and coordination between energy producers and

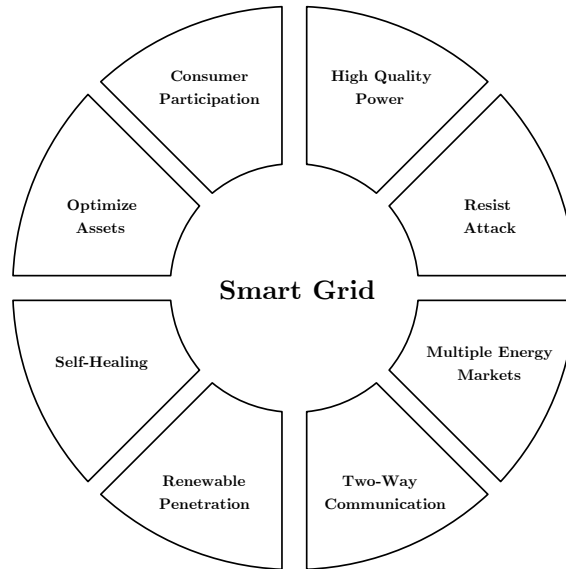


Figure 1.1: Functions of the smart grid.

consumers, improving demand-side management (DSM) [168, 38].

Demand Response. The smartness of the grid, however, is hitting some limits due to the uncertainty caused by the massive introduction of renewable energy sources. One relatively recent, yet promising form of DSM is demand response (DR), which has been extensively explored in existing literature (*e.g.*, [72, 186, 84, 93, 184, 153, 2]). In general, based upon the response time, DR programmes can be classified into two broad categories, direct and indirect control. The direct approach responds to requests and signals from the power grid quickly, which provides system operators with accurate and fast control over the power grid (different levels of DR are introduced in Section 2.4.2). Thus, most of the ancillary services opt to direct DR control. Conversely, it is computationally and communicatively more intensive. In contrast, indirect DR is cheaper and more flexible but embodies a greater degree of uncertainty. As opposed to relying upon the direct requests from the system operators or upon constant monitoring of the power grid, it financially incentivizes prosumers to modify their energy consumption and/or production following signals from the energy markets and/or local utilities. Al-

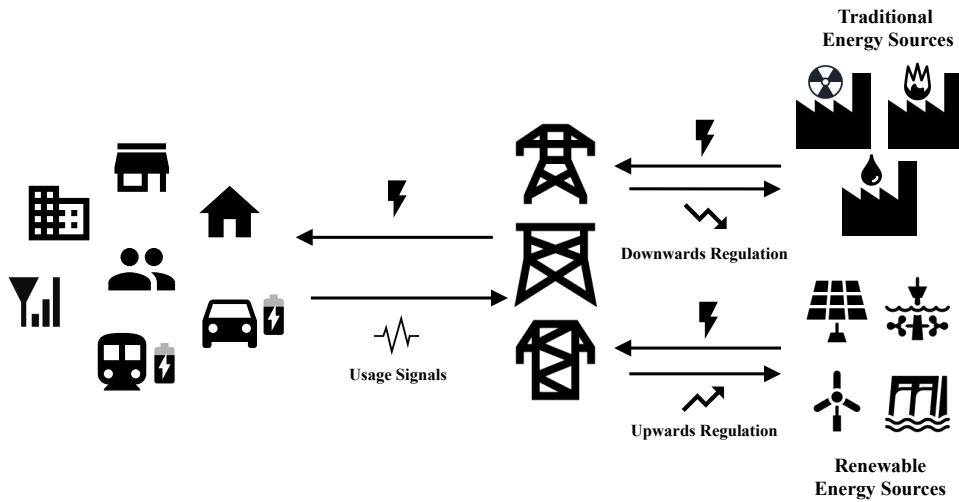


Figure 1.2: Demand response in the smart grids.

though unlike the frequency level of the power grid that directly reflects the status of the current demand, energy prices of the financial markets are good indicators and predictors of the demand-supply balance. Thus, indirect DR has the potential to effectively reduce the market volatility caused by wind and solar energy generation, mitigate network problems (*e.g.*, congestion, voltage), and respond to failures (*e.g.*, avoiding blackout), helping both the customers and the power grid dramatically reduce costs [145]. Although many grids currently do not have congestion problems, as more distributed solar generation and more erratic consumption resulted from mobility, *e.g.*, electrical vehicles (EVs), are introduced, a large fraction of the distribution grids is expected to suffer from congestion in the near future [70]. To tackle this, intensive research has been conducted to explore the interrelationship between the mobility of EVs and DR, *e.g.*, using vehicle-to-grid (V2G) technology [172, 89, 169]. Throughout the years, people have also resorted to refrigerators [117], fans [166], laundry machines and dishwashers [133, 129], or even boilers [45], searching for DR resources to ameliorate the intensive demand during peak times. During these periods, the energy delivery is much more costly since less efficient energy plants (usually powered by fossil fuels) are employed.

1.1 Problem Statement

Over the past decades, datacenters have grown ever larger and more powerful [118]. According to recent reports, datacenters consume more than 1.5% of global electricity use [100, 118, 178], and, with a growth rate of 10–12% per year [160, 19, 58, 95], this figure is expected to grow to 3% by 2025 [162]. In the US, about 2.2–3.5% electricity use can be attributed to datacenters [95]. Similar phenomena appear in the Netherlands, where the power consumption of datacenters amounts to 2.7 billion kWh by 2020, which is 2.3% of the total Dutch energy use [43]. A decade ago (2011), Google datacenters already use almost 260 MW of power, exceeding the consumption level of Salt Lake City [131]. Similarly, a single datacenter of Microsoft in Washington DC consumes 48 MW of power, which is the equivalent of around 40k households in the US [168]. Due to their carbon footprints and energy consumption, datacenters have become a front-line target in the battle against climate change [64].

1.1.1 Opportunities

The increasing capacity of datacenters may well be a blessing in disguise for the power grid. A well-managed datacenter of 30 MW has approximately the same capability for regulating the power grid as huge energy storage of 7 MWh [172]. Besides its capacity, unlike other subjects such as lighting and residential power, datacenters feature its elastic load. In other words, datacenters are capable of curbing their power demand without service degradation. Also, many Cloud service providers employ spot pricing mechanisms to manage demands (*e.g.*, [73]). To put this into perspective, a study from Lawrence Berkeley National Laboratory (LBNL) shows that 15% of the load can be shed within 15 minutes without adjusting temperatures or any other building managements [58]. Additionally, the flexibility of datacenters can be leveraged at a finer-grained level through energy-saving techniques such as power capping [32] and dynamic voltage & frequency scaling (DVFS) [99]. Note that DR is not mean to save energy, but DVFS can. This enabling configuration will be described in detail later in Section 2.3. Furthermore, datacenters are built for extremely reliable and available services. For instance,

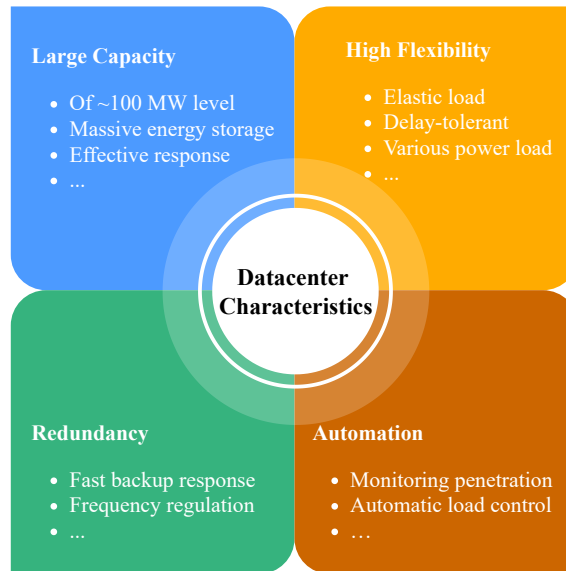


Figure 1.3: Characteristics of datacenters.

two availability classes to which most modern datacenters belong are Tier III and Tier IV. Datacenters that fall into the former category insure a 99.982% availability, and for the latter, the figure is 99.995% [114]. To accomplish promised uptime and performance guarantees, datacenters generally have a considerable amount of redundancy throughout their power systems, as well as large battery capacities in their primary power support. Lastly, datacenters are complex but highly automated systems. Ubiquitous monitors and controls empower datacenters' participation in the DR programme. Therefore, datacenters are well-suited candidates for DR programmes [31, 172].

1.1.2 Challenges

Albeit with great potentials of taking part in DR programmes, datacenters nowadays provide the power grid with little, if any, response to the power grid [58, 60, 109, 108]. Firstly, the current market designs [85, 176] and DR programmes [7, 145] (i) are not particularly suitable for datacenters and (ii) can barely fully ex-

tract the flexibility of datacenters [109]. Secondly, datacenters may incur charges if no response is offered during DR programmes. No profit would be produced either if the peak in energy demand did not happen during the coincident periods. Thirdly, the proposed DR strategies for datacenters often require cooperation between the energy market, geographically distributed datacenters, and their utility companies [138, 110, 57, 170, 182, 168, 135, 138, 110, 106, 57, 170, 182]. In other words, they can hardly be carried out without structural changes to the energy market and/or substantial adjustments in datacenter operations. Hence, the complexity of the orchestration and the potential risks therein hinder the participation of datacenters in the power grid. Furthermore, experimenting, testing, and evaluating energy-aware techniques tend to be costly or sometimes even unrealistic in large-scale, modern datacenters. This could lend hesitancy for datacenter shareholders to embrace energy-saving techniques, *e.g.*, DVFS since the critical guarantee of performance and availability specified in their service-level agreement (SLA) always takes precedence over *unknown* benefits brought by enabling energy-saving configurations. This challenge further imposes potential risks on datacenters when participating in the DR programmes. One way of bridging this gap is to measure energy consumption at the hardware level, for example, installing system/component power meters to monitor the power usage [27, 51, 39]. Although such a direct approach is becoming a common practice, it only applies to facilities that have already been built [47]. Consequently, it is not portable, scalable or informative for future planning and responsive scheduling. Another method is to employ software instruments, specifically, datacenter simulators, to model the energy consumption of thousands of servers as well as non-IT infrastructure, *e.g.*, cooling and ancillary equipment. In comparison to methods at the hardware level, datacenter simulators are more flexible, reproducible and cost-effective, playing an instrumental role in facilitating energy-aware decision-making [80]. As a result, datacenter simulation has been widely adopted for both academic research and industrial use [8, 143].

1.2 Research Questions

As the first step towards addressing the aforementioned challenges, this work aims to provide individual datacenters with insights on the feasibility and profitability of

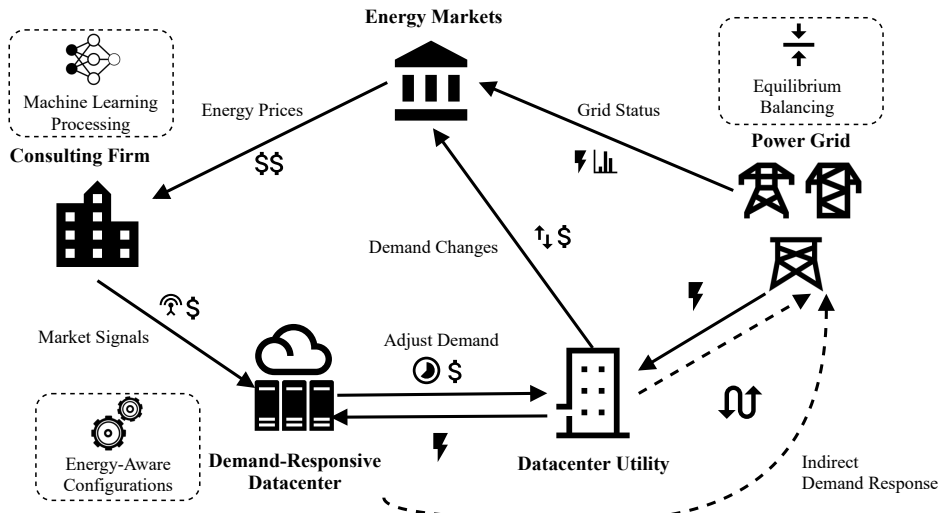


Figure 1.4: Aerial view of the project.

directly participating in the energy market whilst offering indirect DR to the power grid (Figure 1.4 demonstrates an aerial of this project). To achieve this objective, we put forwards the statement of this thesis:

Thesis Statement — Individual datacenters can and should directly participate in the energy market both to save their energy costs and to curb their energy consumption, whilst providing the power grid with indirect DR.

To offer evidence for the thesis statement, we raise the main research question (**MRQ**) followed by a sequence of research questions (**RQs**).

MRQ How feasible and beneficial is it for individual datacenters to directly participate in the energy market whilst providing the power grid with indirect DR?

1.2.1 Research Question 1

To answer the **MRQ**, we need to first estimate the energy consumption of datacenters. There are, however, a myriad of factors that influence the energy consumption

of datacenters. To name a few such factors: different topologies of the datacenters, various types of hardware as well as many configurations thereof. In an effort to overcome this challenge, we resort to whole power system modelling in order to capture the heterogeneity of those factors. To this end, an energy resource chain, starting from the power source and ending at the IT infrastructure, is needed. This resource chain passes through several supporting subsystems and components of datacenters, for example, the automatic transfer switches (ATSs), the uninterrupted power supply (UPS) systems, and the floor/rack power distribution units (PDUs). Inside the servers, power supply units (PSUs) transform AC power to DC power (in a typical AC architecture). A detailed introduction to these various components can be found in Section 2.4.5. Besides these hardware components and subsystems, fine-grained manoeuvres happen in the CPUs, adjusting the power consumption in real time. One of such techniques is the DVFS, an active power management technique whereby the frequency of a microprocessor can be automatically adjusted on spot based upon the computing loads. By exploiting DVFS, machines can save energy and, in turn, reduce operational costs (an elaborate introduction of DVFS can be found in section 2.3). Neither building the resource chain nor incorporating and synergizing various DVFS policies is a straightforward task, which gives rise to the first research question:

RQ1 How to model the power system of datacenters?

1.2.2 Research Question 2

IN **RQ1**, we estimate the energy consumption of datacenters through simulation. Now we are interested in the extent to which datacenters are able to benefit from directly participating in the energy market. We focus on benefits resulting from taking part in different markets, specifically, the day-ahead and the balancing markets. We do not, however, quantify the savings for the electricity grid resulted from a lower need to generate energy from fossil fuels and/or storing it for datacenter use, or the operational changes in the electricity grid (*e.g.*, supply curtailment [142]). Should datacenters participate? If so, which market(s) should be given particular attention? To answer these questions, we post the second research question:

RQ2 Is it beneficial for datacenters to participate in the energy market in the first place?

1.2.3 Research Question 3

RQ2 seeks answers to the question of why datacenters should participate in the energy market. The next step is to inquire about the most economical way of procuring energy in the energy market, *i.e.*, how to participate. Referring back to Section 1.1.1, reliability and availability are of paramount importance for datacenters. In fact, load forecast (*e.g.*, [6, 159, 54, 148, 86]) is commonly resorted for energy-planning in datacenters. Therefore, load-forecast-based procurement strategies are of particular interest. In turn, we ask the third research question:

RQ3 How to procure energy in the energy markets according to forecasted power load?

1.2.4 Research Question 4

Leveraging different procurement strategies in **RQ3** may bring in substantial profits for datacenters, the energy consumption level, however, cannot be improved by only doing so. To provide indirect DR whilst reduce the carbon footprint, we turn to the novel orchestration between machine learning (ML) methods and the low-level energy-saving technique DVFS. We adjust the DVFS policies in accordance with the predicted market signals produced by ML models. Moreover, since almost every new advancement in computer science comes at a cost, mitigating the overhead introduced by using DVFS when responding to market signals is of the same importance. These challenges beg the fourth research question:

RQ4 How to optimize energy consumption when participating in the energy market using DVFS, based upon ML methods?

1.2.5 Research Question 5

There is much to explore when it comes to the interrelationship between datacenter simulation and the energy markets/power grid. However, as the old saying goes, “difference in the profession makes one feel worlds apart”. To assist further explorations and bridge the gap in domain knowledge, we raise the last research question:

RQ5 How to create an exploratory tool for problems in this domain, to be used by experts in both the IT and the energy industry?

1.3 Research Methodology

To achieve **RQ1**, we employ quantitative research, specifically, system modelling and simulation [92]. We design [79, 132] and prototype [66] the energy modelling and power management subsystem of OpenDC in which a set of power models and a number of generic DVFS algorithms are built and integrated. As a result, these components should be integrated and work in concert. On the basis of these models and management techniques, we model and simulate the power system of a typical datacenter.

To answer **RQ2**, we carry out experimental research [83], conducting discrete-event simulations on real-world datacenter workload, quantifying and comparing the energy costs of participating in different markets. Note that the prediction of the workloads in production will be more precise as the time approaches the energy delivery period. However, we do not take into account such an increasing precision in workload predictability since it is the controlled variable in this work.

To address **RQ3**, we conduct a case study, collaborating with partners and experts in the energy market. We obtain the predicted energy prices of the intraday market as well as details concerning trading and operations in the energy market. We regard the predicted prices of the balancing market as indicators of energy demand, with which we adjust the fractions of energy bought in during the day-ahead market period. By doing so, we leverage energy cost whilst still sufficiently satisfy energy needs in the balancing market.

To attack **RQ4**, we further employ various frequency scaling algorithms in response to the market signals, developing a DVFS scheduling algorithm powered by ML methods. By designing workload-level benchmarks [69, 127], we conduct various experiments using trace-based simulations, focusing on the impact of various data-center phenomena, such as specialized/mixtures of scenarios and correlated forms of performance variability.

To tackle **RQ5**, we honour open-science guidelines [16, 174] and build open-source scientific software, following rigorous software engineering methods and PR-review software development cycle [180]. Agreeing on various specifications with our partners and experts, we adhere to the standard format of market data from official websites and containerize the deployment of our research instruments. We make our OpenDC datacenter simulator together with its market extension an out-of-box tool that is ready to be used by experts in both the energy and the IT industry.

1.4 Thesis Contribution

By addressing all research questions with our best efforts, we endeavour to transfer our knowledge and experience/lessons learned as well as to deliver developed software and experimental results to the community as much as possible without reservation. In this section, we list our scientific contributions (**TC**) as well as technical contributions (**TC**). Also, we identify potential societal impact (**SI**) and economic impact (**EI**) of this research. Additionally, we honour the FAIR data principles[174] (**FAIR**), exercising the best practices for sharing data. The contribution of this work is nine-fold:

1. We are the first to demonstrate the substantial financial incentive for individual datacenters to directly participate in both the day-ahead and the balancing market **{SC, EI}**.
2. We suggest a new short-term, direct scheme of energy market participation for individual datacenters, in place of the current long-term, inactive market participation **{SC}**.

3. We develop a novel proactive DVFS scheduling algorithm that is able to both reduce the energy consumption **{SI}** and save the energy cost of datacenters **{EI}**.
4. We propose an innovative combination of machine learning methods and the energy-management technique DVFS **{SC}** that can provide the power grid with indirect DR in an effort to overcome the increasing challenges brought by renewable energy sources **{SI}**.
5. We are the first to achieve whole power system modelling in datacenter simulation **{TC}**.
6. We create a user-friendly and ready-to-use tool for experts in both the IT and the energy industry to further explore the research potential lies at the intersection between the two fields **{TC}**.
7. We publish our code as open-source projects, facilitating future scientific explorations and collaborations **{TC, FAIR}**.^{1, 2}
8. Alongside the code, we publish our datasets and raw experimental results, which, in turn, are findable, accessible, and do not subject to any ethical, legal or contractual restrictions **{FAIR}**.
9. Besides the datasets, we build documentation and make tutorials with examples to ensure reproducibility, interoperability, and reusability of both the software and the data **{TC, FAIR}**.^{3, 4}

¹<https://github.com/atlarge-research/opencdc>

²<https://github.com/hongyuhe/opencdc-eemm>

³<https://opencdc-eemm.rtfid.io>

⁴<https://opencdc.org>

1.5 Plagiarism Declaration

I confirm that this thesis work is my own work, is not copied from any other source (person, Internet, or machine), and has not been submitted elsewhere for assessment.⁵

1.6 Thesis Structure

Firstly, key terms pertinent to this work are covered in Chapter 2. Then, in Chapter 3 we present the design of the energy modelling and management system, including its subsystems and the market extension. In addition, it also introduces the development pipeline and requirement engineering process. Next, we detail the implementations of the system in Chapter 4. After that, in Chapter 5 we employ the developed infrastructure and tools to conduct experiments. In the last chapter, we answer research questions, summarizing our key findings and results. Lastly, we identify limitations and envision future research (Chapter 6).

⁵<https://www.vu.nl/en/about-vu-amsterdam/academic-integrity/index.aspx>

Background

An overview of related subjects is laid out in this chapter. Key terms covered include the metrics used for evaluating energy consumption and for the saturation of critical loads (§2.1), energy proportionality and existing solutions (§2.2), frequency scaling (§2.3), and the power grid (§2.4), especially, the architecture of the power system in datacenters that we simulate is also presented (§2.4.5). Last but not least, we discuss the energy modelling for datacenters in Section (§2.5).

2.1 Metrics

PUE & CPE. Two fundamental energy metrics are commonly used in datacenters, a) Power Usage Effectiveness (PUE) first proposed by Malone and Belady [113] (Equation 2.1), which can be used for both benchmarking and energy estimation, and b) Compute Power Efficiency (CPE) for measuring the computational efficiency of datacenters (Equation 2.2).

$$\text{PUE} = \frac{P^{\text{total}}}{P^{\text{IT}}} \quad (2.1)$$

$$\text{CPE} = \frac{U}{\text{PUE}} = \frac{U \cdot P^{\text{IT}}}{P^{\text{total}}}, \quad (2.2)$$

where P^{total} is the total facility power, P^{IT} is the power draw of the IT infrastructure, and U denotes the utilization of the IT equipment.

Reflecting on the continuous growth of energy consumption from 1992 to 2014, and the raw performance per watt that is radically doubling every year, Malone and Belady point out that this trend is unlikely to stop in the next years. Also, they suggest a significant change in the past decades that the major cost of the datacenters is moving from operation-driven [177] to an infrastructure- and energy-oriented model.

The average PUE value of datacenters worldwide is still close to 2.0 in 2021 [77]. This is a clear indication of widespread energy-inefficiency in the sense that to produce 1 W of computational power, around 2 W of power is consumed by supporting infrastructures such as cooling and other ancillary facilities. In 2008, only 0.4 improvements of PUE can result in about 350,000\$ reduction in energy cost, which is equivalent to about 430,000\$ in 2019. To better capture the power usage of IT infrastructure, the authors propose the adoption of CPE that directly reflects the actual fraction of energy used for computing. Note that a slight increase in PUE can boost the CPE significantly. For example, a typical PUE of 2.0 for a well-managed datacenter has a CPE of 10%, whilst a PUE of 1.6 corresponds to a CPE of 50%. Thus, CPE is a relatively more sensitive and intuitive indicator of energy losses for both the IT and supporting infrastructure.

TUE & ITUE. Patterson et al. [130] developed IT-power usage effectiveness (ITUE) and total-power usage effectiveness (TUE) to improve PUE and CPE respectively, by taking into account the power distribution and cooling losses inside IT equipment [130]. These metrics aim to tackle the challenge of estimating and tracking the total efficiency of the entire energy stack for datacenters. The computation of the metrics are captured by Equation 2.3 and 2.4:

$$\text{ITUE} = \frac{P^{\text{total}}}{P^{\text{compute}}} \quad (2.3)$$

$$\text{TUE} = \text{ITUE} \cdot \text{PUE}, \quad (2.4)$$

where P^{compute} is the power draw of only the equipment related to computing service.

Resource Utilization Metrics. The debate of performance metrics is known as one of the “Rat Holes” in systems research [82]. Gauging the saturation of the CPUs is never a trivial task, and over time, a number of metrics have been widely used in various contexts. In this section, we focus on three commonly employed metrics, namely, CPU load, CPU usage, and CPU utilization. First and foremost, it is worth noting that the definition of CPU load/usage/utilization varies with different use cases, platforms, and organizations [164, 76, 41, 173]. Sometimes they are even used interchangeably by many people, albeit quite different. This is partially due to the fact that there is no single, universally standard way of defining these metrics. Nevertheless, the CPU load (l) is generally used in the context of the Linux scheduler, in which the run queues of the processors are accessible. The calculation of l along the lines of Equation 2.5; computing the average CPU load (l_a) is also a common practice (Equation 2.6).

$$l = N_{\text{run}} + N_{\text{queued}} + N_{\text{blocked}} \quad (2.5)$$

$$l_a = \frac{l}{c}, \quad (2.6)$$

where N_{run} is the number of tasks that are being processed, N_{queued} is the number of tasks in the run queue, N_{blocked} is the number of tasks blocked by I/O, and c is the (logical) core count of the machine.

When it comes to CPU utilization/usage, the definitions often become a bit more nebulous. In this work, we distinguish these two in accordance with different viewpoints in the context of server architecture. With regard to the host, neither the run queue nor scheduler is visible at the firmware level. Therefore, from the standpoint of bare-metal machined, it is impractical to evaluate any metrics other than the ratio of the current CPU speed to the CPU capacity. In our case, we simulate Type-I hypervisors directly on top of the physical interface of the host machines. Hence, we employ the concept of the CPU usage (u) illustrated by Equation 2.7

throughout our energy modelling and management system.

$$u = \frac{f}{F}, \quad (2.7)$$

where f is the instant CPU frequency, and F is the CPU capacity.

In contrast, from the virtual machine (VM) or operating system (OS) point of view, the run queue, the scheduler and the wall time for every task are available. In turn, evaluating the CPU load is a feasible work. Thus, in the case of this work, we specify the CPU utilization *at the VM/OS level* following Equation 2.8. In addition, the average CPU load (also known as Per-entity Load Tracking [40]) is used in the `schedutil` scaling governor, and the concept of CPU utilization u_{os} (called ‘‘CPU load’’ in the context of the Linux kernel) is used in other governors.

$$\begin{aligned} u_{os} &= \frac{t_{\text{CPU}}}{t_{\text{wall}}} \\ &= 100\% - \frac{t_{\text{idle}}}{t_{\text{wall}}}, \end{aligned} \quad (2.8)$$

where t_{CPU} denotes the CPU time, t_{wall} denotes the wall time, and t_{idle} is the accumulative time in which the CPU is idle.

Besides the available abstraction offered by the current architecture of our infrastructure, another reason for using the CPU usage rather than the other two metrics is that the time measurement in simulation is generally at a coarser granularity than that of run-time systems in the real world. In other words, datacenter simulations are of a high abstraction of the real-world scenarios in that the traces, on which the experiments are conducted, usually have an interval of several seconds or even minutes between records. Conversely, timing is critical in gauging the CPU utilization, so much so that pitfalls often occur if the interrupts were to happen at undesirable points in time [37]. In addition, when assessing the CPU utilization, the Linux kernel does not take account of the total CPU capacity and the actual CPU speed since both of which are out of reach at the OS level; they are, nevertheless, available through the hardware interface in our simulation infrastructure. Hence,

we decide to take advantage of the resources available in our instrument, basing the following development of our energy modelling and management system upon the CPU usage.

2.2 Energy Proportionality

The inclinations demonstrated in previous reports [20, 96] show that the rising energy consumption steadily dominates the total cost of ownership (TCO) including both computing and infrastructure costs. Barroso and Hölzle [10] claim that, in order to prevent energy footprints of datacenters from exploding, the improvement of energy efficiency should keep up with the growth of computing power. From the result of their study [113], datacenters' raw performance has been growing five times faster than their performance-watt ratio. Moreover, Malone and Belady demonstrates a fatal mismatch of datacenters in which the most used working mode often runs in the least energy-aware way with fairly low CPU utilization. To combat this challenge, energy proportionality, *i.e.*, datacenters consume nearly no energy when standing by and gradually raise power consumption as workloads increase, is proposed as an ultimate design goal. Consequently, operating datacenters at the near-peak performance level with a high utilization is preferred as the higher utilization, the better energy efficiency.

That been said, energy proportionality is not an easily achievable objective. By the nature of distributed systems, not only computation but data is allocated amongst hundreds of, if not thousands of, geographically distributed nodes. The purpose of such an architecture is to create duplications, increasing service availability and, in turn, reducing risks caused by disastrous situations. One of many services that hinge on such an architecture is the Google File System (GFS) [59]. Conversely, such settings come at a cost — distributed servers are expected to be always up and running, even if no heavy workloads are hosted [113]. As a result, servers in a datacenter are often neither fully idle nor operates at their maximum utilization. Instead, they mostly operate in the 10 to 50% utilization range [144]. Moreover, as Malone and Belady have suggested, because of the need for performing constant, small operations in the background, networked nodes can hardly enter deep sleep states. When servers are running at their lowest operational states, they

consume more than half of their full power [113] (approximately 70% of their full-speed power [137, 98, 161, 55]). In addition, the impact of the energy loss incurred during the wake-up stage is not as significant compared to the energy consumption under normal utilization level. These characteristics can rarely be found in other systems, such as mobile and embedded systems. Besides the peak-power range, energy efficiency, therefore, should be optimized at all frequency steps.

With the growing energy consumption of modern datacenters and the increasing concern of global warming, more and more studies are being conducted in the search for sustainable solutions. With regard to the energy efficiency of modern datacenters and cloud systems, a full taxonomy has been built by Beloglazov et al.. Methods for energy management are generally categorised as two major types, a) static power management (SPM) and b) dynamic performance scaling (DPS) [15, 80]. In each of them, both the hardware and the software solutions have their important role to play [104]. The major objectives thereof are two-fold: (1) improving hardware design such as energy-efficient computing as well as models of cooling systems [62, 152, 147], and (2) creating better resource management algorithms, including workload scheduling [5, 3, 122], policies of power management [120], etc.

Beloglazov et al. investigated the solutions at a more fine-grained level. They separate the case studies into four levels namely, the hardware and firmware level, the OS level, the virtualization level and the datacenter level. To achieve a better TOC and to speed up Returns On Investments (ROI), and most importantly, to mitigate the carbon footprints of datacenters and cloud services, an integrated approach detailed at each of the four levels is proposed [15]. Similarly, Chien et al. proposed a set of computing models, the Zero-Carbon Cloud (ZCCloud) that features a bottom-up approach from the selection of sustainable sites to high-level infrastructure design [35].

2.3 DVFS as Mechanism to Manage Energy

In computing machines, different hardware components usually operate at various states, and generally speaking, binary states (*e.g.*, active and inactive) are applicable to most of these components. The idle state and inactive state should be clearly distinguished, as they are rather different from each other. When a machine is idle, it operates at its lowest active state without undertaking any useful work, whereas, in the inactive states, it is in one of the standby or in sleep states [68]. Note that, in some literature, the inactive states are regarded as idle states whilst, in this work, the idle state is referred to as the lowest available working state. Also, CPU frequency scaling introduced in this section is an active power management technology.

2.3.1 Frequency Scaling

Operational states differ from component to component. For example, memory could be in states such as precharging, refreshing, writing, reading, etc., whilst for network switches, the ports therein can operate at various rates, which correspond to another set of states different from that of the memory. Each of these operational states corresponds to a different energy state, *i.e.*, their energy efficiency differs in accordance with the type of the operations and/or the speeds of that operations are carried out. This, in turn, makes the energy management of the CPUs relatively more intricate.

$$P \propto C \cdot V^2 \cdot F + P^{\text{idle}} \quad (2.9)$$

Modern CPUs are capable of running at various speeds/frequencies, each of which has a certain rate of power consumption; Equation 2.9 illustrates this well. The power consumption (P) of a CPU is proportional to the capacitance (C), the voltage (V) and the clock frequency (F) of the CPU, whereas its idle power P^{idle} is an additive term. Clearly, the frequency is linearly correlated with the power consumption, whilst the voltage has a quadratic correlation with it. However, this might give the impression that halving the CPU frequency will make the running time of the hosted workloads twice as long, whilst they still consume the same

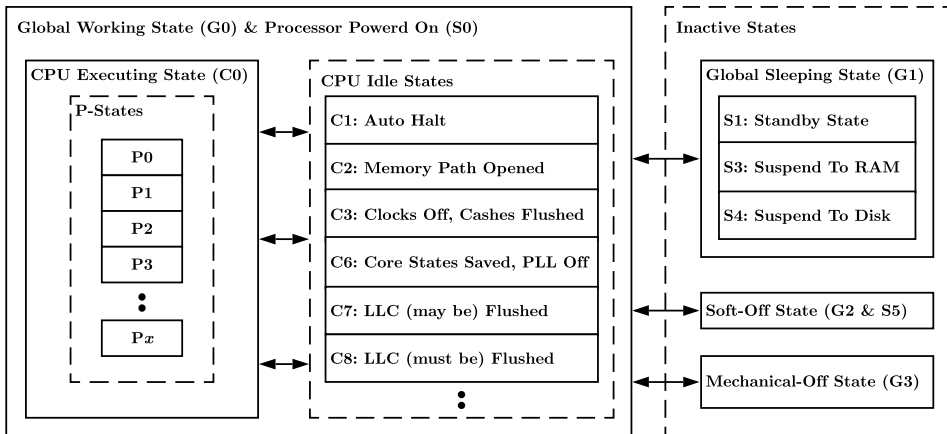


Figure 2.1: Transitions between the power states defined in the ACPI standard.

amount of energy. If so, the best option would always be race-to-halt, *i.e.*, execute all tasks (threads in the view of the kernel) as fast as possible and then, put CPUs into sleep.

This impression is not actually the case, due to the fact that the CPU frequency and the voltage applied to the CPU move together. In other words, by the law of physics, it is impossible to acquire a higher CPU frequency without increasing the voltage because boosting the voltage level requires frequency uplift. Therefore, adjusting CPU frequency has a quadratic impact on its energy consumption. In turn, race-to-halt is not ideal since the P^{idle} only has a linear effect that is not able to offset the impact of voltage variation. Hence, if possible, the CPUs should be reduced to a lower speed with a lower voltage in order to save power, which is where dynamic frequency and voltage scaling (DVFS in short) comes into play. Furthermore, the purpose of using DVFS is twofold [78]: besides saving power, DVFS helps processors reduce the peak thermal load. The effectiveness of the cooling system relies upon the peak power instead of the average power. Thus, capping the peak power decreases the cost and the size of the cooling equipment.

2.3.2 Power States

To achieve DVFS, Operating Performance Points were introduced, representing several levels of voltage and clock frequency at which the CPUs are able to operate. P-states (performance states) is the terminology used for Operating Performance Points in the Advanced Configuration and Power Interface (ACPI) standards [1]. As shown in Figure 2.2, the higher the P-state, the lower the core frequencies and voltage, in turn, ultimately saving more energy. Besides P-states, ACPI also defines a set of other types of states, namely, G-states for global system states, D-states for device power management, C-states for the CPU power states, and S-states that entail a number of sleep states. In respect of the CPU, P-states and C-states (CPU-states) are of particular interest. Figure 2.1 further illustrates how P-states and C-states work with other power states. In the context of G0 – Global Working State, it refers to the active mode in which the CPU is executing instructions, whilst from C1 to C8, the power usage of the CPU is gradually reduced in order to sequentially save more energy. In the context of C0 – CPU Executing State, the processor can operate at different P-states to further curb the energy consumption. As explained above, P-states modulate both the CPU frequency (in MHz) and the voltage at the same time, and the P-state P_x depends on the number of frequency steps available in different platforms. By the virtue of the quadratic relationship demonstrated by Equation 2.9, noticeable energy saving can be achieved by enabling P-state scaling. In addition, in some machines of older generations that predate C-States and P-States, throttling states (T-States) are used under thermal emergency where the processor is overheating. This is achieved by gating the CPU clock — the higher the temperature, the higher the T-State and, in turn, more CPU cycles will be gated. As T-States are outdated technology, they are not considered in this study.

2.3.3 DVFS in Linux

P-states are effective in saving power, they, however, are handled differently by different platforms. In this work, we focus on the DVFS implementation of Linux. The `CPUFreq` subsystem is responsible for CPU frequency scaling in the Linux kernel, offering a basic infrastructure and user-mode interface for all devices that

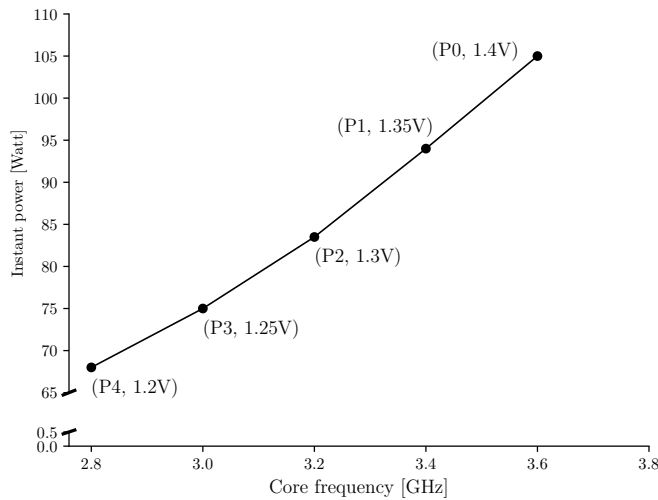


Figure 2.2: P-states and corresponding power consumption levels (data source: [75]).

support P-states. It not only provides other components with a framework in which they operate but also gives the opportunity to implement various frequency-scaling mechanisms in accordance with the (estimation of) CPU capacity demanded by workloads. To this end, a number of generic scaling **Governors** and **Drivers** are provided by the Linux kernel [91]. A **Governor** is a piece of software in which the algorithms/policies for adjusting the CPU frequencies are implemented. The scaling rules thereof are based on the estimation of the required CPU capacity. Each of the **Governors** implements one set of frequency scaling rules, and these policies are located under the directory `/sys/devices/system/cpu/cpufreq/`. The scaling **Governors** are independent of specific CPU architectures. A **Driver** is another piece of software that is responsible to interact with hardware directly. They offer available **Governors** a set of machine-specific P-states and apply the frequencies proposed by **Governors** to the machine via hardware-dependent interfaces. If the scaling algorithm implemented by a **Governor** is per-policy as opposed to system-wide/global, **Drivers** will find the corresponding tunable attributes (`sysfs`) in the subdirectory of the policies (`/sys/devices/system/cpu/cpufreq/policy{#}`). In spite of the existence of driver-specific properties, **Drivers** and **Governors** are designed to be orthogonal, *i.e.*, they are supposed to be used in any combinations.

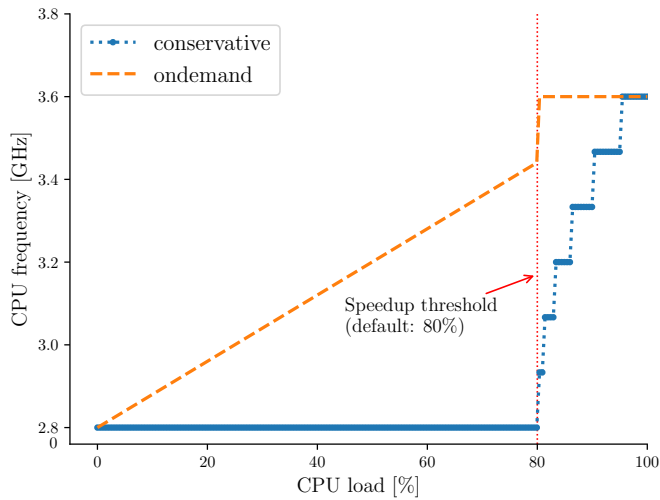


Figure 2.3: Different behaviours of the `ondemand` and the `conservative` governors in the Linux kernel.

This design is achieved via a set of `struct cpufreq_policy` objects, each of which is associated with one or several CPUs. The type of Governors can be altered during runtime, and in turn, several Governors attached to the CPUs can share the same policy object by setting the `scaling_governor` attribute in `sysfs`.

2.3.4 Governors & Drivers

Six generic Governors are available in the Linux kernel: 1) the `performance` governor, which immediately request the highest frequency within the limit specified by the `scaling_max_freq` attribute of each policy, 2) the `powersave` governor, which proposes the lowest frequency above the threshold specified by the `scaling_min_freq` of each policy, 3) the `userspace` governor, which does nothing aside from allowing the `scaling_setspeed` attribute of each policy to be set in user mode, 4) the `schedutil` governor, which runs in the context of the Linux scheduler and uses the scheduler for estimating the next frequency to which the CPUs ought to be adjusted, 5) the `ondemand` governor, which uses the CPU load

as the metric for selecting the CPU frequency, and 6) the `conservative` governor, whose policy resembles the two-stage frequency scaling process of the `ondemand` governor but requests changes in small steps. The `schedutil` governor is developed to tackle the challenge of estimating the requested CPU capacity. It employs data obtained from the scheduler instead of only the data from the CPUs since the scheduler is more informed by the run queue, information of I/O blocking, etc. Such a mechanism can be of great help in DVFS because, for instance, a task that is not running but waiting or blocked by I/O also contribute to the load of the system, whereas a long-running task that accumulatively consumes a large number of resources may not be as demanding at the moment. In addition, the difference between the `ondemand` governor and the `conservative` governor are subtle but important. Both of them are running in the process context asynchronously, which causes little overhead to the scheduler but generates more context switches. The interrupts triggered by them for updating the P-states can be irregular, and the idle time of the CPU is, thereby, reduced. As illustrated in Figure 2.3, before reaching the (configurable) speedup threshold, the `ondemand` governor will propose the next frequency proportional to the current CPU load ¹ (Equation 2.10). Conversely, as for the `conservative` governor, no changes in frequency will be requested at the first stage ².

$$f = f_{\min} + l \cdot \frac{f_{\max} - f_{\min}}{100}, \quad (2.10)$$

where f is the next frequency to propose, f_{\min} and f_{\max} denote the maximum and minimum frequency specified in the scaling policy, respectively.

Once the threshold is met, the `ondemand` governor will jump straight to the maximum frequency limit (`scaling_max_freq`), whilst the `conservative` governor will request frequency changes continuously (both increase and decrease) in small steps in order to avoid significant frequency fluctuations over short periods of time. This mechanism is particularly useful when drastic changes in CPU frequen-

¹https://github.com/torvalds/linux/blob/master/drivers/cpufreq/cpufreq_ondemand.c

²https://github.com/torvalds/linux/blob/master/drivers/cpufreq/cpufreq_conservative.c

cies are not supported or suited for the machine. The default threshold in the Linux kernel is 80%. In other words, if the idle time to wall time ratio is less than 20%, the two governors will start boosting core frequencies. In addition, the minimum step size of the conservative governor is 5% of the maximum frequency limit.

Note that the P-state scaling provided by Intel (`intel_pstate` [90]) is rather different from that of the generic policies. It comes with its own algorithms and bypasses the built-in drivers of the Linux kernel. Specifically, the Intel SpeedStep[®] and the Speed Shift[®] technologies [74] are available by the time of this study. The former switches P-states based on certain algorithms that are not open-sourced, and the latter is an improved version of the former. Instead of changing P-states in a discrete manner, Speed Shift[®] enables a full multiplier range or narrow window. It is able to fully ramp up a core speed in response to a lower P-state faster (30-35 ms) compared to that of the SpeedStep[®] (100-150 ms). However, it is limited to the Skylake architecture and needs support from operating systems. Thus, `intel_pstate` is not considered in this work.

2.3.5 Multicore

As elaborated above, although a `struct cpufreq_policy` object can be assigned to multiple CPUs, a single CPU is able to occupy a policy object itself. However, in old-generation machines, all (logical) cores of a package are managed under the same power domain [150], which means they are all in the same P-states at any time based upon their maximum load [74]. This is akin to the mechanism of the multicore C-states. Furthermore, under the OS C-states, there are 1) CC-states, which offer a set of idle states for each physical core, and 2) PC-states for the idles states at the package level covering various shared resources. Both the CC-states and the PC-states are set by taking the minimum level of their respective state (which has the highest frequency value) over all its components. Notably, in Intel x86 processors that support hyper-threading [116], per-thread C-states can be obtained, managed in the same manner as that of the PC-states and CC-states. For P-states, however, only until recent generations (after Haswell [65]), processors have just started supporting P-states for each (logical) core, known as multicore-aware P-state coordination [74]. Nevertheless, few datacenter simulators support this feature (§2.5),

as it is too detailed for high-level datacenter simulation and may introduce substantial overhead as well [25]. Also, these datacenter simulators or their respective energy extensions [68, 42, 26] reuse the core-level infrastructure previously built to achieve this feature so that backward compatibility is better maintained. As we do not have such an issue in our OpenDC simulator, we take package-level decisions when switching between P-states to circumvent the unnecessary overhead (§4.1).

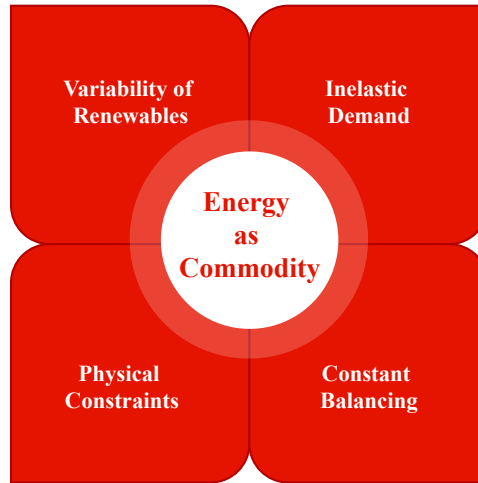


Figure 2.4: Challenges faced by energy transmission and trading.

2.4 Power Grid

Energy, or specifically, electricity, is a special type of commodity. This section covers the basics of the power grid (in the EU) with regard to both its financial perspectives and real operations. As of the time of this work, once it has been generated, it is neither practical nor economical to store them on a large scale for a long time. To function well and to avoid issues such as blackout and power outage, a balance between power generation and consumption must be kept at all times. Its transportation and distribution are performed on a power network, governed by specific physical rules of mother nature. Also, it features inelastic demand as the power load can affect the energy prices but not the other way around. The majority of the end-users is of rubric of the society (*e.g.*, residential, production, hospital, etc.), whilst the roots of the generated energy are not differentiable, *i.e.*, the power produced by burning coal and the power produced by solar panels have no difference in the eyes of the consumer. Nevertheless, the origin of energy makes a huge difference in the producer side as well as to the energy market. Determined to tackle the climate crisis, the EU has turned its market into a massive entry of renewable energy sources. However, this comes at costs that greatly diminish the competence of renewable energy compared to fossil fuel.

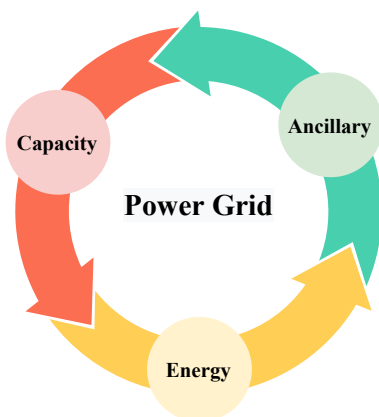


Figure 2.5: Markets around the power grids.

2.4.1 Markets Around the Power Grid

As shown in Figure 2.5, when it comes to the power grid, there are three types of markets in the EU. In the capacity market, system operators ensure sufficient capacity of power generation is retained for sustaining competitive prices and reliable operation in the coming years. They also provide many services such as Primary/Secondary/Tertiary reserves and voltage control, of which the ancillary market is comprised. In this work, we focus on the energy market in which optimal scheduling and power exchanges take place. In respect of the energy market, four major parts arranged in a sequence interact with one and another (Figure 2.6), keeping a delicate balance for both the financial markets and the actual operations in the power grid. Firstly, contracts for physical delivery of energy, price hedging and risk management are made in the Forward & Futures markets. Secondly, treating the current trending as the spinal core for predicting the matching of everyday demand and supply, participants buy or/and sell energy for the next 24 hours in a closed auction held in the Day-Ahead market, which is the most important market in the EU. In other words, according to the previous market-clearing outcomes (price and volumes for each market time unit), market operators will dispatch themselves concerning the pattern of power generation and consumption of the next day (often in a 15-minute increment).

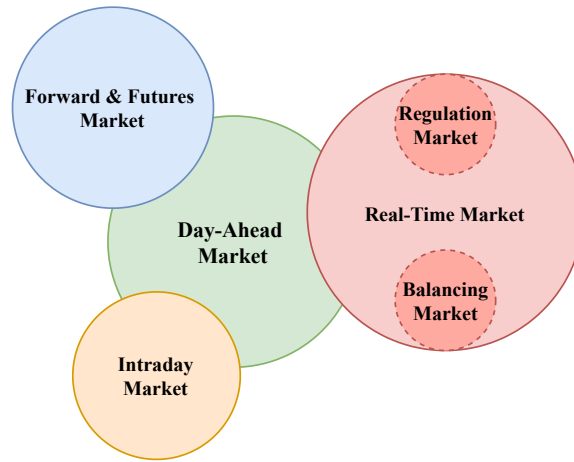


Figure 2.6: Types of energy markets in the EU.

After this blind auction, a spot price is settled at the intersection where the demand meets the supply. This can be regarded as the first balance reached in terms of the financial market of the power grid, it, however, does not imply any obligation towards prosumers, *i.e.*, no one is forced to produce or consume the promised quantities. Then, as a continuation, the Intraday market provides prosumers with a bilateral trading platform by which they can adjust their self-dispatched units in the Day-Ahead market based on the newest status (*e.g.*, updated weather/market predictions) before the actual operation/delivery commence. Finally, the Real-Time market serves as the final guard during physical operations, which is where the system operators (in the EU) take over the market, offering regulation to counteract the remaining imbalance and charging the prosumers based upon actual figures in the power meters against the contracted volume.

2.4.2 Balancing the Power Grid

There are two levels of (im)balance in the power grid, the positive/negative/no imbalance at the overall system level and that of the prosumer level. Even if not all prosumers have reached a balance, it is possible that the overall system of the power grid can still strike a delicate balance. Although balancing the grid is the central

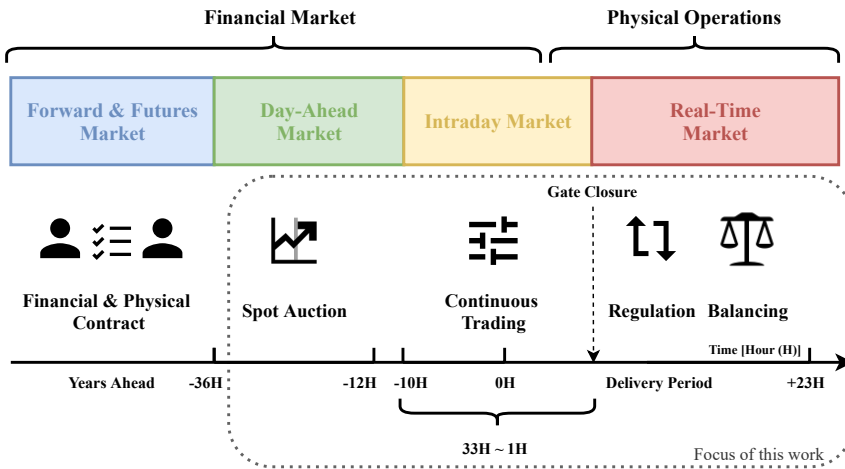


Figure 2.7: Sequence of energy markets.

duty of the system operators, it is neither solely pertinent to the system operators nor only take place in the Real-Time market. Instead, it is relevant to every participant and is performed at almost every stage of the market sequence demonstrated in Figure 2.7. Firstly, as described in the previous section, in the EU, the day-ahead auction is the most important market whose clearing yields a spot price at which demand and supply meet for the first time. This, however, does not imply any actual obligation towards prosumers. In other words, they are free to (intentionally or accidentally) break the contract during real-time operations.

As shown in Figure 2.7, there is still a maximum 33-hour interval in time between the clearing of the spot auction and the start of the delivery. Market operators can take advantage of this period to adjust their initial commitments in the Intraday market until one hour before the delivery phase, hedging the risks of under-/over-production. This can be treated as the second balancing procedure of the power grid. Last but certainly not least, after the physical operations commence, the system operators take control to ensure the grid is balanced during energy transmission. The Real-Time market consists of two sub-markets, namely, the obligation market and the balancing market. Prior to the start of transmission operations, prosumers who participated in the regulation market can offer to buy and/or sell regulation power to the system operators. In other words, these are the participants

who are capable of helping balance the grid, declaring to the system operators the willingness of altering their set-points on the consumption side and/or the production quantity. Also, the system operators can actually purchase the regulation resources from neighbouring countries. This is usually in the form of commercial entities buying, selling and transporting electricity across a border, and participating in markets in both countries as any other energy prosumer would. BritNed, the cable between the United Kingdom and the Netherlands, is an example of this. In contrast, the prosumers who do not respect their original agreements and, in turn, induce imbalance in the power grid will be charged by the system operators on the basis of the difference between their original schedules and the actual readings from the power meter. On the contrary, if a participant is helping the system operator balance the power grid, *e.g.*, a consumer is in positive imbalance and the system is in positive imbalance as well, then the prosumer would be rewarded by the system operator as they are helping balance the grid (in spite of not following their schedule). Via these settlements and other energy reserves, the balance of the power grid is maintained in real-time operations by the system operators.

Note that the elaboration above is tailored to the EU energy market in which the transmission system operator (TSO) has the ultimate responsibility to keep its transmission system in balance. For other places such as the USA that relies on the independent system operator (ISO), the energy market, especially, the balancing mechanisms are different. Nevertheless, the market sequence is shown in Figure 2.7 still applies and the objectives of the balancing phase stay the same. An excellent load/price forecast and flexible operational responses to the market will bring substantial benefits to the prosumer. For example, if it is predicted that the demand/price will be lower during the Real-Time market, then one can sell (short) in the Day-Ahead market and buy in more energy in the Real-Time market. This can facilitate energy-aware scheduling in industry production. As a result, the prosumer will receive the profits along the lines of Equation 2.11:

$$G = (Q_a - Q_s) \cdot \zeta, \quad (2.11)$$

where G denotes the financial gain, Q_a represents the actual quantity of energy transmission, Q_s is the scheduled quantity in the spot market, and ζ denotes the

locational marginal pricing.

Bidding in the Day-Ahead market is going to push the prices up; good forecast and flexible responses to the market foster smooth pricing in the Real-Time market, which will bring the prices down to what they ought to be. Consequently, it can bring the prices of today and tomorrow closer, and ultimately help the power grid strike a balance that would greatly reduce the electrical system charges and, in turn, lower the grid balance cost.

Network Code	Definition	Activation Time
FCR	Primary control (automatic activation)	> 15 minutes
aFRR	Secondary control (automatic activation)	30 seconds – 15 minutes
mFRR	Tertiary control (semi-automatic or manual activation)	≥ 15 minutes
RR	Optional control (semi-automatic or manual activation)	≥ 15 minutes

Table 2.1: Balancing reserves in the EU.

Furthermore, the imbalance energy trading happens in the regulation & balancing markets are of particular interest of this work. During the trading, Balancing Service Providers (BSPs) offer bids to counter the imbalance in power supply and consumption during the delivery hour. The imbalance is introduced by Balance Responsible Parties (BRPs), which are the prosumers responsible for the deviation between their actual delivery and their self-dispatched volume concluded in the day-ahead/intraday clearings. Specifically, in the Netherlands, these self-dispatched commissions are in the form of “E-programs” in a 15-minute resolution. Having been reported the readings of the meters, TSO will ascertain the deviation for every BRP. Then, BRPs are obliged to pay their energy deficit or surplus by participating in the imbalance settlement. In terms of timing, BRPs can change their “E-programs” until one hour before the start of the delivery through, *e.g.*, the Intraday market, whilst BSPs are able to change their bids 30 minutes before the start of a delivery hour. The 15-minute interval of an “E-programs” is known as an Imbalance Settlement Period (ISP), and every day has 96 ISPs. As summarized in Table 2.1, various mechanisms are employed by TSO in the imbalance system, including Frequency Containment Reserve (FCR), automatic Frequency Restoration Reserve (aFRR), manual Frequency Restoration Reserve (mFRR), and Replacement Reserve (RR). They differ in their rate of ramping up. For example, the activation

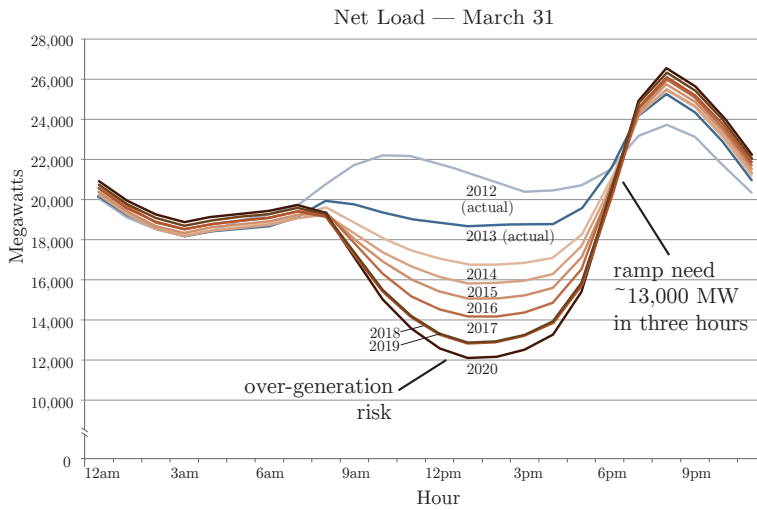


Figure 2.8: Duck curve showing the net load of the power grid for 11 January from 2012 to 2020 in California (figure source: [81]).

time of an aFRR is required to be no longer than 15 minutes with the minimum bid size of 1 MW. In this work, we do not distinguish the types of means by which the BSPs provide balancing energy to the grid since they are paid from the same ground irrespective of the type of reserves activated.

2.4.3 Challenges of Renewables

As elaborated in the previous section, balancing the power grid is a win-win for both the grid system and the prosumers. However, it faces increasing challenges brought by the introduction of renewable energy sources, for example, solar, wind and tidal power (its current contribution is diminutive to the grid). This is due to the fact that these origins of energy are erratic and, thus, hard to predict. This brings risks to both the prosumers and our environment. For prosumers, the energy market incurs more fluctuation and, in turn, is more precarious that even the Intraday market will not give enough operational headroom. Concerning the environment, renewable energy generation needs constant backup from traditional energy sources, which exacerbate the reliance on fossil fuels. One of many quintessential exam-

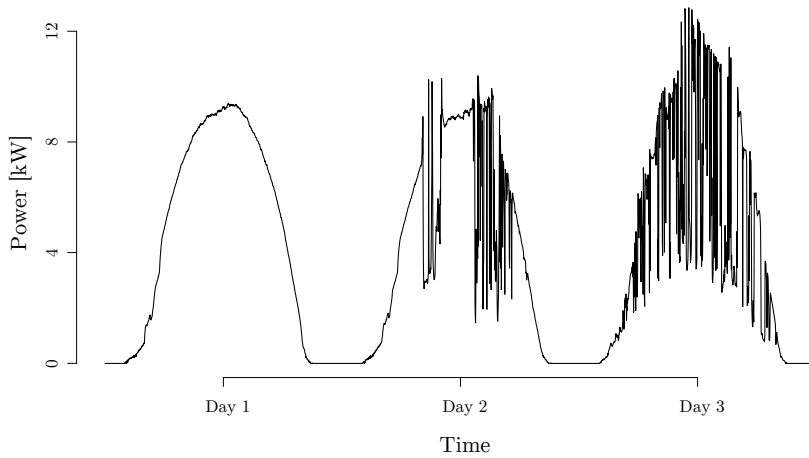


Figure 2.9: Solar power generation at the University of Queensland over three days in September 2018 (data source: [81]).

ples for variable generation resources is the duck-curve phenomenon (Figure 2.8) first introduced by California independent system operator (ISO). As the degree to which the power grid is dependent on renewable sources grows from 2012 to 2020, the daily variation of the load within increasingly boosts over the years.

To further illustrate the challenges brought by renewables, Figure 2.9 shows the power generated by the solar panels at the University of Queensland, St Lucia campus, Prentice building in three consecutive days. When there is no sunshine at all, a drastic drop will occur in solar power generation; if the weather is overcast, the reflection of light on the clouds may actually bring more solar power. The third day in this graph is obviously rather cloudy. The challenges brought by renewable energy sources in balancing the grid put further emphasis on the paramount importance of flexible and quick demand-side responses to the requests from the system operators and price/load of the energy market. This will not only enable prosumers to stay competent and profitable but also provide the power grid with sufficient energy reserves in the long run.

2.4.4 Demand Side Management

Intensive work has been conducted to explore the datacenter's participation in DSM programmes. A synergistic control strategy and a flexibility factor have been developed to increase the capacity of frequency regulation in a recent study [52]. Fu et al. presented a real-time, multi-market optimization framework for the datacenters [53], facilitating datacenters to participate in both the energy market and the regulation operations. Taking into account the energy cost, demand costs and regulation revenues, the framework optimizes the bid in each hour, helping datacenters meet energy and demand goals whilst retain minimum cost through obtaining maximum regulation revenues. However, this work focuses on enlarging the FR capacity of datacenters, which is of rarity in the industry [12].

Via dynamic pricing, Li et al. proposed a collaborative framework for optimizing the overall costs of several datacenters in their position paper [102]. The framework employs collaborative efforts across multiple geographically distributed datacenters that communicate via dedicated network fabrics to negotiate mutually optimal energy prices. In this collaborative system, an optimization platform [107] is used to enable constraint optimization problems (COPs). However, to benefit from these schemes, multiple datacenters and their utility companies have to be involved in order to optimize the price updates in the energy market.

In respect of market design, two categories of programmes for DR are commonly available. The first category entails bidding/supplying a certain amount of demand flexibility into the market. In other words, consumers bid their flexibility via supply functions that are parameterized (*e.g.*, [85, 176]). As for the second category, consumers buy in or respond to published prices that were selected according to predictions on (potentially) available flexibility; examples of the second sort include [38, 101, 136, 109]. Specifically, programmes such as the Coincident Peak Pricing (CPP) [7] and price-based incentive programmes [145] are available for datacenters to participate in. These programmes charge much higher prices for electricity during peak hours (usually over 200× higher than the base price [109]). Such peak-hour costs can account for 23+% of the consumers' energy bill [158], which is a strong motive for consumers to reduce their energy usage during peak hours. However, when it comes to datacenters, Liu et al. argued that the market

designs and models in the second category outperform those in the first. Also, the current available DR programmes, such as CPP, are not particularly suitable for datacenters since they can barely fully extract the flexibility of datacenters due to many reasons. Chief amongst them is that although datacenters may incur charges if no response is offered, no profit would be produced either if the peak in energy load does not happen during the coincident periods. As a result, datacenters nowadays provide little, if any, response to market signals [58, 60, 109, 108]. In addition, a large amount of effort has been devoted to the optimization of workload management of datacenters (*e.g.*, [34, 56, 71, 105, 121, 175, 179, 183]), especially, via workload distribution [138, 110, 57, 170, 182, 168]. When such optimization is considered, time for shedding the load is able to be reduced and, in turn, the flexibility of datacenters can be exploited further. Moreover, by leveraging the regional difference of the energy cost, datacenters can leverage their distributed nature to further optimized [135, 138, 110, 106, 57, 170, 182]. In fact, these management strategies require cooperation between the power grid and geographically distributed datacenters together with their utility companies. Thus, they are far more complex to orchestrate than what is introduced in this work, which, in turn, could potentially hinder their adoption. In this work, we develop a straightforward, short-term scheme, whereby individual datacenters can participate in the energy market, both saving their energy cost and curbing their energy consumption, whilst providing the power grid with indirect DR.

2.4.5 Datacenters and the Power Grid

Whilst system-level challenges have been brought into the grid by renewable energy sources, increasing numbers of opportunities are being created in the meantime. Especially, as the proportion of worldwide electricity attributed to datacenters is sky-rocketing, these challenges could well be a blessing in disguise for datacenters. Over recent years, more and more studies regarding the participation of datacenters in the power grid have been carried out. Several identified characteristics therein of datacenters underpin the increasingly important role played by datacenters. To name a few, firstly, the capacity of datacenters are huge in the sense that the nameplate load of a single datacenter is able to reach as large as 50+ MW [19]. Moreover, for large cloud providers, the critical power of a single datacenter can

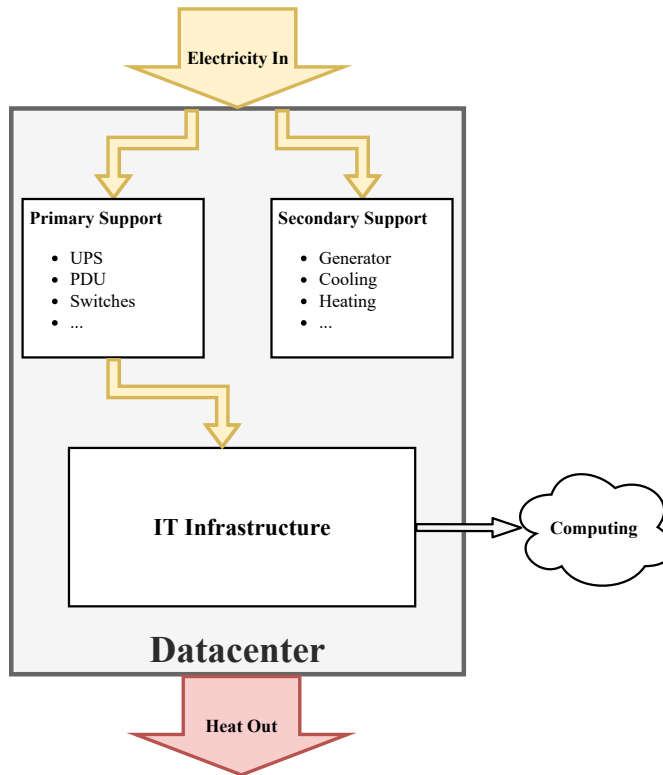


Figure 2.10: Energy flow in datacenters.

even exceed 100 MW [9]. To put this into perspective, Wierman et al. suggest that a well-managed datacenter of 30 MW has approximately the same capability for regulating the power grid as massive energy storage of 7 MWh. Thus, overlooking datacenters in balancing the grid means missing out on a huge amount of capacity. Secondly, datacenters are built for extremely reliable and available services. For instance, two availability classes to which most modern datacenters belong are Tier III and Tier IV. Datacenters that fall into the first category insure a 99.982% availability, and for the latter, the figure is 99.995% [114]. To accomplish promised uptime and performance guarantees, a considerable amount of redundancy is introduced in the power system of datacenters.

As demonstrated in Figure 2.10, a large portion of the energy goes into the power support system that has a high degree of redundancy, and almost all energy

ultimately turns into heat in the end. To ensure a secure power supply, datacenters usually have two energy sources. As accessing to two utility power sources is not often feasible, most datacenters use a backup generator as the second power source. These generators are powered by gas, diesel or flywheel, providing power when the utility power fails. In the event of power failure, the automatic transfer switch (ATS) will start the generator to power the uninterruptible power supply (UPS) load. The UPS system supports servers, data communication systems and other equipment during a sudden power failure or voltage drop. It provides clean power to sensitive data equipment by eliminating power surges, noise, spikes, etc. The UPS system also constantly conditioning and monitoring utility power to protect the load. Batteries therein are constantly been charged for emergency support during a utility outage. Note that one of the most important factors to ensure proper UPS performance is the battery quality in the sense that one bad battery is able to bring the entire system down during a power interruption. Serving as the last layer of adjusting the electricity, Power distribution units (PDUs) provide the ability to control and monitor how the power is distributed to the IT infrastructure. The facilities in the mechanical yard are critical for environmental controls (*e.g.*, heating, cooling and humidity). They maintain a proper environment for electronic equipment by tolerating fluctuations in moisture and temperature. One of the most efficient forms of cooling for datacenters is a close-coupled in-row, chilled water cooling system, also known as a computer room air conditioner (CRAC). Furthermore, the appropriate placement of air returns and the use of perforated floor tiles/sensors can help eliminate hot spots and gain efficiencies in the building. All these equipment are built for facilitating the computing services provided by the server farm, in which racks designed specifically for datacenters offer a modern enclosure with strength and stability for any server environment. Components such as the backup generators, UPS, transformers, chillers/Computer Room Air Conditioner units (CRACs), Computer Room Air Handler units (CRAHs), etc., can all be regarded as redundant equipment.

Further, perhaps what counts more is the significant flexibility of datacenters, making them extremely elastic power loads for the grid. For example, the wide range of temperatures under which the datacenters can operate results in various power loads [50]. Also, many workloads of modern datacenters are delay-tolerant. In other words, the schedule of these workloads can be shifted in response to the en-

ergy market and requests from the system operators to optimize profit as elaborated in previous sections. Power management techniques such as frequency scaling [33], power capping [183, 29, 105] and different levels of energy-saving configurations [71] further boost the flexibility of datacenters. Lastly, datacenters are complex but highly automated systems. Ubiquitous monitors and controls empower datacenters' participation in the power grid.

In general, datacenters participate in the power grid in two ways, Demand Response (DR) and Frequency Regulation (FR) [12]. Note that FR can be categorized as a special means of DR [171] as the overlap between the two is substantial. LBNL proposes potential DR resources that reside at different components of datacenters [114], including supporting equipment such as backup generators and UPS, programmable power managements such as DVFS and power capping, server consolidation by virtualization, and load (re)scheduling & migration. Note that some of these methods are not particularly environment-friendly, for example, the using the backup generators [50]. But in terms of effectiveness, even methods like power capping are relatively coarse-grained, it can potentially enable datacenters to tuck about 25% more servers into the same amount of space [146]. In respect of FR, its resources have a huge overlap with that of the DR. The main difference between the two is that FR requires a way faster timescale, usually at the second level. Moreover, both the supply side and the demand side need to constantly work together to automatically lubricate the small frequency fluctuations in the power grid. By the virtue of fast charging and discharging, UPS becomes a good candidate in FR; the "UPS-as-a-Reserve" [46] is one pilot project of such. Similarly, because of their quick response, power management such as DVFS [28, 30] and dummy workload [171] are also well-suited resources for FR.

FR enables datacenters to take part in the real-time market as well as the ancillary market. Nevertheless, DR is far more common than FR because, for datacenters, FR is too fast and risky to visualize and control [12]. In this work, we focus on datacenters' participation in the day-ahead and the intraday markets. Thus, DR is the focal point of this study.

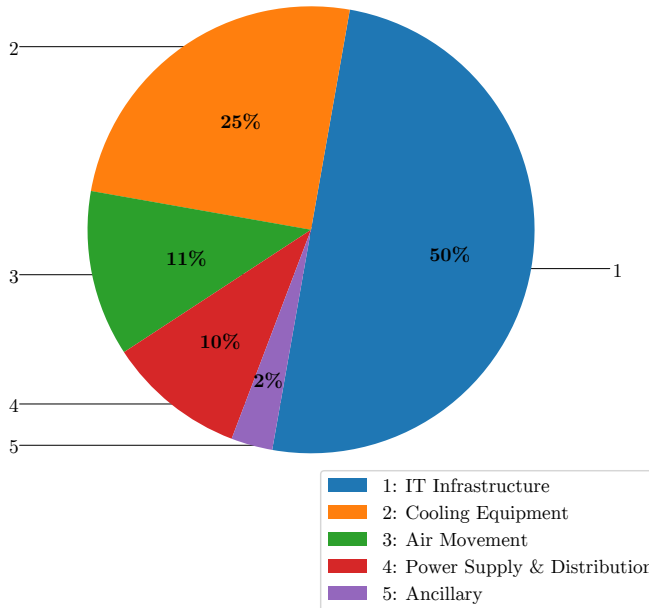


Figure 2.11: Approximate distribution of energy usage in a datacenter with PUE value of 2.0 (data source: [11]).

2.5 Energy Modelling for Datacenters

Methods for modelling energy consumption in datacenters can be broadly classified into two categories: measuring energy usage at the hardware level [27, 51, 39] and modelling energy consumption using simulation [13, 18, 36, 63, 140, 163, 167, 181]. Hardware measurement has a huge advantage over the whole-system simulators in terms of speed, so much so that the latter can hardly be used for long-term applications and very large dataset without applying reduction configurations [47]. However, online energy monitoring and metrics collection systems are of rarity incurring substantial costs in practice. To our knowledge, the work from Fan et al. [48] is the first to use theoretical energy models for very large-scale datacenter power provisioning on live, production workloads.

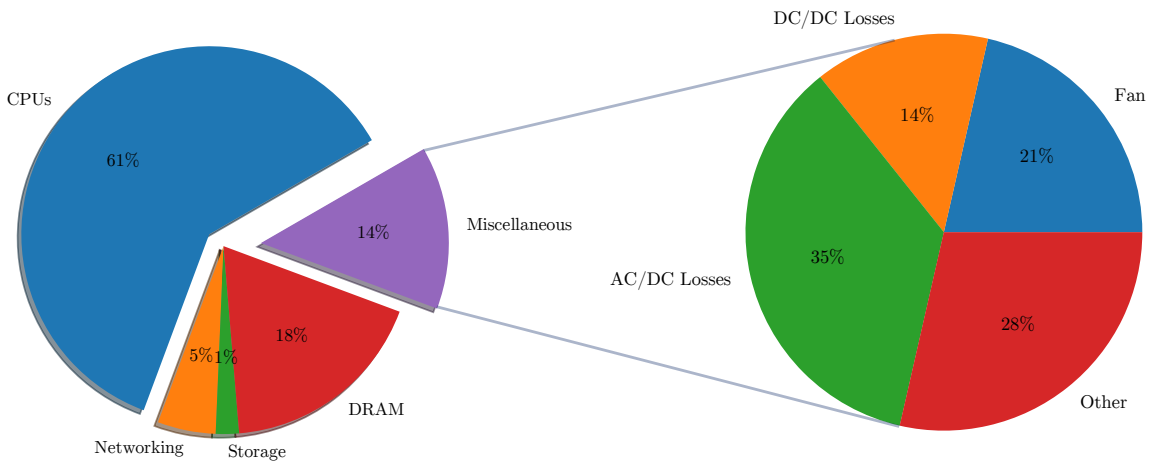


Figure 2.12: Approximate distribution of energy losses in IT equipments (data source: [11, 75]).

Similar to the approach from Economou et al. [47], instead of considering the full system with fine-grained models, Fan et al. take a helicopter view using metrics such as CPU utilization and I/O activity at a coarser granularity to estimate energy consumption. The authors focus on critical power without taking into account energy losses and cooling power consumption at the datacenter level. Economou et al. [47] suggest that the power consumption of non-IT infrastructure is by no means negligible because they amount to 30-50% of the total energy consumption. Figure 2.11 gives an approximate distribution of energy usage in a datacenter with a PUE value of 2.0 [11], showing that about 50% of the total energy consumption is used by IT infrastructure. Figure 2.12 offers an overview of the rough proportions of the energy losses in terms of various IT components operating at peak power [75, 9].

Although non-IT components account for around half of the total power usage, Fan et al. [48] argue that in modern datacenters, critical power can accurately capture the energy consumption for other non-IT facilities in the sense that the dynamic power of the non-IT components can be modelled as a static tax proportional to the critical power in modern datacenters. This estimation can be further facilitated by proper calibrations. The authors proposed a power model (Equation 2.12a and 2.12b) where u is the CPU utilization that is obtained via the operating systems, which averages across all CPUs, and r , which was set to 1.4 in the original

experience, is the calibration parameter chosen to optimize the mean square error (MSE). This model, albeit simple, has been widely adopted in simulating energy consumption in datacenters.

In this model, the CPU utilization is employed as the single indicator for estimating the power consumption of the critical load. Fan et al. demonstrated that CPU utilization is able to serve as an extremely accurate signal providing veracious results for thousands of machines. Consequently, measurements for additional loads, e.g., hardware performance meters, are complementary yet unnecessary. Furthermore, a sub-/super-linear correlation between power and the CPU frequency is assumed in this power model. Studies [14, 137, 98, 161, 55] have illustrated that validity of this assumption lies in the fact that DVFS is only applied to the CPU, not to other components, and the number of states defined in DVFS for the CPU is finite.

$$\begin{cases} \mathbb{P}(u) &= P^{\text{idle}} + (P^{\text{max}} - P^{\text{idle}})u & (2.12a) \\ \mathbb{P}_{\text{MSE}}(u) &= P^{\text{idle}} + (P^{\text{max}} - P^{\text{idle}})(2u - u^r) & (2.12b) \end{cases}$$

In addition to this, whilst advocating the usage of actual peak power as opposed to the nameplate power (because the latter is so conservative that 80-130% more machines can be deployed in the case that the nameplate power is targeted), Fan et al. [48] emphasize the importance of energy management overall power range, e.g., DVFS. To support such optimization, the authors set a predefined threshold for the CPU utilization and half the power provisioning of the CPU (whilst leaving others unchanged) whenever the utilization drops below the threshold. Despite such DVFS strategy used in their simulation is rather simple, in a cautious measurement, this strategy results in about 30% reduction in peak power and is capable of saving around 23% system energy. In addition, servers are idle for a large fraction of the running time, which consumes about 50-60% of the actual peak power, whilst they are barely completely inactive, i.e., sleep or standby. Therefore, in terms of large-scale datacenters consisting of thousands of machines, different sleep states (C-states) have a limited impact on the energy consumption at the datacenter level.

Simulator	IT Infrastructure		Primary Support		Secondary Support	Energy Market Integration
	Critical Load	DVFS	UPS	PDU		
<i>DCSim</i> [62]	✓	✗	✗	✗	✗	✗
<i>CloudSim</i> [24]	✓	✓	✗	✗	✗	✗
<i>GDCSim</i> [62]	✓	✗	✗	✗	✓ ⁺	✗
<i>CloudSched</i> [154]	✓	✓	✗	✗	✗	✗
<i>DISSECT-CF</i> [115, 88, 87]	✓	✓	✗	✗	✓ ⁺	✗
<i>GreenCloud</i> [17, 156, 94]	✓ ⁺	✓ ⁺	✗	✗	✗	✗
<i>iCanCloud/E-mc²</i> [26, 124]	✓ ⁺	✓ ⁺	✗	✗	✗	✗
<i>SimGrid</i> [25, 68, 42]	✓ ⁺	✓ ⁺	✗	✗	✗	✗
<i>OpenDC</i> [80, 119]	✓	✓ ⁺	✓ ⁺	✓ ⁺	✓	✓ ⁺

Table 2.2: Overview of the nine surveyed datacenter simulators, where the ✓ symbol means that the corresponding energy model is available, the ✗ symbol means that it is unavailable, and ⁺ represents advanced support.

Furthermore, over the last decade, many datacenter simulators [154, 24, 61, 155, 115, 88, 87, 62, 17, 156, 94, 26, 124, 68, 42, 25, 97, 112, 111] have been developed. Some of them embed energy models for different components, including both IT and non-IT infrastructure. Such advancements foster the development of energy-saving algorithms and energy-aware decision-making acknowledging the benefits brought upon by simulation, such as simplicity, reproducibility, and cost-friendliness. Table 2.2 gives an overview of nice state-of-the-art datacenter simulators in which the power consumption has been captured to various degrees. As shown in the table, we are the first to achieve whole power system modelling in datacenter simulation.

Energy Modelling & Management

In this chapter, we reason and describe our design of the energy modelling and management system. We begin by introducing the development pipeline closely followed in this work (§3.1). Then, in Section 3.2, we detail the requirement engineering process carried out in this work. Finally, we present the architecture of the entire system as well as its subsystems in Section 3.3.

3.1 Development Pipeline

To achieve a coherent design and reproducible results, we rigorously follow a development pipeline shown in Figure 3.1 from the outset, adhering to the code of conduct in the AtLarge research group.

The first stage in this pipeline is requirement engineering, which will be described in detail in the next section (§3.2). Based on the requirements elicited and documented in the first stage, we then draft our design, outlining more elaborate specifications. During the implementation stage, we put up and organize ideas into different categories using Kanban. These ideas are constantly being reviewed and updated together with the initial design. Furthermore, instead of developing

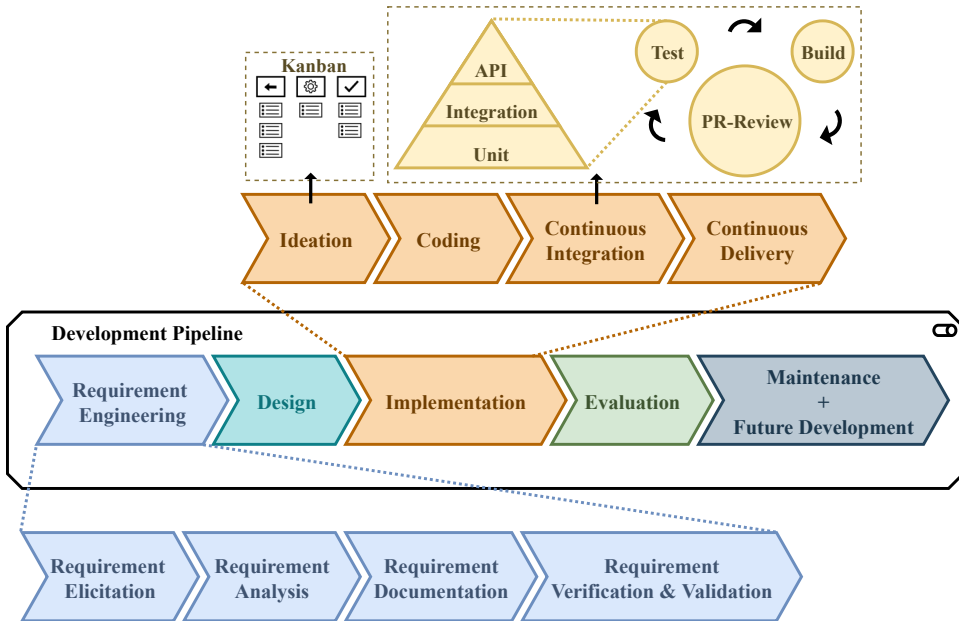


Figure 3.1: Development pipeline.

the system in a separate branch, we practice Continuous Integration and Delivery (CI/CD), which is commonly employed in modern software development. During CI/CD, we develop the system in small, measurable steps, every integration of which is reviewed by at least one senior engineer from our research group. Testing, building, and deployment are all instrumented and automated so that the system is a useable instrument at all times. Next, we conduct regular evaluations of the system, which provides frequent feedback for previous stages. At last, any changes to the codebase are carefully documented, facilitating maintenance and future development.

3.2 Requirement Engineering

In this section, we detail the first stage of the development pipeline – Requirement Engineering, which is further modulated into several steps. We first begin with

eliciting requirements from a variety of facets in the views of stakeholders using the Six Thinking Hats technique [44] (3.2.1). Next, we conduct analysis on the elicited requirements through use case modelling in Section 3.2.2. In Section 3.2.3, we then formally document the analysed requirements, categorized into function requirements (FRs) and non-function requirements (NFRs). Last, but certainly not least, we verify and validate the documented requirements in Section 3.2.4 with experts in both datacenters and the energy market.

3.2.1 Requirement Elicitation

In this section, we first identify and classify potential stakeholders, and then, employ the Six Thinking Hats approach [44] to explore system requirements.

3.2.1.1 Stakeholders.

The importance of stakeholders is paramount, we, therefore, start with identifying some potential stakeholders in Table 3.1.

Industry	Stakeholders
IT	datacenter managers, datacenter operators, datacenter technicians, cloud architects, cloud tenants
Energy	consulting firms, energy market operators, power grid system operators, renewable energy suppliers
Others	legislators, end-users of cloud services

Table 3.1: Potential stakeholders.

Having identified potential stakeholders, we classify them into two categories in Table 3.2: active stakeholders of whom the successfulness of this research is the predominant interest, and passive stakeholders who more care about whether the outcome of this work abides by their agreements and rules.

Category	Stakeholders
Active stakeholder	datacenter managers, datacenter operators, consulting firms, energy market operators, renewable energy suppliers
Passive stakeholder	datacenter technicians, cloud architects, power grid system operators, legislators, end-users of cloud services

Table 3.2: Stakeholder classification.

3.2.1.2 Change of Perspective.

Now, we employ the Six Thinking Hats method [44] seeking system requirements. With the purpose of creating new ideas, the colour sequence of hats applied in the following elicitation is: blue, white, red, green, yellow, black, and blue.

Blue Hat. The role of the blue hat is to manage and control the elicitation. Thus, we commence the process by explaining the objectives. We aim to find possible requirements from various perspectives so that the design of this work caters for the needs of identified stakeholders.

White Hat. Firstly, we summarize facts (**F**s) regarding the power grid and data-center energy management from previous sections (§1.1, §2.4).

- F1** By virtue of smart grid functions, activities from energy prosumers are able to have bidirectional influence in regulating and balancing the power grid.
- F2** The capability of the power grid is hitting some limits due to the massive introduction of renewable energy sources, which features intermittency and stochasticity caused by a range of sporadic environmental factors.
- F3** Datacenters are well suited for regulating and balancing the power grid because of their unique characteristics such as large capacity, high flexibility and redundancy, etc.
- F4** Datacenters nowadays rarely actively participate in the energy market, providing little regulation capacity such as DR to the power grid.

Red Hat. According to the perspective of the red hat, we next express emotions and feelings of the initiatives. Referring back to Chapter 1, the environmental crisis is taking its toll at a worrying pace and consequently, the calling for a lower carbon footprint is at an all-time high. Datacenters play an essential role in our day-to-day life, whilst with their ever-increasing energy consumption, they are at the front line of curbing environmental issues. Also, the operational cost resulted from energy consumption incurs heavy bills on datacenters. People from both the society and the computing industry are longing for further exploration regarding datacenters' participation in the energy market.

Green Hat. Chief amongst the spirits of the green hat is creativity, promoting innovative solutions and new ideas. With this initiative in mind, we propose the following potential solutions (**Ss**) to the aforementioned challenges.

- S1** We can model the whole datacenter power system, providing datacenters managers and operators the trending of energy consumption by means of datacenter simulation.
- S2** We can estimate the energy costs in different markets given the (predicted) workload, supporting the active participation of datacenters in the energy market.
- S3** We can take advantage of available resources such as price forecasts produced by machine learning methods, facilitating the decision-making in the energy market.
- S4** We can simulate fine-grained energy management configurations such as the DVFS technique, assisting datacenters in optimizing operational scheduling and further enabling proactive demand response.

Yellow Hat. Optimism and positivity of the proposed solutions are brought to the table by the yellow hat, which puts emphasis on the possible advantages and opportunities. Firstly, regarding **S1**, we believe knowing is power – providing the status of the datacenters' energy consumption visually and quantitatively will empower datacenter managers and operators. For **S2**, diversifying the types of markets

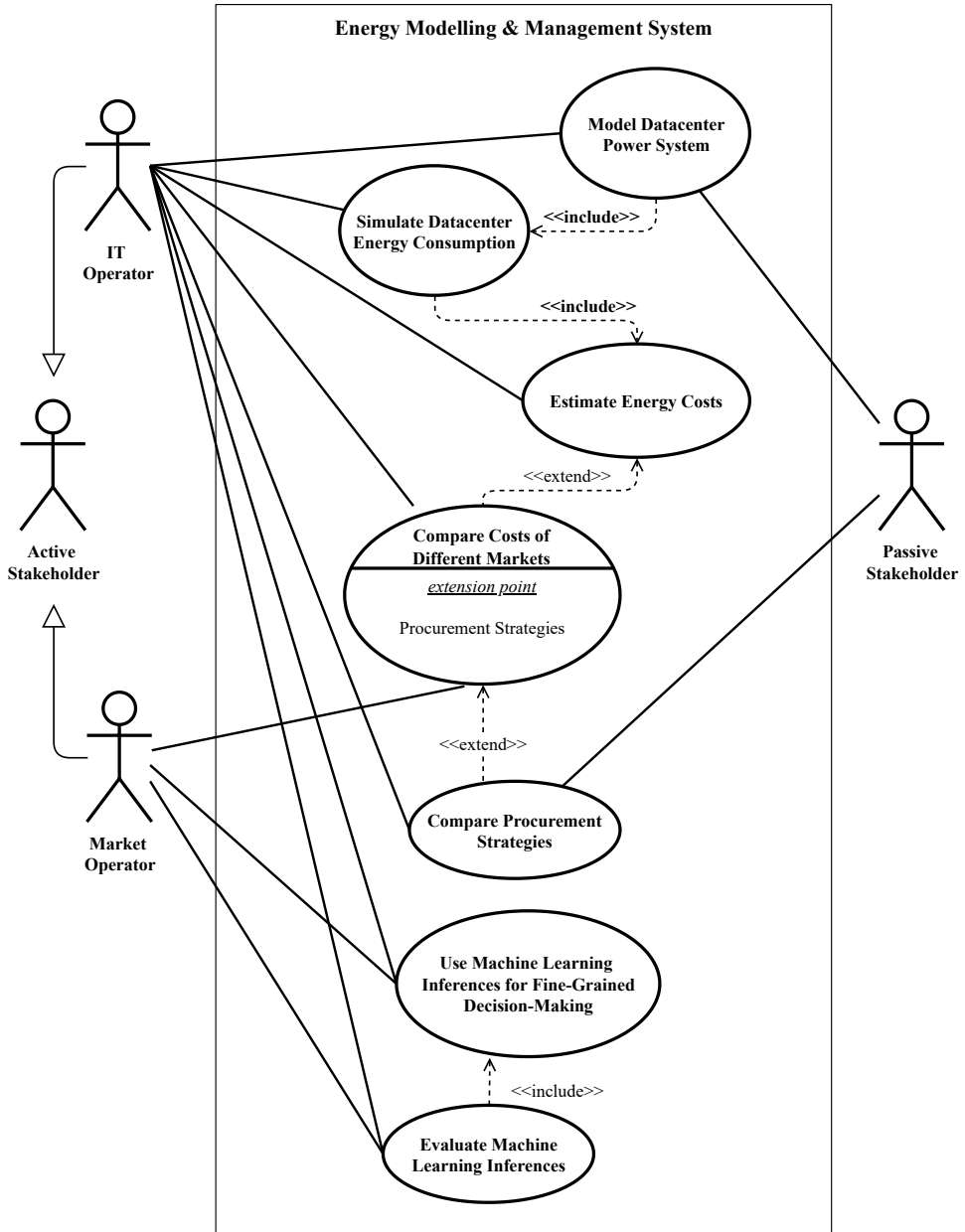
in which datacenters can participate can give market operators more options when certain constraints occur in one market, such as the available bids in the intraday market cannot meet the current energy needs. Also, **S2** can leave datacenters with more leeway under situations where, for example, the workload is not as delay-tolerant. **S3** can prompt enabling cooperation between datacenters and consulting firms, creating a level playing field in the energy market for datacenters to participate. Lastly, **S4** can further add values to the synergetic collaborations in **S3** by employing proactive, fine-grained optimizations through low-level energy configurations.

Black Hat. To provide early criticism and judgment before the requirement verification & validation stage (§3.2.4), the black hat plays devil’s advocate, bringing up potential difficulties and risks that the proposed solutions could face. With regard to **S1**, when it comes to datacenter energy modelling, there is a myriad of factors that have substantial impacts on the energy consumption of datacenters, for example, the heterogeneity of supporting equipment and machines, the various topologies, etc. Different from detailed emulations, simulations incur less overhead but can hardly capture many of those factors. Concerning **S2**, load prediction is by no means trivial and can have a great effect on the participation of datacenters in different markets. Thirdly, the performance of the machine learning methods from the consulting firm can be both beneficial, if the inference is accurate, and detrimental if it is not. Thus, it is hard for datacenter managers and operators to embrace **S3** without bounded estimations of the performance impact. In addition, fine-grained energy management techniques such as DVFS may be too low-level that could barely make a sizable influence. Further, such optimization could introduce additional overhead to the system. Lastly, experts in the energy industry and the IT industry may well have drastically different familiarity with datacenter simulation tools. Consequently, the usability of the developed tool could significantly vary amongst different user groups.

Blue Hat. To conclude, we first reiterate four major observations (**F1 – F4**) in the stage of the white hat. Then, we add on to the facts the emotional elements entailed by the perspective of the red hat. Next, via the green hat, we propose four possible solutions (**S1 – S4**) to overcome the challenges. From the perspective of

the yellow hat and that of the black hat, we reflect on the potential pros and cons of the proposed solutions respectively. As a result, we recognize the potential benefits the proposed solutions can bring to the stakeholders whilst, in the meantime, be aware of the possible side effects and risks that come along with the positive facet.

Figure 3.2: Use case diagram of the system.



3.2.2 Requirement Analysis

In this section, we take results from the previous stage, requirement elicitation, and put them into the context of the energy modelling and management system that we are trying to build. Figure 3.2 is a generalized use case diagram, showing the interactions and communications between the system and various actors. The IT operator represents the active stakeholders in the IT industry summarized in Table 3.1 and 3.2. As a primary actor, the IT operator should be able to initiate all four major use cases in the system, namely, modelling datacenter power system, simulating datacenter energy consumption, estimating energy costs, comparing costs of different energy markets, making decisions based on ML inferences, and evaluating the ML model performance. Another primary actor is the energy market player, including the active stakeholders in the energy industry (Table 3.1 and 3.2), and can initiate the last three use cases. The associations between the first three use cases are “include” since the system should only be able to simulate the power consumption given a specific power system model and to provide cost estimations based on simulation results. Similarly, the association between the last two is also “include” as the performance of the ML inferences should always be measured. In the case those procurement strategies are provided, the system should allow users to compare them based on a load prediction in order to mitigate the concern raised by the Black Hat regarding **S2**. Unlike active stakeholders, passive stakeholders, such as legislators and cloud architects, serve as the secondary actor who is primarily concerned with, for example, the validity of the power system model and the legitimacy of the energy procurement strategies.

3.2.3 Requirement Documentation

In this section, we extract the marrow of the previous analysis and formally document the requirements into functional requirements (FRs) and non-functional requirements (NFRs).

Based on the use case analysis in the previous section, we summarize **FRs** as follows:

(we use “the system” to indicate the energy modelling and management system)

- FR1** The system should enable users to model the power system of datacenters.
- FR2** The system will simulate the datacenter energy consumption given a power system model when hosting user workloads.
- FR3** The system will estimate the energy cost of datacenter operations based on the results of the workload simulation.
- FR4** The system shall demonstrate to users the potential costs in participated markets.
- FR5** The system shall provide users with the ability to compare various procurement strategies.
- FR6** The system will empower fine-grained decision-makings for users according to ML inferences.

Now, referring back to the concerns raised by the Black Hat, we further specify a list of NFRs below:

- NFR1** When modelling the power system of datacenters, the system should incorporate heterogeneous topologies, hardware components, etc.
- NFR2** When comparing various load-forecast-based procurement strategies, the system should be able to assess their impact to an extensive extent.
- NFR3** When employing inferences from machine learning methods, the system should conduct bounded evaluations on the effectiveness of using the predicted prices.
- NFR4** When developing the user interface of the tool, the system should respect inclusive design, *e.g.*, experts from the energy and the IT industry should be able to run the system with less than four steps/click and with minimum prerequisite knowledge.

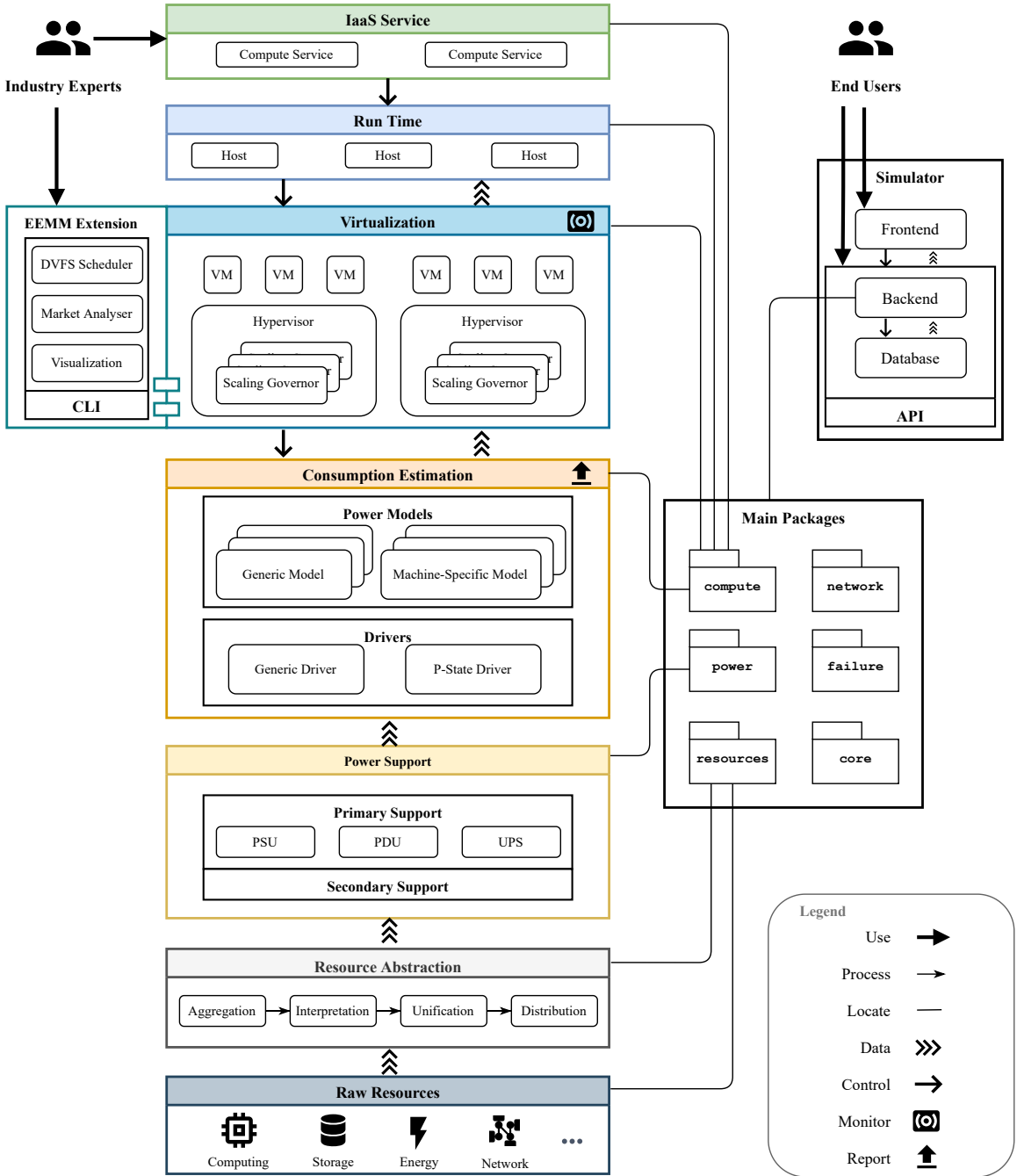
3.2.4 Requirement Verification & Validation

Last, but certainly not least, we request experts in both the IT and the energy industry to examine our requirements and corresponding analyses, seeking any flaws and conflicts. Table 3.3 lists the experts involved in the verification and validation process.

Industry	Role	Size of Infrastructure Dealt w/
Energy	Machine Learning Engineer	Large
IT	Engineer and Scientist	Medium
IT	Researcher	Medium & Small

Table 3.3: Experts involved in requirement verification and validation.

Figure 3.3: Overview of the architecture of the entire system.



3.3 System Architecture

In this chapter, we elaborate on the design of the system architecture, teasing it down layer by layer in the following sections. Firstly, we give an overview of the architecture of the entire system in Section 3.3.1. Then, in Section 3.3.2, we present the blueprint of the power support subsystem. Finally, the market extension EEMM is described in the last section of this chapter (§3.6).

3.3.1 System Overview

Referring back to Chapter 1, datacenter simulators have been widely adopted in both academia and the industry. OpenDC is one of such simulation tools, which is easy-to-use with a wide range of state-of-the-art features, for example, capacity planning [4], modelling serverless computing and hosting ML workloads [119]. This work develops advanced power models and a unified energy resource chain integrated into the infrastructure of the simulator.

Figure 3.3 dissects the system in a layered manner. Most of the high-level functionalities are readily available to end-users through the frontend UI and the code API. IT professionals and experts in the energy industry can access more functionalities with fine-grained control over the simulator by directly invoking the infrastructure as a service (IaaS) interface. Such services are provided by one of six packages, the compute package, which resides in the backend of the simulator. Via the IaaS interface, users are able to specify detailed simulation setup. The simulator supports heterogeneous hardware types, topologies, and scenario portfolios. Moreover, this work further enables users to configure the energy modelling and management system in a flexible way.

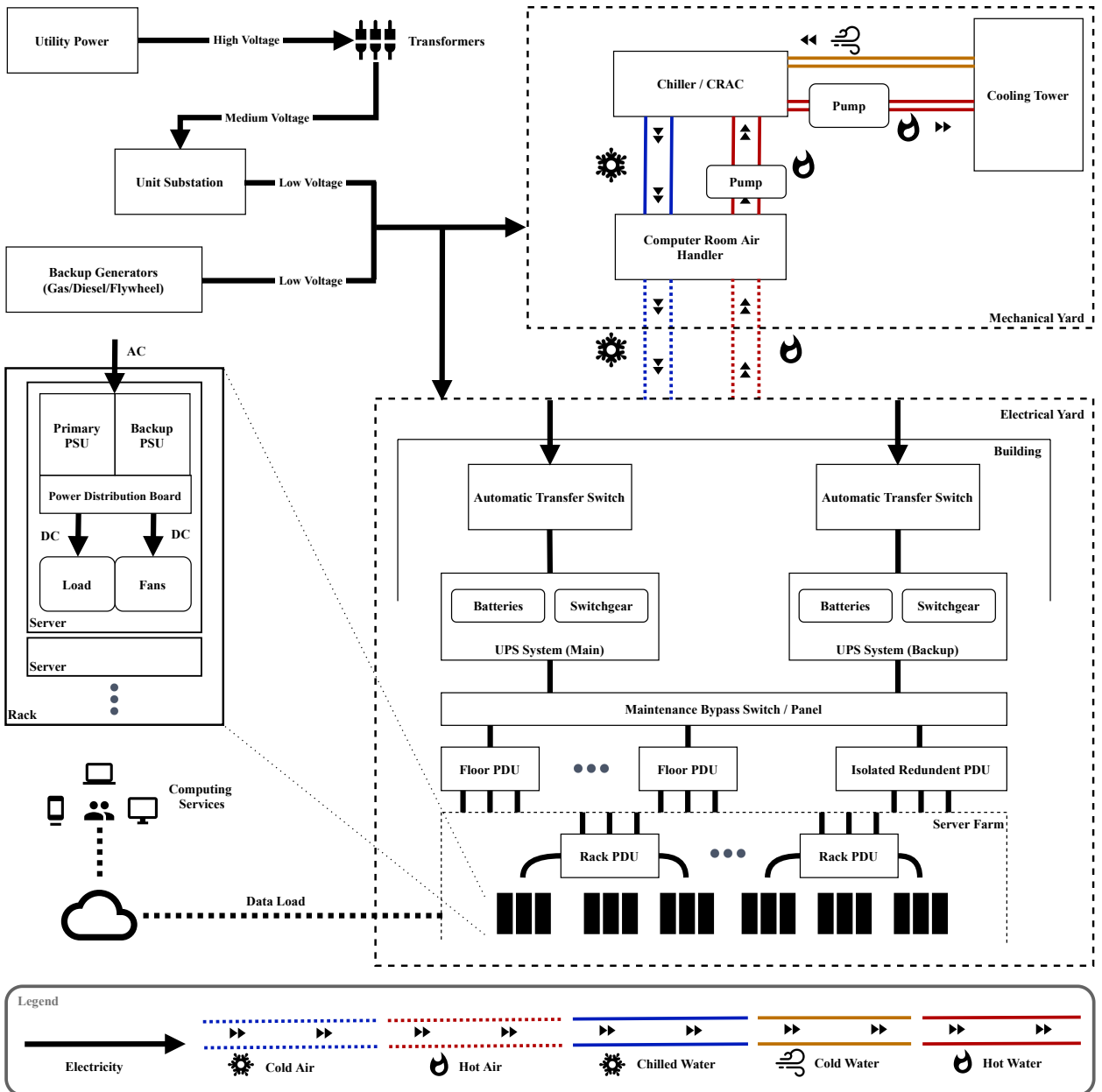
Experts can enable platform-dependent energy models and energy-saving configurations, such as a set of frequency scaling governors in the Linux kernel and various power estimation models (either generic or machine-specific), as well as different scaling drivers with customizable P-states. On top of that, users can configure the topology of the power support subsystem, which locates in the power package and is underpinned by a unified resource chain penetrated throughout the

system stack. This power support subsystem will be detailed in Section 3.3.2.

The bottom layer of the system represents the physical resources assigned by users. Instead of using their raw form directly, we abstract all resources into one single representation with a common unit. This process is realized by first aggregating the resources from all sources (*e.g.*, computing, storage, energy and so on), next interpreting different resources to unify them into a common unit, and finally, distributing the resources to support subsystems. Such an abstraction happens in the `resources` package, forming the backbone of the power modelling and management system. In this way, various physical resources with rather different units can be utilized throughout the system with one common view.

Furthermore, working towards accomplishing native support for functionalities related to the energy market, this work provides an extension of the energy modelling and management system (EEMM), which will be introduced in Section 3.3.3.

Figure 3.4: Architecture of a quintessential AC power system of datacenters modelled in this work.



3.3.2 Power Support Subsystem

To achieve the **FRs** documented in Section 3.2.3, the very first premise is reliable power modelling and flexible energy management, underpinning **FR1** and **NFR1**. This section describes the core design of the power support subsystem.

First and foremost, according to the study presented in Section 2.4.5, we construct the skeleton of a quintessential datacenter power system, upon which the power modelling subsystem is built. Figure 3.4 demonstrates the AC power system modelled in this work. First, the electricity comes into the datacenter from the utility power station. Then, the energy is distributed to two parts, the primary support in the electrical yard and the secondary support in the mechanical yard. Inside the datacenter building, the energy passes through the transfer switches and reaches the UPS systems. The UPS systems will next distribute the energy to several floor PDUs. These floor PDUs deal with a higher voltage than the rack PDUs which are directly connected with the servers. Inside of the server, the PSU transform AC power to DC and power both the computing load and the internal cooling system. Note that almost all input electricity will ultimately be exiled in the form of heat. We recognize that such an architecture may well vary from datacenters to datacenters. With this caution in mind, we design the power support subsystem that is highly customizable.

Figure 3.5 shows the design of the power support subsystem. Components of the system form four layers, namely, the raw resource layer, the aggregation layer, the distribution layer, and the resource consumption layer. To support various power system topologies, the number of components at each layer can be customized by users. The middle two layers are responsible for the core abstraction of resource interpretation and unification elaborated in Section 3.3.1. Moreover, from a higher point of view, components in the first two layers serve as power inlets that generate and aggregate energy, and the last two serve as power outlets that distribute and consume energy.

From left to right, the power chain therein starts at various power sources in the first layer and ends at the IT infrastructure, the server farm. On the contrary, the reporting of energy use follows exactly the opposite direction. This architecture facilitates the energy monitoring and data collection mechanisms one abstraction

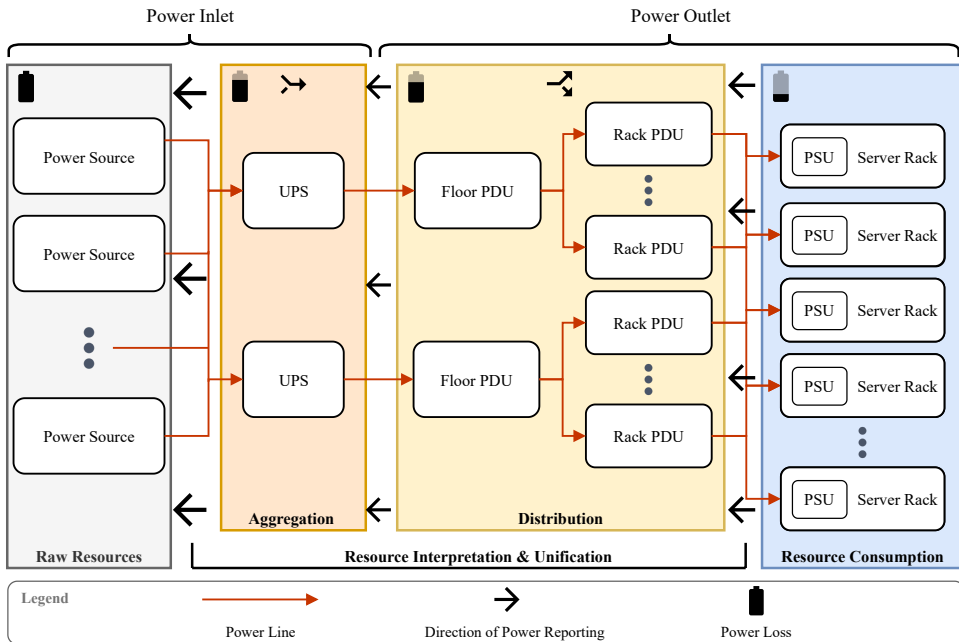


Figure 3.5: Architecture of the power support subsystem.

above, directly addressing **FR2**. The implementation of the power support subsystem described in Chapter 4 closely follows this design.

3.3.3 Market Extension

In this section, we turn our attention to addressing the market-related requirements, *i.e.*, **FR3** to **FR6**, and **NFR2** to **NFR4**.

Upon the basis of the simulation infrastructure presented in the previous two sections, we build an extension of the energy modelling and management system (EEMM). Figure 3.6 demonstrates the architecture of the extension, in which five core modules work in concert.

Industry experts can directly interact with the extension via the command-line interface provided by the `cli` modules in a Unix system. By providing the extension with the simulation results and market data, including the energy prices in

various markets, the `preprocess` module will convert data and feed them into the `market` module. Note that this step does not require any extra manual data processing from users other than the standard market data from corresponding official websites ¹², which aim at fulfilling **NFR4**.

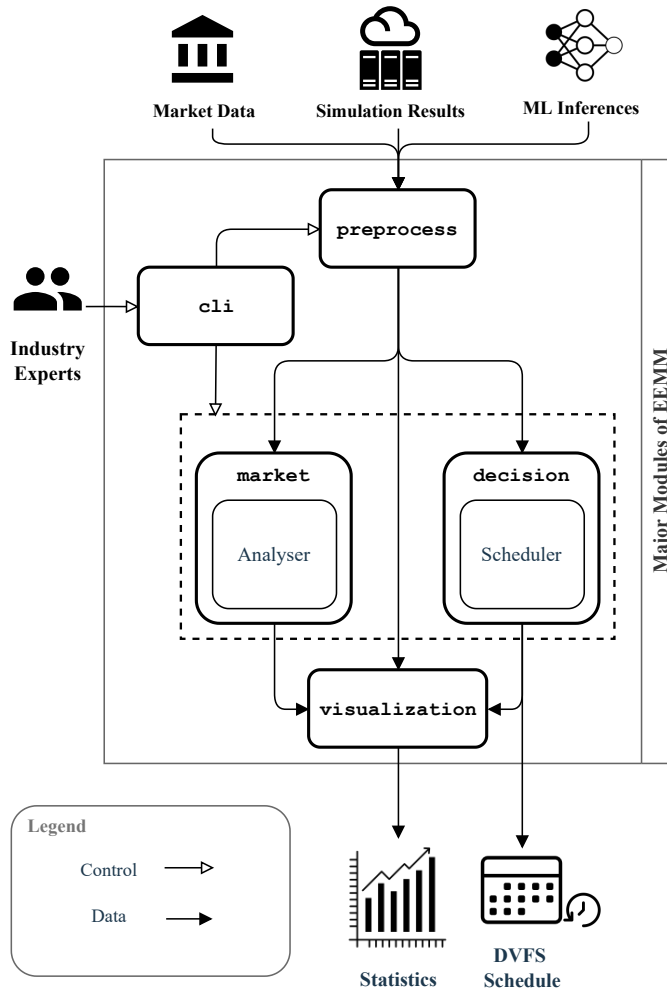
The `market` module estimates the energy consumption and provides insights of the energy costs in different markets, addressing **FR3** and **FR4**. Moreover, it contains an analyser that inspects the market data and simulates procurement strategies. Next, these analyses will be passed to the `visualization` modules that visualizes the statistics for users to compare different strategies (**FR5**).

Furthermore, to address **FR6**, the `decision` module is developed to process the ML inferences as market signals. It embodies a proactive DVFS scheduler, supporting users to make fine-grained decisions in response to the market signals.

¹<https://transparency.entsoe.eu/dashboard/show>

²https://www.tennet.org/english/operational_management/system_data_relating_processing/settlement_prices/index.aspx

Figure 3.6: Architecture of the market extension of the energy modelling and management system.



IV

Implementation

In this chapter, we take a deep dive into the implementation of the system. Firstly, the development of the energy modelling and management system is introduced in Section 4.1. Then, we move on to the implementation of the market extension EEMM in Section 4.2.

{	ConstantPowerModel	$\mathbb{P}(s) = s$	(4.1a)
	LinearPowerModel	$\mathbb{P}(u) = P_{\text{idle}} + (P_{\text{max}} - P_{\text{idle}})u$	(4.1b)
	SquarePowerModel	$\mathbb{P}(u) = P_{\text{idle}} + (P_{\text{max}} - P_{\text{idle}})u^2$	(4.1c)
	CubicPowerModel	$\mathbb{P}(u) = P_{\text{idle}} + (P_{\text{max}} - P_{\text{idle}})u^3$	(4.1d)
	SqrtPowerModel	$\mathbb{P}(u) = P_{\text{idle}} + (P_{\text{max}} - P_{\text{idle}})\sqrt{u}$	(4.1e)
	MsePowerModel	$\mathbb{P}(u) = P_{\text{idle}} + (P_{\text{max}} - P_{\text{idle}})(2u - u^r)$	(4.1f)
	InterpolationPowerModel	$\mathbb{P}(u) = P(u_1) + (P(u_2) - P(u_1))\frac{u-u_1}{u_2-u_1}$	(4.1g)
	AsymptoticPowerModel	$\mathbb{P}(u) = P_{\text{idle}} + \frac{(P_{\text{max}}-P_{\text{idle}})}{2}\left(1 + u - e^{-\frac{u}{a}}\right)$	(4.1h)
	(AsymptoticPowerModelDvfs)	$\mathbb{P}(u) = P_{\text{idle}} + \frac{(P_{\text{max}}-P_{\text{idle}})}{2}\left(1 + u^3 - e^{-\frac{u^3}{a}}\right)$	(4.1i)

4.1 Energy Modelling & Management System

Closely following the design described in Section 3.3, we implement and integrate the energy modelling and management system into the OpenDC simulator. Figure 4.1 is a simplified UML diagram of the system.

Power models occupy the lower part, of the diagram. These models include both the generic models and the machine-specific ones that can be further tuned towards a particular computing platform. These models align with their mathematical formulation shown in Equation 4.1a to 4.1i. Note that these models are implemented during the Honours research of the author (please see the report [67] for more details) but are reorganized and improved during this work. These power models are controlled by two power drivers, and the `PStatePowerDriver` utilizes the P-states provided by users to adjust power estimation in discrete steps.

Algorithm 1 illustrates the detailed mechanism of this driver. As described in Section 2.3.5, we take package-level decisions in line 10 to 12. CPU saturation is

measured by the metric specified by Equation 2.7 (§2.1) in line 27.

Algorithm 1: P-state scaling algorithm applied in the PStatePowerDriver.

Input:

A list of scaling context C associated with the CPUs of *machine*;
 A map M that contains a set of P-states as keys and power models at each
 respective frequency level as values;

Output:

The next P-state;

Data: A power consumption table T for each corresponding P-state in M ;

```

/* Initializing M according to T. */
1 foreach state, powerLevel ∈ T.entries do
2   | Choose a power model m; // Model types can vary from level to level.
3   | Instantiate m based on powerLevel;
4   | M.put(state, m);
/* Updating the current P-state. */
5 initially target ← 0;
6 initially currentUsage ← 0;
7 initially isUpdated ← false;
8 if ∃c ∈ C : c has been updated by governors then
9   | isUpdated ← true;
10  | foreach c ∈ C do
11  |   | target ← max(c.requested, target); // Take package-level decisions.
12  |   | currentUsage ← c.cpu.speed + currentUsage;
/* Locating the appropriate P-state. */
13 initially pstate ← 0;
14 if isUpdated then
15  | // The following can be simplified via a tree map instead of a normal hash map.
16  | upperBound ← max(M.getKeys());
17  | target ← min(upperBound, target);
18  | levels ← sort(T.getKeys()); // Sort in ascending order.
19  | foreach l ∈ levels do
20  |   | if level ≥ target then
21  |     | pstate ← level;
22  |     | break;
22 else
23  | pstate ← P-state from the last update;
24  | foreach cpu ∈ machine.cpus do
25  |   | currentUsage ← cpu.speed + currentUsage;
/* Computing the instant power consumption. */
26 model ← M.get(pstate);
27 u ←  $\frac{\text{currentUsage}}{\text{pstate} * C.size}$ ;
28 model.computePower(u);
29 return pstate;

```

Moreover, four generic scaling governors are developed, corresponding to the four governors found in the Linux kernel (§2.3), namely, the powersave, the performance, the conservative, and the ondemand governors. Note that the schedutil governor is not included since it is scheduler-dependent, and the governor userspace is also excluded as it is nothing but a static governor that can be easily realized by either the PowerSaveScalingGovernor or the PerformanceScalingGovernor.

Algorithm 2 shows the scaling mechanism realized in the OnDemandScalingGovernor, and that of the ConservativeScalingGovernor is presented in Algorithm 3.

Algorithm 2: Scaling algorithm in the OnDemandScalingGovernor.

Input:
 The load threshold t with default value being 0.8;
 The scaling policy P ;

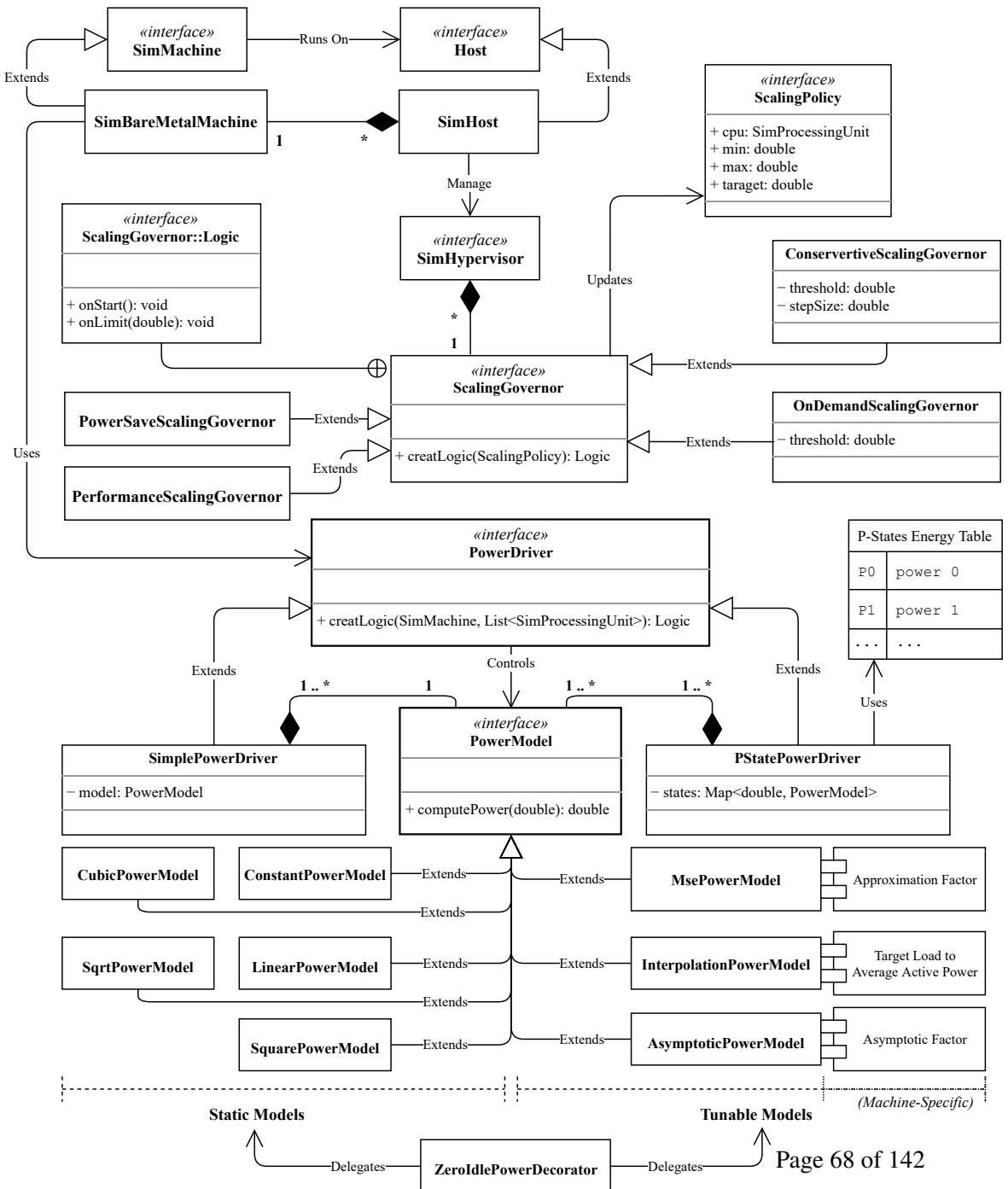
Output:
 void;

```

1 if the associated CPU has not been initialized then
2   |  $P.target \leftarrow P.min$ ;
3   | return void;
4 endif
5 if  $l > t$  then
6   | /* Proportional scaling.a */
7   |  $P.target \leftarrow P.min + l * \frac{P.max - P.min}{100}$ ;
8 endif
9 else
10  |  $P.target \leftarrow P.min$ ;
11 endif
12 return void;
```

^ahttps://github.com/torvalds/linux/blob/master/drivers/cpufreq/cpufreq_ondemand.c

Figure 4.1: Simplified UML diagram of the energy modelling and management system.



Algorithm 3: Scaling algorithm in the ConservativeScalingGovernor.**Input:**

The load threshold t with default value being 0.8;
 The step size s with default value being -1.0 ;
 The previous CPU load o ;
 The current CPU load l ;
 The scaling policy P ;

Output:

void;

```

1 if the associated CPU has not been initialized then
2   |  $P.target \leftarrow P.min$ ;
3   | return void;
4 endif
5 if  $s \leq 0$  then
6   |  $s \leftarrow P.max * 0.05$ ;           // Set the step size to the default value in the Linux kernel.a
7 endif
8 else
9   |  $s \leftarrow \min(s, P.max)$ ;
10 endif
11 initially,  $step \leftarrow -1$ ;
12 if  $l > t$  then
13   | /* Checking for load increase.                                     */
14   | if  $l > o$  then
15   |   |  $step \leftarrow +s$ ;
16   | else if  $l < o$  then
17   |   |  $step \leftarrow -s$ ;
18   | else
19   |   |  $step \leftarrow 0.0$ ;
20   | endif
21 endif
22  $P.target \leftarrow \min(\max((P.target + step), P.min), P.max)$ ;
23  $o \leftarrow l$ ;
24 return void;

```

^ahttps://github.com/torvalds/linux/blob/master/drivers/cpufreq/cpufreq_conservative.c

The `ScalingPolicy` is implemented to model the struct `cpufreq_policy` object in the Linux kernel, as introduced in Section 2.3. It contains essential information concerning the associated CPU, as well as the target frequency modulated by the scaling governors. This object is initialized by the `SimHypervisors` which operates scaling governors thereof. To achieve a flexible initialization for various types of CPUs as an effort to address **NFR1**, the inner interface `ScalingGovernor::Logic` is introduced, serving as a factory to assist this initialization process. `SimHypervisors` are managed at runtime by `SimHosts` that run on a `SimBareMetalMachine`. A bare-metal machine is able to invoke `PowerDrivers` at the physical layer for energy estimation.

With regard to the implementation of the power support subsystem, its development adheres to the UML diagram shown in Figure 4.2. Referring back to the design described in Section 3.3.2, the `SimPowerOutlet` and the `SimPowerInlet` classes represent the two top categories in Figure 3.5, inherited by various component classes at different layers. A `SimPsu` that resides in a bare-metal-machine directly interacts with the aforementioned `PowerDrivers`. Furthermore, power distribution is achieved by a simple max-min fair-sharing policy, illustrated by Algorithm 4. Note that the system has been integrated with `OpenDC`, which is fully containerized and can be easily run via `docker` (**NFR4**).

Algorithm 4: Max-min fair-sharing power distribution algorithm.

Input:
 A list of demands D ;
 The capacity c of the power sources;

Output:
 The remaining capacity of the power source;

```

1  $D \leftarrow \text{sort}(D)$ ; // Sort  $D$  in ascending order.
2  $n \leftarrow \text{sizeof}(D)$ ;
3 initially,  $\text{ration} \leftarrow \frac{c}{n}$ ;
4 initially, allotments  $A \leftarrow$  an empty list of size  $n$ ;
   /* Initial assignment. */
5 for  $i \leftarrow 0$  until  $n$  by 1 do
6   |  $A[i] \leftarrow r$ ;
7 endfor
   /* Handling overloaded demands. */
8 for  $i \leftarrow 0$  until  $n$  by 1 do
9   |  $\text{ration} \leftarrow A[i]$ ;
10  | initially,  $\text{share} \leftarrow 0.0$ ;
11  | if  $i < (n - 1)$  and  $\text{ration} \geq D[i]$  then
12  |   |  $\text{share} \leftarrow \frac{\text{ration} - D[i]}{n - (i + 1)}$ ; // Fair-share over-supplied capacity.
13  |   | for  $j \leftarrow i$  until  $n$  by 1 do
14  |   |   |  $A[j] \leftarrow A[j] + \text{share}$ ;
15  |   |   endfor
16  |   |  $A[i] \leftarrow D[i]$ ;
17  |   endif
18  |   else
19  |   |  $A[i] \leftarrow \min(\text{ration}, D[i])$ ;
20  |   endif
21 endfor
22 return  $c - \sum_i^{n-1} A[i]$ ;
```

4.2 Market Extension

The core of the extension EEMM is the DVFS scheduler. Referring back to Section 2.3, the major drawback of enabling DVFS is the prolonged execution time. By the virtue of the dominating quadratic relationship shown in Equation 2.9, the benefit

of the power saving brought by DVFS will ultimately outweigh the linearly scaled overhead, *i.e.*, the longer duration. That been said, such an overhead incurred by using DVFS is by no means negligible and should be further mitigated. Therefore, we take this factor into account when developing the scheduler.

When hosting traces of the virtual machines (VMs), our simulation currently does *not* prolong the execution time of the jobs in the case that the capacity of the host is reduced. Instead, the overhead of DVFS is reflected by the over-commission of the CPUs. To elaborate on this further, when the frequency of the CPUs are restrained to a much lower level, the total capacities of all VMs may well (temporarily) exceed the maximum capacity of their host. We do not scale the VMs downwards to accommodate the reduced hosting capacity but capture the over-commissioned CPU cycles instead. Similar circumstances could also occur when new VMs are spawned, which can lead to exceeding the total capacity of the host. In contrast, as part of the capacity of the host is released (*e.g.*, some VMs have finished the execution or the frequency of the CPU is increased), the scaling driver in our simulator does not stick to the proposed CPU frequencies by the scaling governor and level down the resource processing speed accordingly.

Thus, the DVFS scheduler should juggle both the energy consumption affected by switching the scaling governor and the overhead of using DVFS, the CPU over-commission. To this end, we introduce the concept of *damping factor*, which is inspired by the Google PageRank algorithm [141] with different implications and implementations. The damping factor, in this case, is not a probability value but a threshold that restrains the increase in the level of CPU over-commission. The lower the damping factor, the more frequent the scheduler ameliorates the constraint on the CPU frequency by switching to a more performant scaling governor. Algorithm 5 illustrates such a scheduling strategy implemented in the EEMM. We will further elaborate on the two decision points (line 11 and 15) in Chapter 5.

In an effort to address **NFR4**, we develop the market extension as a Python library, which can be easily installed by using the following command in a Unix system.

```
$ pip install openc-eemm
```

Algorithm 5: DVFS scheduling algorithm implemented in the market extension.

Input:

The spot price p^S of the next ISP;
 The forecasted imbalance shortage price p^F of the next ISP;
 A list of available scaling governors G ;
 The damping factor d ;
 The current damping factor counter c ;

Output:

The next scaling governor to use;

Data: A series of datacenter traces T up until now;

```

1 initially  $prev \leftarrow$  get previous over-commission level from  $T$ ;
2 initially  $curr \leftarrow$  compute current over-commission level from  $T$ ;
3 initially  $governor \leftarrow null$ ;
   /* Gauging the over-commission status. */
4 if  $curr > prev$  then
5   |  $c++$ ; // Record the increase of the current over-commission level.
6 endif
7 else
8   |  $c--$ ;
9 endif
10  $prev \leftarrow curr$ ;
   /* The first decision point. */
11 if  $p^F \leq 0$  then
12   |  $governor \leftarrow G.performance$ ;
13 endif
14 else
   /* The second decision point. */
15   if  $p^F > p^S$  then
16     |  $governor \leftarrow G.powersave$ ;
17   endif
18   else
19     if  $c \geq d$  then
20       |  $governor \leftarrow G.conservative$ ;
21       |  $c \leftarrow 0$ ; // Relax the over-commission meter.
22     endif
23     else
24       |  $governor \leftarrow G.ondemand$ ;
25     endif
26   endif
27 endif
28 return  $governor$ 

```



Evaluation

In this chapter, we elaborate on the experiments conducted to answer **RQs** described in Section 1.2, as well as to meet the **FRs** and **NFRs** listed in Section 3.2.3. We start with detailing the setup in Section 5.1. Then, in Section 5.2, we analyse the results regarding the participation of datacenters in the energy market. Lastly, in Section §5.3, we evaluate the performance of the proactive DVFS scheduler powered by ML methods.

5.1 Experiment Setup

In this section, we describe the simulation model employed in the experiments, specifically, the market model (§5.1.1), the specifications of machines simulated (§5.1.2), and the energy models employed for estimating the energy consumption of the datacenter (§5.1.3).

5.1.1 Market Model

Firstly, referring back to Section 2.4, there are two financial markets prior to the start of the real-time delivery, the day-ahead market and the intraday market. The energy prices settled in these two markets are often referred to as the spot price.

The major difference between the two is that the day-ahead market has one single settlement for all participants, whereas the intraday market consists of continuous bilateral trading carried out through the trading day and, in turn, has no common settlement. As described in Section 2.4.2, the intraday market is where prosumers conduct the final adjustments to their self-dispatched quantities. Besides the varying prices, the number and types of products are highly dependent on the surplus/deficit situation of every participant. Therefore, in our experiment, we do not make any assumption in the participation of the intraday market, *i.e.*, we do not make any adjustment in the quantity of energy ordered in the day-ahead market before the start of the actual delivery period.

Secondly, in respect of the day-ahead market and the balancing market, we focus on the energy market of the Netherlands. In other words, all prices used in the following experiments are the energy prices of the Netherlands only. Similarly, although the energy trading system introduced in Section 2.4 applies across the EU, there are still nuances between different countries. For this matter, we also only focus on the trading system of the Netherlands, but the experiment results do not lose their generality as these nuances are rather minute.

Lastly, datacenters normally either have long-term contracts with their utility companies or buy energy in an on-demand scheme. Since it is not feasible to obtain/disclose the energy prices of the bilateral, long-term contracts, we only focus on the on-demand scheme. Note that, in practice, datacenters generally do not have significant discounts through long-term contracts due to the lack of intermittency, but we, nevertheless, want our experiment results to be inclusive and representative. Therefore, we consider three on-demand energy prices from low to high summarized in Table 5.1.

Price Level	Price [€/MWh]	Source
Low	38.0	NieuweStroom B.V. (2021, average) [23]
Medium	56.5	PricewaterhouseCoopers (2017, average) [134]
High	80.4	Essent N.V. (2021, fixed) [125]

Table 5.1: On-demand energy prices considered in the experiments.

5.1.2 Machine Model

DVFS technology is one of the major interests of this work. Although the frequencies and voltages of P-states in different CPUs are commonly available, their corresponding consumption levels, however, need to be specially measured by either software or hardware power meters. That said, since neither developing an instrument for measuring the P-state power nor testing the accuracy of the measurements is within the scope of this work, we resort to the existing literature for this matter.

Frequency Steps [MHz]	1600	1867	2113	2400	2670
Idle Power [Watt]	82.70	82.85	82.95	83.10	83.25
Max. Power [Watt]	88.77	92.00	95.50	99.45	103.0

Table 5.2: P-states consumption levels of the old machine model.

For the consumption levels of P-states, the latest reputable reference that we found is [61]. However, the CPU therein is old and, therefore, might not be representative in terms of its power consumption. To understand the impact of using this old machine model, we also include a recent machine with a type of CPU released this year (2021) from the SPEC benchmark [21]. These two machines models are summarized in Table 5.3, and the P-states consumption levels of the old machine model are in Table 5.2.

Machine Model	Year of Release	CPU	Base Frequency	Cache	#Cores	#Threads
Old	2007	Intel® Core™2 Quad Q6700	2.66 GHz	8 MB	4	4
New	2021	Intel® Xeon® Platinum 8380	2.30 GHz	60 MB	40	80

Table 5.3: Machine models used in the experiments.

In the experiments, we adopt a set of Business Critical Workload (BCW) traces [151] from the Dutch IT service provider Solvinity, containing records monitored over a course of one month. To host the traces properly without overloading/underutilizing the hosts, it is necessary to carefully calculate the resources needed by each of the two machine models in order to set up the compute service at the IaaS layer (Figure 3.3). The number of hosts (N_{Hosts}) for each machine model is

computed following Equation 5.1, where Ψ_d denotes the maximum instant CPU demand of the traces, the Ψ_f represents the maximum frequency of the machine, and c is the core count of the CPU package. Similarly, the number of memory units N_{Units} populated in each machine is calculated by Equation 5.2, where Ψ_m is the maximum instant memory request of the traces, and m denotes the size of the memory unit.

$$N_{\text{Hosts}} = \left\lceil \left\lfloor \frac{\Psi_d}{\Psi_f} \right\rfloor / c \right\rceil \quad (5.1)$$

$$N_{\text{Units}} = \left\lceil \left\lfloor \frac{\Psi_m}{N_{\text{Hosts}}} \right\rfloor / m \right\rceil \quad (5.2)$$

Note that, although the results are rounded up at each step, CPU over-commission can still occur since the original traces contain requests that exceed the capacity of the VMs in the first place [149]. Also, the computing node used in the SPEC benchmark of the new machine model contains two packages; we set the number of cores per host to the total number of logical cores, which determine the actual capacity, as opposed to the number of dies/chips. The detailed setup for each machine model is summarized in Table 5.4.

Machine Model	#Cores/Host	#Host	Size of Memory Unit [MB]	#Memory Units/Host
Old	4	284	4,000	4
New	160	9	3,200	48

Table 5.4: Host setup for each machine model.

5.1.3 Energy Model

Critical Load. According to the results of the author’s Honours research [67], when no particular computing platform is assumed, the linear power model (Equation 4.1b) and the square root power model (Equation 4.1e) are able to properly bound the total energy consumption of the critical load for a specific machine. Therefore, we use these two power models for the old machine model, which are

referred to as “LINEAR” and “SQRT” in the following experiments.

In the case that platform-specific data is available, machine-specific models are preferred over generic ones. Since we have the load-to-power data of the new machine model from the SPEC benchmark, we use the interpolation power model (Equation 4.1g) for the new machine model, which is referred to as “INTERPOLATION” in the experiments. The load-to-power data of the new machine model is presented in Table 5.5.

Model \ Load	0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
New	118.0	188.0	224.0	258.0	273.0	298.0	333.0	380.0	430.0	492.0	633.0

Table 5.5: Target loads to average active power values (in Watts) for the new machine model.

With regard to the PSU, we employ the power model from another state-of-the-art datacenter simulator, iCanCloud/E-mc² [26, 124]. To model the power losses of PSU, we first compute the percentage of load (η_l) following Equation 5.3 based on the power draw of the server (P_{server}) and then, maps it to the corresponding energy efficiency (η_e) by Equation 5.4, where both the rated output power (τ) and the mapping for the (η_e) are specified in manufacturers’ datasheets. In this work, the energy efficiency of the PSU adhere to the 80 Plus Titanium standard, and the rated output power τ is set to a value 870 W, which is arbitrary but can be commonly found in PSU products of datacenter servers.

$$\eta_l = \frac{100}{\tau} \cdot P_{\text{server}} \quad (5.3)$$

$$\eta_e = \begin{cases} 90\% & \text{if } 0 \leq \eta_l \leq 10\% \\ 94\% & \text{if } 10 < \eta_l \leq 20\% \\ 96\% & \text{if } 20\% < \eta_l \leq 50\% \\ 91\% & \text{if } 50\% \leq \eta_l < 100\% \end{cases} \quad (5.4)$$

Having computed the energy efficiency, we then accumulate the energy consump-

tion of the PSU (E_{PSU}) following Equation 5.5.

$$E_{\text{PSU}} = \int_{t_0}^{t_{\text{PSU}}} \frac{(P_{\text{server}} \cdot 100)}{\eta_e} - P_{\text{server}} dt \quad (5.5)$$

Primary Support. To estimate the energy consumption of the primary support, specifically, the UPS systems and the PDUs, we employ the power model proposed by Rasmussen [139], which has been widely adopted in estimating the power consumption of the UPS and PDU over the years. Rasmussen suggests that it is insufficient to model the consumption of components of the power support system by using only the single parameter, the nameplate loss coefficient, specified in the datasheets from the manufactures. Instead, their energy consumption should be captured by two values, the tare loss coefficient (π) that is independent of the load, and a polynomial loss coefficient that is load-dependent. In the case of PDU, the load-dependent part of consumption is captured by the proportional loss coefficient (α), and that of the UPD is captured by the square-law loss coefficient (β). The coefficients are summarized in Table 5.6, based on the work from Rasmussen.

Component	λ	α	β	π
UPS	0.040	0.050	—	0.090
PDU	0.015	—	0.015	0.030

Table 5.6: Coefficient values of the UPS and PDU power model.

$$\left\{ \begin{array}{l} \alpha = \pi_{\text{UPS}} - \lambda_{\text{UPS}} \\ P_{\text{UPS}}^{\text{tare}} = \lambda_{\text{UPS}} \cdot P_{\text{UPS}}^{\text{rated}} \\ P_{\text{UPS}}^{\text{loss}} = P_{\text{UPS}}^{\text{tare}} + \alpha \cdot (\sum_i^{N_{\text{PDU}}} P_{\text{PDU}_i}^{\text{in}}) \end{array} \right. \quad (5.6)$$

$$\left\{ \begin{array}{l} \beta = \pi_{\text{PDU}} - \lambda_{\text{PDU}} \\ P_{\text{PDU}}^{\text{tare}} = \lambda_{\text{PDU}} \cdot P_{\text{PDU}}^{\text{rated}} \\ P_{\text{PDU}}^{\text{loss}} = P_{\text{PDU}}^{\text{tare}} + \beta \cdot (\sum_i^{N_{\text{server}}} P_{\text{server}_i}^{\text{in}})^2 \end{array} \right. \quad (5.7)$$

We compute the energy consumption of the UPS and PDU using Equations 5.6 and 5.7, where P^{in} denotes the inlet power, P^{tare} denotes the tare power loss, P^{loss} denotes the total power loss, P^{rated} denotes the nameplate power, N_{server} denotes the number of active servers, and N_{PDU} denotes the number of attached PDUs.

Secondary Support. Without presuming the consumption rate of every single piece of equipment, such as the cooling tower, the CRAC, the backup generator, etc. (Figure 3.4), which can differ greatly from datacenter to datacenter, we make use of the PUE value to estimate the energy consumption of the entire datacenter as a whole. We compute the power consumption of the secondary support ($P^{2\text{nd}}$) using Equation 5.8, based upon the power draw of the servers (P_{server_i}), the UPS systems (P_{UPS_i}), and the PDUs (P_{PDU_i}).

$$P^{2\text{nd}} = \text{PUE} \cdot \sum_i^{N_{\text{server}}} P_{\text{server}_i} - \left(\sum_i^{N_{\text{PDU}}} P_{\text{PDU}_i} + \sum_i^{N_{\text{UPS}}} P_{\text{UPS}_i} \right) \quad (5.8)$$

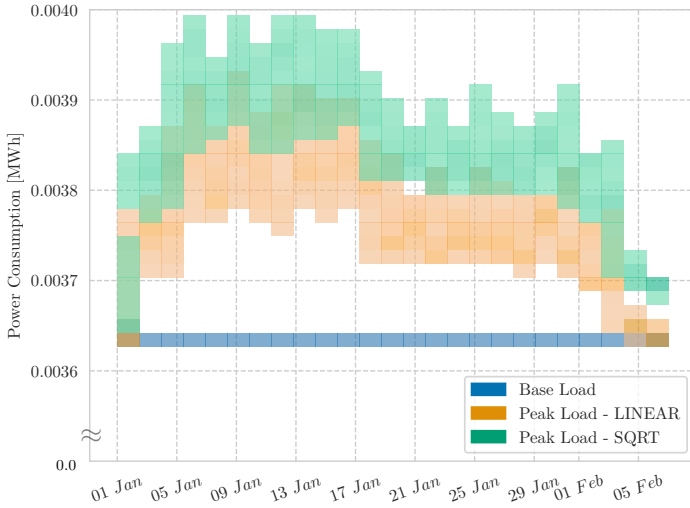
5.2 Energy Market

In this section, we present the results regarding the participation of datacenters in the energy market. First and foremost, for **RQ2** we ascertain whether it is financially beneficial for datacenters to participate in the energy market in the first place and if so, which market to participate in (§5.2.1 and §5.2.2). Next, in Section 5.2.3 we investigate the implications of using the old machine model. Then, in Section 5.2.4 we research the impact of different energy procurement strategies (**RQ3**), *i.e.*, how to participate in the day-ahead and the balancing market? Finally, we investigate why ML methods could be of help in further leveraging profits during market participation.

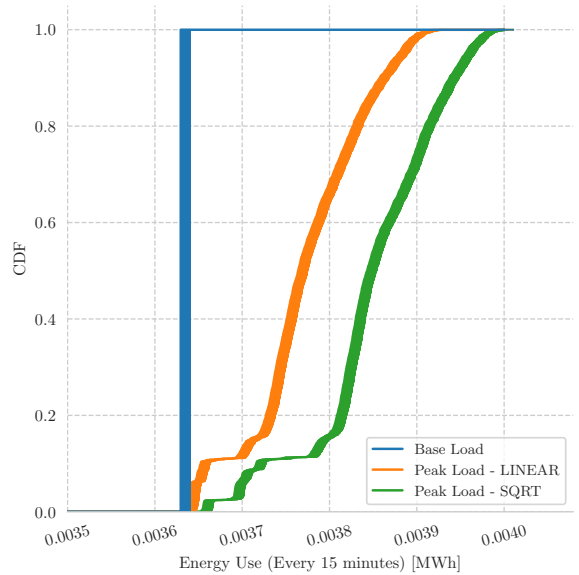
5.2.1 Power Loads

The concept of power loads plays a key role in the following experiments. Power load is the quantity of energy consumed by all equipment in a datacenter. It can be further categorized into two sub-classes, the base load and the peak load. Base load is the bare minimum active power required to keep the datacenter up and running, whereas peak load is demand-dependent. In other words, the base load is the constant power draw at all times, and peak load can be perceived as the margin between the power draw at peak demands and the constant base load.

Specifically, we define the active idle power of the datacenter as the base load, and the proportion of power draw catering for the demands as the peak load. Figure 5.1 demonstrate these two loads estimated by the two power models, where Figure 5.1a shows the instant power loads, and Figure 5.1b shows their cumulative distribution functions (CDF). Note that we add on the base load to the peak load to demonstrate their relationship, the peak load would be below the base load otherwise. As described in Section 5.1, the two generic power models, the LINEAR and the SQRT are able to serve as the lower and upper bounds for a specific computing platform. Therefore, the actual power load should lie between the two curves.



(a) Instant power loads.



(b) CDF of the power loads.

Figure 5.1: Different power loads.

5.2.2 Energy Costs

Now, we investigate the energy costs of the power loads in the day-ahead and the balancing markets. The market data analysed runs from 30 November 2019 to 11 May 2021, which is about 1.5 years in duration. Figure 5.2 illustrates the distribution of the energy prices of the two markets. Note that the positive price is the price that a datacenter, as a balance responsible party (BRP), would have to pay for its energy use. On the contrary, if the price is negative, datacenters will be paid by the energy market or the system operator to consume energy (mostly with the purpose of balancing the power grid). We can see from Figure 5.2 that the balancing market generally has a much higher positive price than that of the day-ahead market. However, the level of negative price in the balancing market is much lower than that of the day-ahead market to a similar extent. Albeit larger variation, the median of the energy price in the balancing market is slightly below the mean of the energy price in the day-ahead market. Hence, we conclude that the balancing market demon-

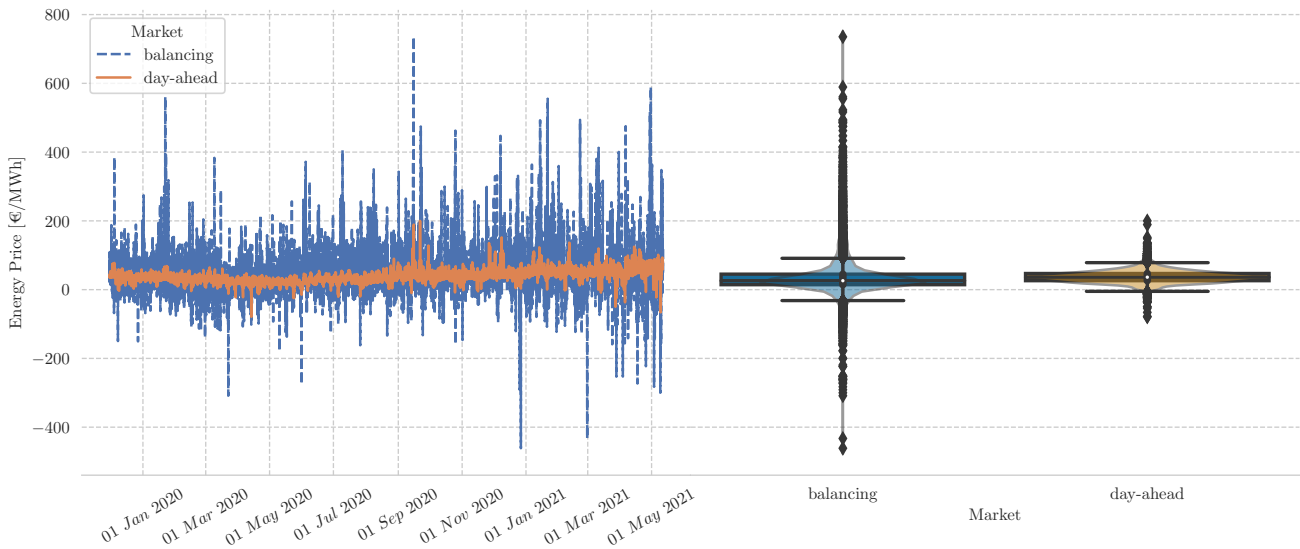
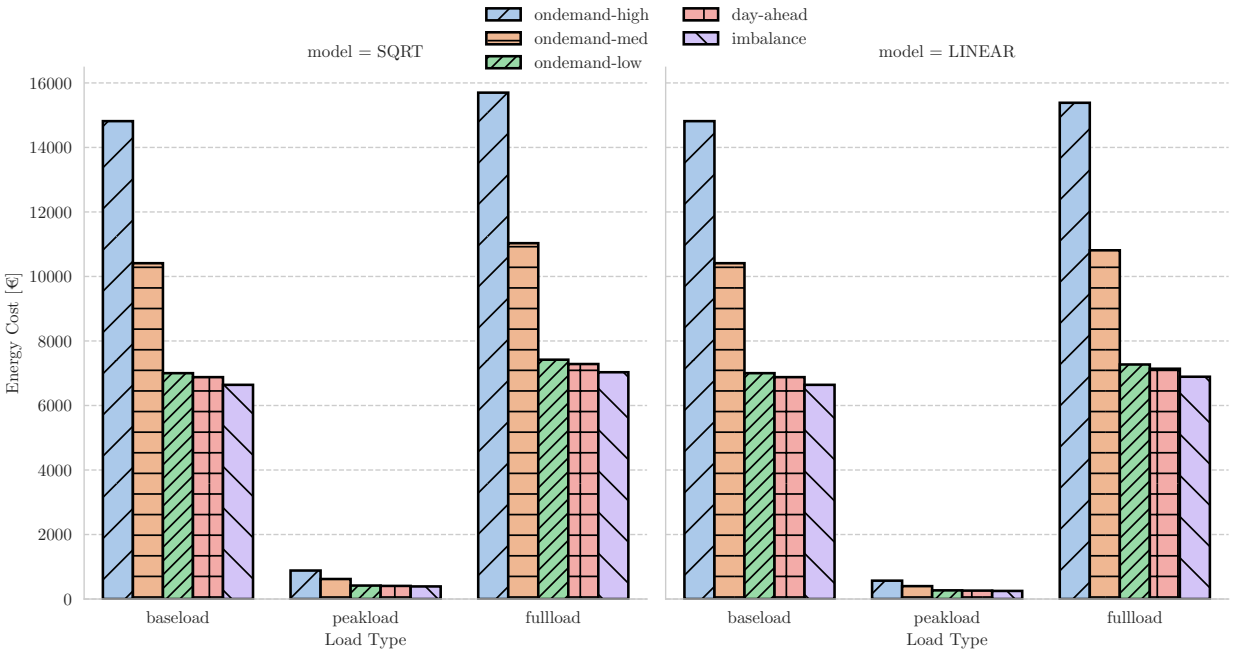


Figure 5.2: Distributions of day-ahead prices and imbalance prices.

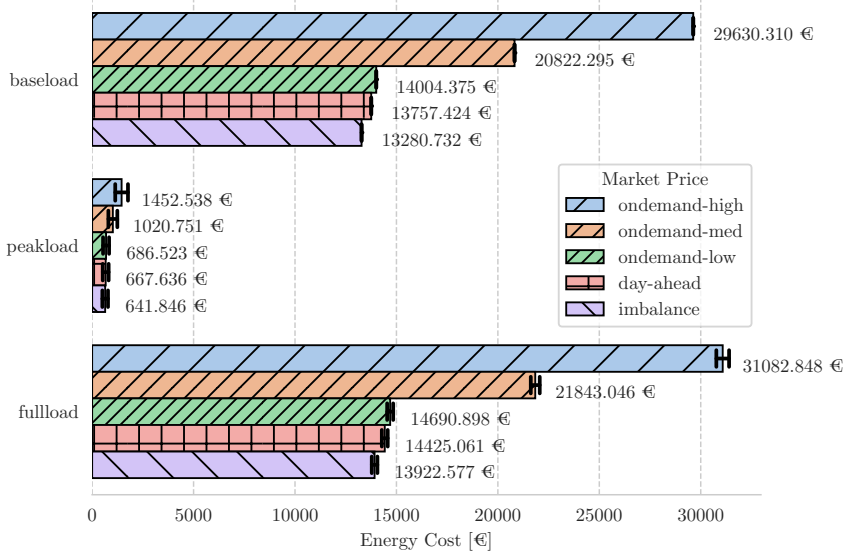
states a coexistence of *high risks* (greater positive price and larger variation) and *high profitability* (lower negative price and lower median value). Having said that, is it beneficial in the long run to take part in the day-ahead and/or the balancing market?

To answer the aforementioned question, we use the EEMM extension to compute the energy costs for the two power loads of the *old* machine model under different market prices. Figure 5.3a shows the costs of the power loads estimated by the two power models in different markets (**FR4**), and Figure 5.3b shows the combined estimation (**FR3**), in which the error bars capture the variation therein. As we can see that participating in either the day-ahead market or the balancing market results in lower energy cost even compared to the energy cost of the on-demand scheme of the lowest price. Also, the cost in the balancing market is even lower than that in the day-ahead market. Therefore, we conclude that it is financially beneficial to participate in both the day-ahead market and the balancing market (note that, referring back to Section 2.4.2, it is even compulsory for a BRP to resolve its unbalance in the balancing market).

Figure 5.3: Energy Costs of the two power loads in different markets.



(a) Energy costs estimated by two power models.



(b) Combined cost estimation from two power models.

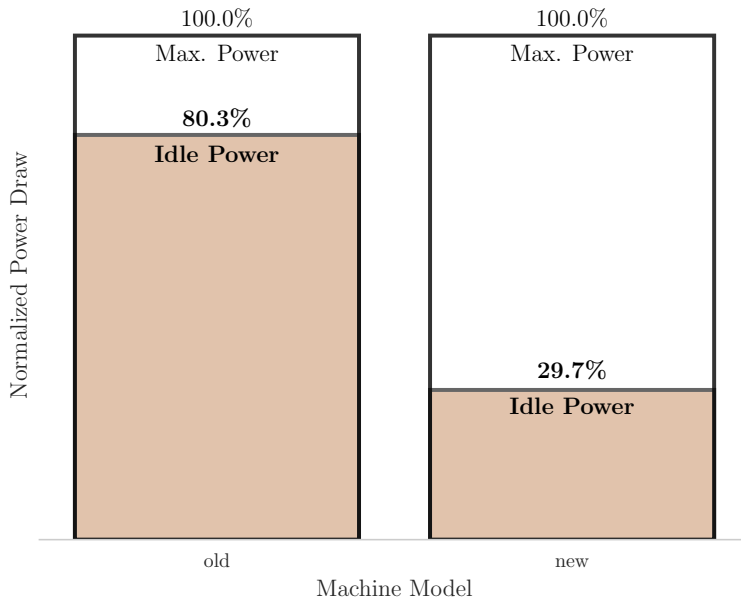


Figure 5.4: Comparison of CPU energy efficiency of two machine models.

5.2.3 Newer Machine Model

Having ascertained the incentive of participating in both the day-ahead and the balancing market for **RQ2**, we now investigate the implementations of using the old machine model.

Figure 5.4 illustrates the difference in the energy efficiency of the two machine models in terms of the idle-maximum power ratio. The idle power of the old machine model takes up more than 80% of its total power capacity. In contrast, the idle power of the new machine model only accounts for less than 30% of the maximum power. Thus, the new machine is drastically more energy efficient compared to the old one.

When it comes to the power loads, the new machine model yields a much lower level of energy consumption whilst, in the meantime, exhibits a much greater variation in the peak load (Figure 5.5). Consequently, the new machine model results in a $\sim 73.6\%$ total cost-saving, as shown in Figure 5.6. Most importantly, the por-

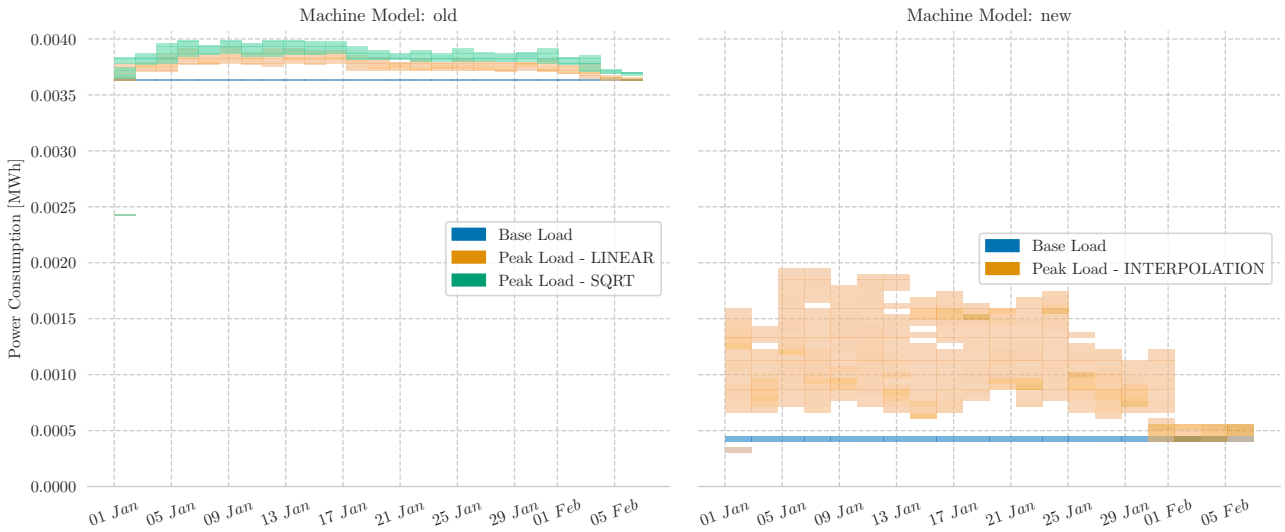


Figure 5.5: Comparison of power loads of the two machine models.

tion of energy costs resulted from the peak load of the new machine model is about $2.5\times$ higher than that of the old machine model. Nevertheless, the conclusion from Section 5.2.2 still hold, *i.e.*, participating in both the day-ahead and the balancing markets is beneficial in terms of the energy cost. These observations are important for generalizing the conclusions in later sections. From the next section onwards, we come back to the old machine model, through which the core experiments, the DVFS optimization, is conducted.

5.2.4 Energy Procurement Strategy

Results from previous sections (§5.2.2 and §5.2.3) illustrate the reason as to why datacenters should participate in the energy market. In the following sections, we investigate the impact of employing different load-forecast-based procurement strategies for **RQ3**, *i.e.*, how datacenters should participate? To answer this question, we start with stating the assumptions (**As**) made for the experiments.

A1 Datacenter operators purchase energy in the day-ahead market based on the

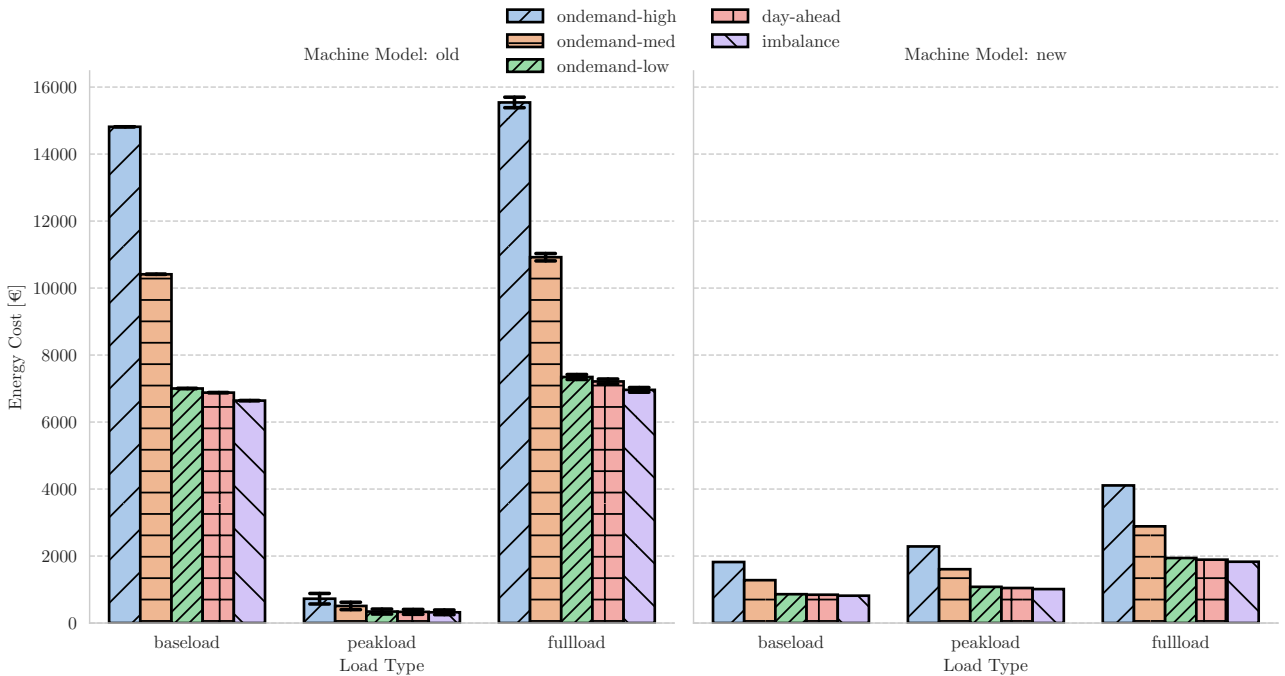


Figure 5.6: Comparison of energy costs of the two machine models.

load forecast of the corresponding *day*.

A2 The load forecast is *perfect*. In other words, the load predictions always precisely match the actual loads of the corresponding day.

A3 Datacenter operators do not deliberately schedule even less energy than the bare-minimum quantity — the base load.

A4 Datacenters' participation is abided by the two-price balancing system.

To elaborate on the two assumptions further, one should recall from Section 2.4.1 that market participants self-dispatch the quantity of energy that they are expected to produce/consume in *each hour* of the corresponding day during the day-ahead market (before the gate closure of the spot market). In turn, **A1** orientates the experiment to the average scenario, as the load forecast can be more frequent (*e.g.*,

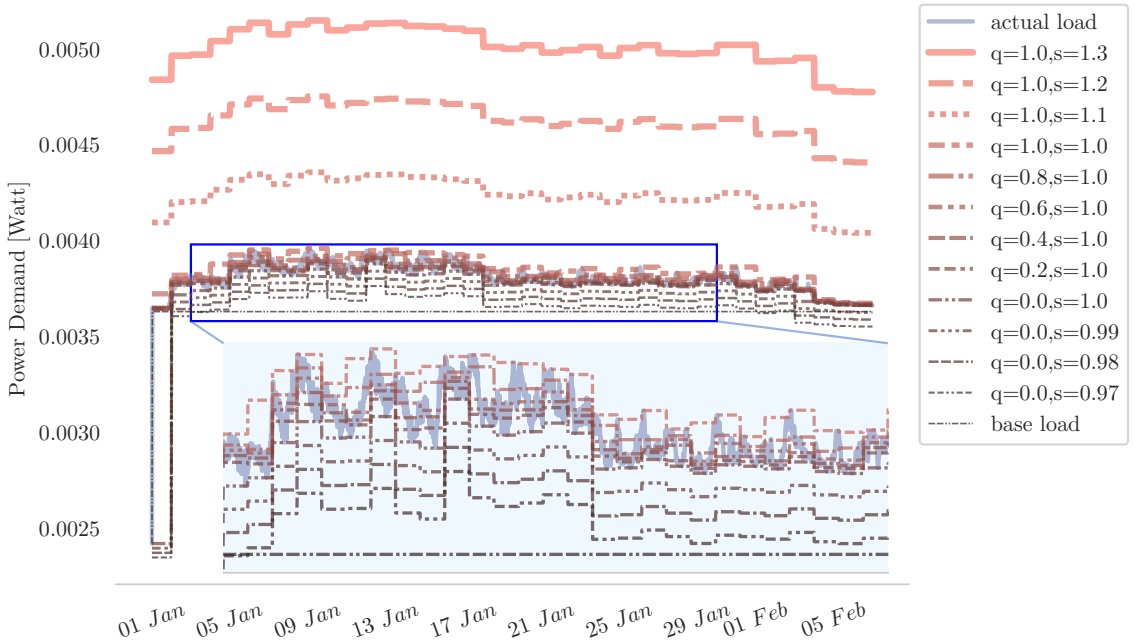


Figure 5.7: Simulated load schedules in the day-ahead market.

on an hour-basis) or be in a less proactive manner (*e.g.*, weekly or even monthly). Concerning **A2**, the performance of the load forecasting of a datacenter is the controlled variable here. Hence, we do not presume its level of accuracy/precision. As for **A3**, datacenter operators are not expected to intentionally introduce a large energy deficit in the first place. Concerning **A4**, the reasoning will be explained in detail in Section 5.2.4.1.

On the basis of **A1** to **A3**, we define the quantity of energy to schedule in the spot market (Q^S) using Equation 5.9:

$$Q^S = \mathbb{Q}_q(l_f) \cdot s, \quad (5.9)$$

where l_f denotes the load forecast of the next day, q is the quantile of the quantile function \mathbb{Q} , and s is a scalar to apply.

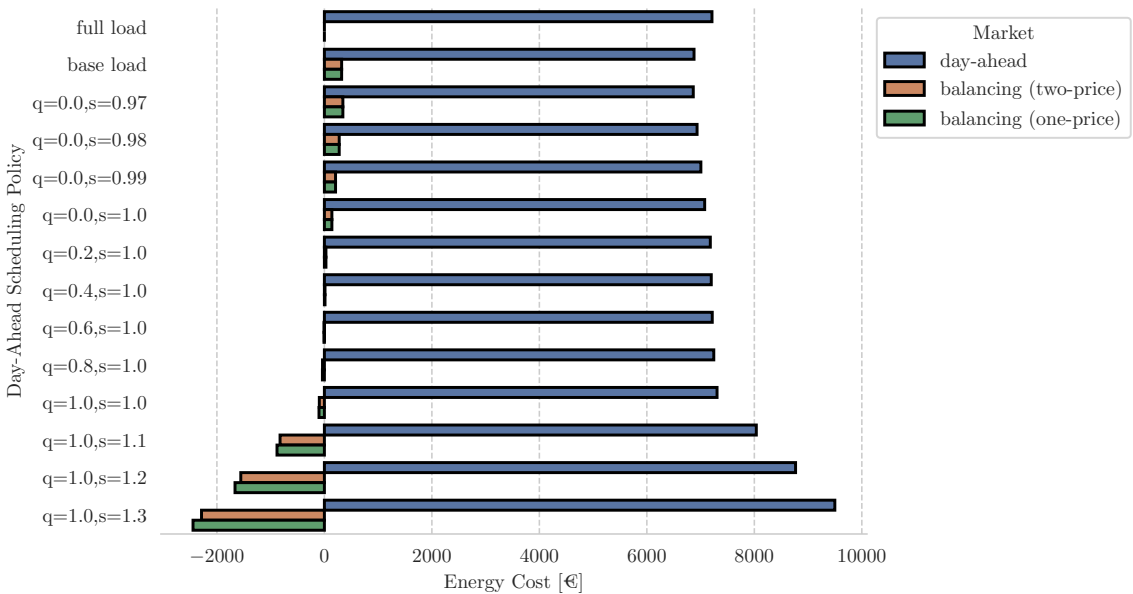


Figure 5.8: Comparison of energy costs of simulated load schedules.

In Figure 5.7, we use the extension EEMM to demonstrate the resulting load schedules in the spot market (**FR5**), where $q \in [0, 1]$, and $s \in [0.97, 1.30]$. Specifically, $q < 1.0$ models the situation in which datacenter operators schedule less power than the forecasted load of the day (under-scheduling), and $s > 1.0$ models the over-scheduling scenario in which the operators schedule more energy than the forecasted load. In addition, the “base load” represents the strategy that in the day-ahead market, the operator only procures the bare-minimum energy whose quantity is *certain*. As shown in the figure, the portion where $s > 1.0$ is sufficiently high (**NFR2**), whilst the part where the $0.0 \geq q < 1.0$ is still lower-bounded by the base load.

Having obtained the load schedules of the day-ahead market, we use the EEMM extension to compute the *total* energy cost for every schedule (**FR5**), where the energy surplus/deficit is resolved in the balancing market. Figure 5.8 shows an unstacked comparison of the energy costs for the simulated schedules. It illustrates that the more energy ordered in the day-ahead/spot market, the more datacenters have to pay. In contrast, as the amount of scheduled energy increases, the energy

cost in the balancing market decreases from positive to negative values. As described in Section 5.2, datacenters will be paid back when the energy price is negative; therein lies the question: can datacenters deliberately (or even maliciously) schedule arbitrarily large amounts of energy in the day-ahead market so that they eventually would gain profit during the imbalance settlement? To answer this question, we need to take a deep dive into the balancing mechanism of the (EU) energy market.

5.2.4.1 Imbalance Pricing Systems

p^S	Spot market price
p_-^B	Shortage price in the balancing market
p_+^B	Surplus price in the balancing market
Q_-	Quantity of shortage energy of a BRP
Q_+	Quantity of surplus energy of a BRP
Q_{\downarrow}	Quantity of energy required by downwards regulations
Q_{\uparrow}	Quantity of energy required by upwards regulations
N_{BSP}	Number of participated BSPs in the balancing market
N_{BRP}	Number of participated BRPs in the balancing market

Table 5.7: Symbols used in defining the two balancing systems.

As indicated in Figure 5.8, there are two balancing systems: one-price and two-price systems [128]. In the case of the one-price system, if the imbalance of a prosumer is (unintentionally) helping balance the grid, the prosumer will in effect earn extra monetary rewards during the imbalance settlement. On the contrary, under the two-price balancing system, such inadvertent assistance is not encouraged since the price level of the compensation is the same as that of the spot market. The mechanism of the one-price system follows Equation 5.10, and the two-price systems adheres to Equation 5.11; Table 5.7 shows the meanings of the symbols therein.

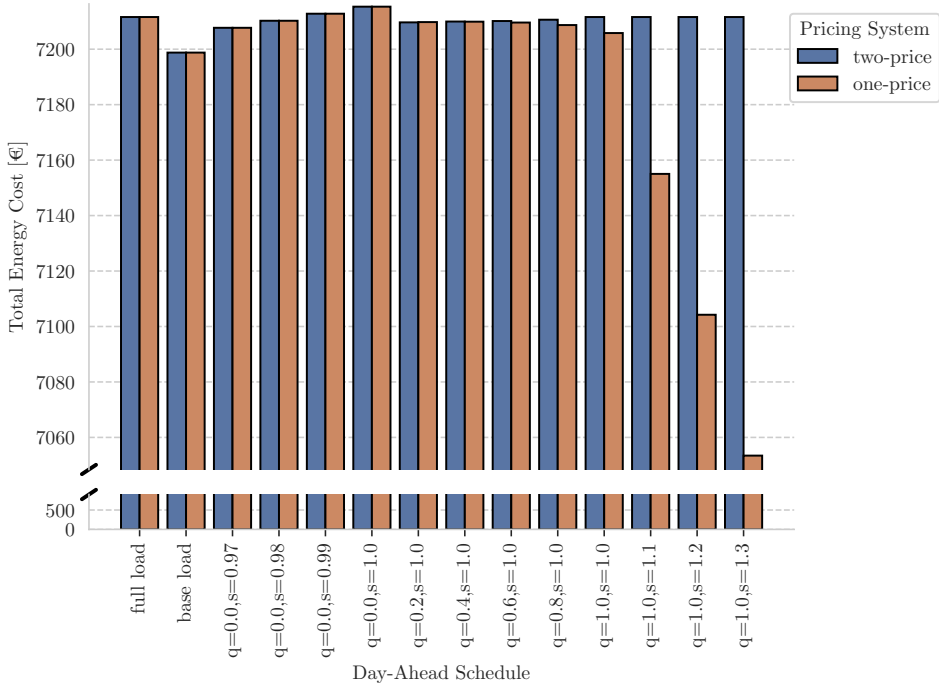


Figure 5.9: Comparison of the two imbalance pricing system.

$$\sum_i^{N_{BSP}} (Q_{\downarrow i} \cdot p_{\downarrow}^B) + \sum_j^{N_{BRP}} (Q_{-j} \cdot p_{-}^B) = \sum_i^{N_{BSP}} (Q_{\uparrow i} \cdot p_{\uparrow}^B) + \sum_j^{N_{BRP}} (Q_{+j} \cdot p_{+}^B) \quad (5.10)$$

$$\sum_i^{N_{BSP}} (Q_{\downarrow i} \cdot p_{\downarrow}^B) + \sum_j^{N_{BRP}} (Q_{-j} \cdot p_{-}^B) = \sum_i^{N_{BSP}} (Q_{\uparrow i} \cdot p_{\uparrow}^B) + \sum_j^{N_{BRP}} (Q_{+j} \cdot p_{+}^S) \quad (5.11)$$

Figure 5.9 demonstrates the difference between the two pricing systems by comparing their resulting total energy costs. As expected, datacenters cannot gain extra profit by scheduling much larger quantities of energy than actually needed under the two-price system, but they can in the case of the one-price system, which will be prohibited by the system operator. Therefore, datacenters are expected to comply with the two-price balancing system, and *do not* intentionally introduce

imbalance to the power grid. In turn, the fourth assumption (**A4**) is valid.

Having clarified all assumptions (**A1** – **A4**), we now conduct the final comparison for the energy costs of the simulated schedules (Figure 5.7). To this end, we stack the energy costs of the two markets in Figure 5.10, where the positive imbalance costs are added on top of the day-ahead costs, whilst the negative ones intrude into the day-ahead costs from the top, forming the overlaps. As highlighted in red, the variation between the total energy costs of the simulated schedules is about 0.2% with the minimum being the base-load strategy. In other words, although the margin is not significant, only scheduling the bare-minimum energy in the day-ahead market is the preferred strategy. Also, as more and more energy is scheduled in the day-ahead market ($s > 1.0$), the total energy cost demonstrates a small yet gradual and steady increase. Note that this conclusion also applies to the new machine model because it is not subjected to load variation due to the assumption of having a perfect load forecast (**A2**).

5.2.5 Relationship between Prices of the Two Markets

In previous sections, we answered the question of why and how to participate in the energy market (**RQ2**, **RQ3**). In this section, we seek answers to the question of why ML methods can be of help for datacenters in terms of leveraging the profit when taking part in the energy market. To this end, we search for potential correlations between energy prices in order to ascertain whether it is reasonable to make decisions based on simple heuristics.

Firstly, we demonstrate the correlation between the imbalance prices (shortage and surplus prices) with corresponding regulation states in Figure 5.11 (the latest, detailed descriptions of the regulation states in the Netherlands can be found in [157]). The probability distribution functions (PDFs) shown at the top and the right part illustrate that a substantial number of values concentrate between -200 and 200 . Furthermore, from the second-order interpolation between the two prices, the red curve, we can clearly see a strong linear correlation between the two imbalance prices. Also, when the regulation state is 2, surplus price is generally higher than shortage price.

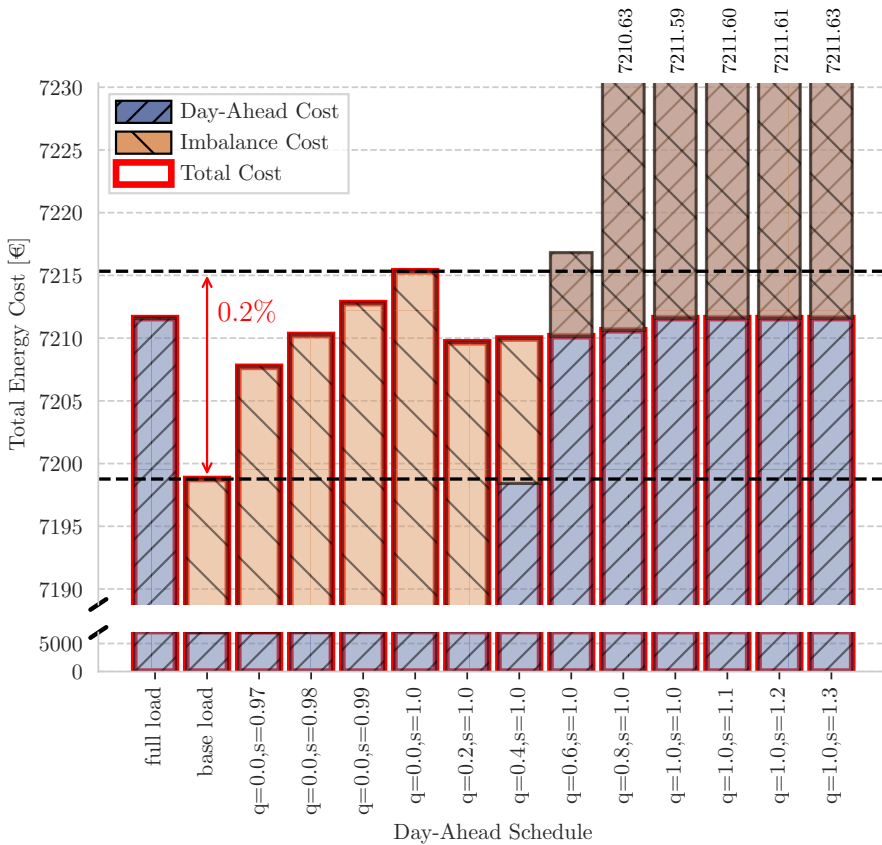


Figure 5.10: Stacked comparison of the two imbalance pricing systems.

Conversely, when it comes to the inter-market relationship, there is little if any correlation between the prices of the two markets. Figure 5.12 shows the Pearson Correlation coefficients (PCC) of the prices. The PCC between the spot price and the imbalance prices is as little as ~ 0.19 , which would provide datacenters with little to no help in making heuristics for leveraging profits. Thus, we conclude that it is not feasible to optimize profits when juggling the two energy markets by making simple heuristics.

Furthermore, currently, our ML inferences of the imbalance prices can only be obtained during the balancing market, which will be described in detail in Section

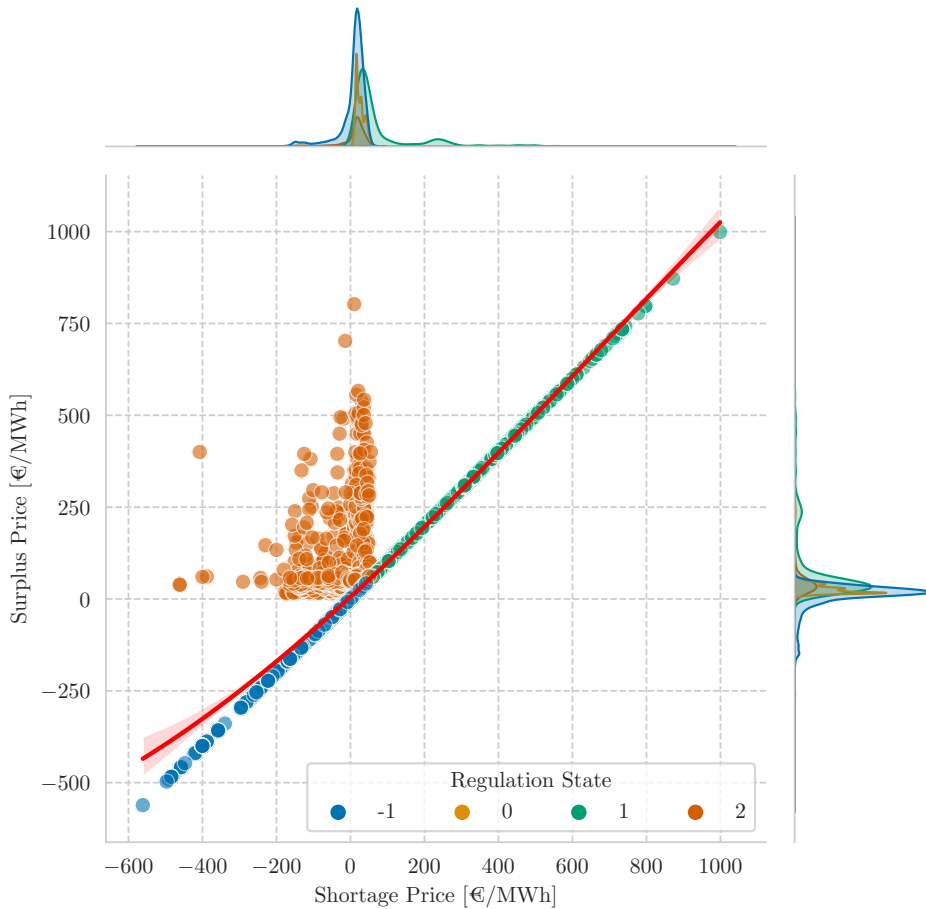


Figure 5.11: Correlation between imbalance prices with second order interpolation.

5.3. Nevertheless, to illustrate the promising potential of employing ML methods, we assume for a moment that (1) the predictions can be obtained during the day-ahead market, and (2) the predictions are perfect. Then, one straightforward procurement strategy could be: purchasing the (forecasted) full load in the day-ahead market if $p^S < p_-^B$, otherwise, using the base-load strategy, *i.e.*, scheduling the base load in the day-ahead market and settling the peak load in the balancing market. This straightforward policy together with the assumption of having the perfect prediction will provide an upper bound for the potential benefit of using the

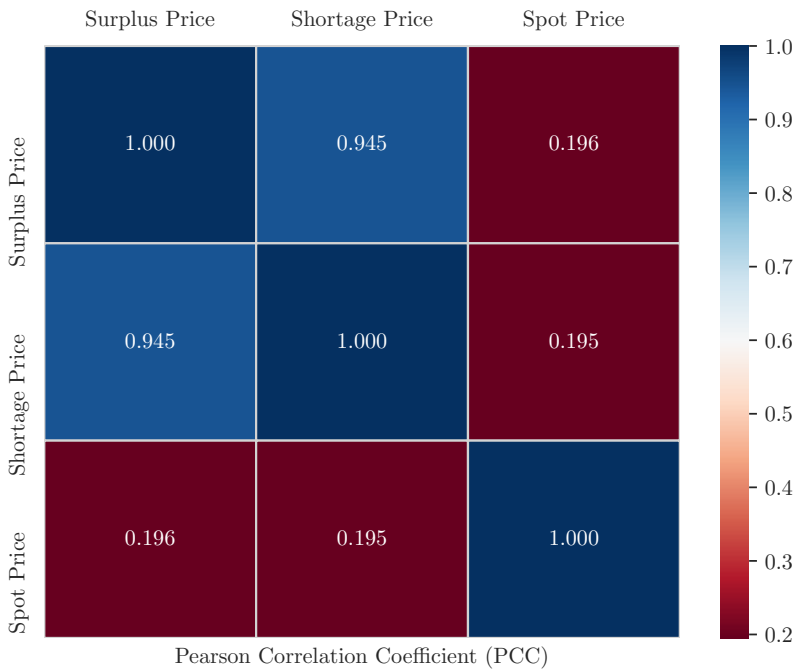


Figure 5.12: Pearson Correlation Coefficients between day-ahead prices and imbalance prices.

early obtained ML predictions for energy procurement. As shown in Figure 5.13, if such an early forecast can be achieved, we can dramatically reduce the energy cost further by up to 51.12% compared to the preferred base-load strategy by using the above mentioned straightforward procurement policy.

5.2.6 Summary

To summarize, in Section 5.2.2 we shed light on the coexistence of high risks and high profitability when taking part in the energy market. Then, for answering **RQ2** we ascertain the high financial incentive for datacenters to take part in both day-ahead and the balancing market. Next, we demonstrate the higher energy efficiency and the large variation of the power load featured by the new machine model (§5.2.3), which is important for later generalizing the experimental results.

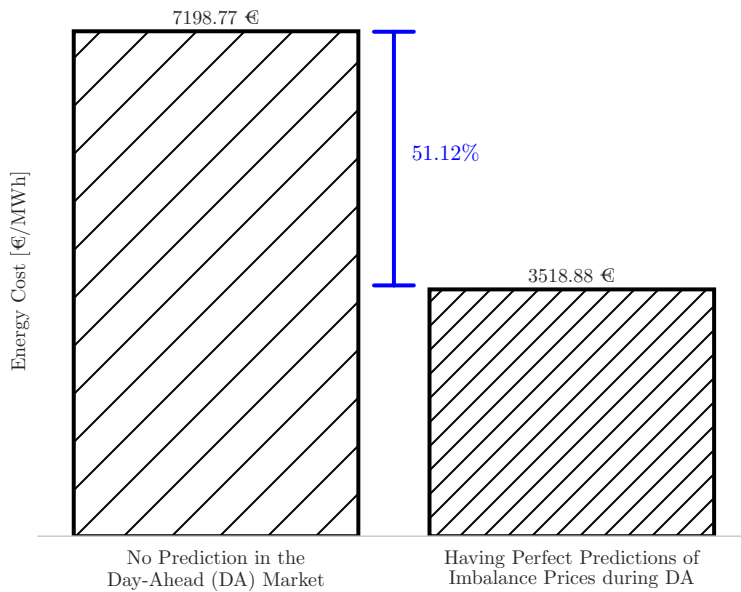


Figure 5.13: Potential benefit of obtaining predictions for imbalance prices during the day-ahead market.

Next, in Section 5.2.4 we illustrate proved that scheduling only the bare-minimum energy, the base load, is the preferred procurement strategy (**RQ3**). In Section 5.2.5, we show that it is infeasible to leverage market participation simply by making heuristics and that substantial profit could be derived by employing early ML predictions.

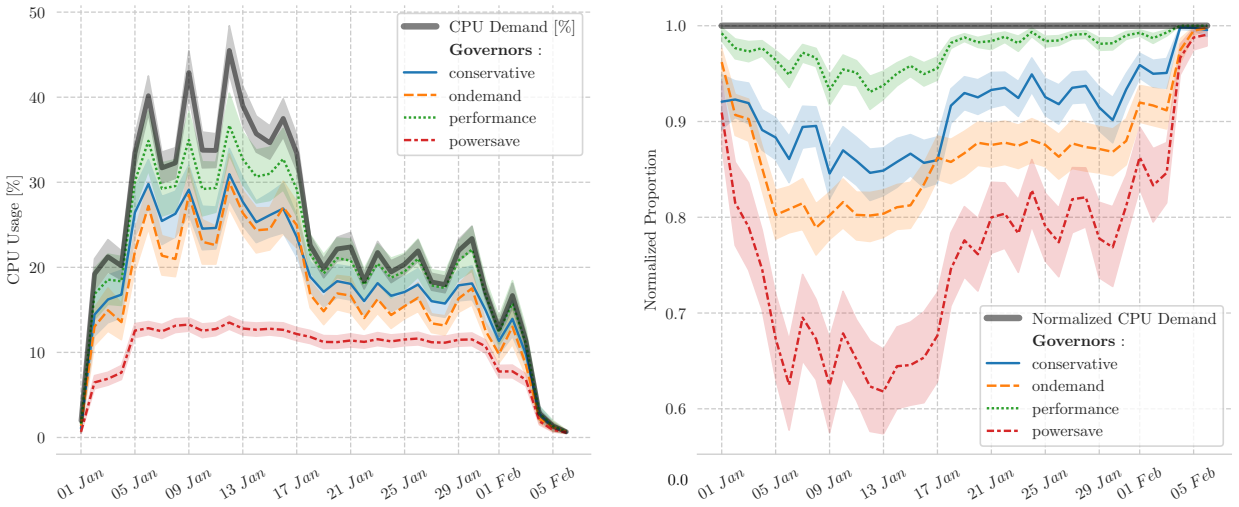


Figure 5.14: Relationships between CPU demand and CPU usage of the four scaling governors.

5.3 DVFS Scheduling

In this section, we present the results regarding the performance and effect of the developed proactive DVFS scheduler, in which ML inferences are employed, to answer **RQ4**. Firstly, we demonstrate the behaviours of the DVFS implemented in our system (§5.3.1). Then, in Section 5.3.2 we investigate the effect of the hyper-parameter of the scheduler, the damping factor. After that, we conduct bounded performance estimation, comparing the ML methods with synthetic estimators in Section 5.3.3.

5.3.1 DVFS Behaviour

First of all, we demonstrate the behaviours of the basic DVFS mechanism described in Section 4.1. The left plot of Figure 5.14 shows the relationships between the CPU demand and the actual CPU usage of the four scaling governors, namely, the conservative, the ondemand, the powersave, and the performance governors. As

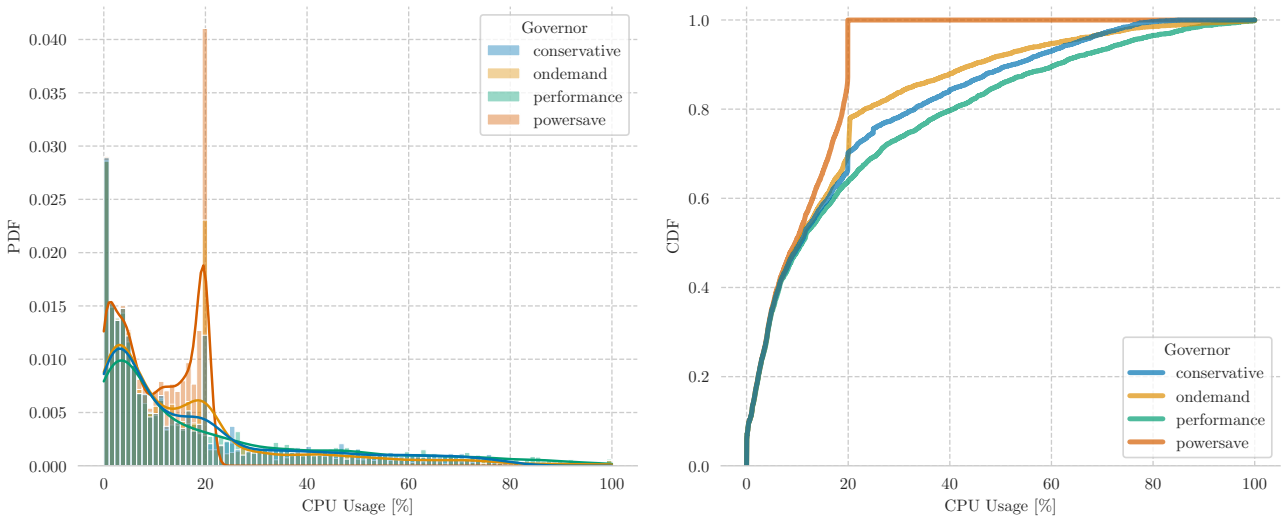


Figure 5.15: PDF and CDF of the CPU usage of different scaling governors.

shown by the dotted green curve, the **performance** governor tries its best to meet the required computation power, whilst the CPU usage of the **powersave** is capped under 20% because of the frequency limit imposed on the CPU. The figure on the right offers a scaled view of that of the right, where the CPU usage is normalized by the CPU demand. The distribution of their CPU usage is further demonstrated by Figure 5.15, where the corresponding PDFs and CDFs are presented. As shown in Figure 5.16, other than the **powersave** governor, the other three governors have about the same median value of CPU usage. Also, the **ondemand** governor exhibits more concentrated CPU usage compared to that of the **conservative**. Conversely, the CPU usage of the **conservative** has never reached 100% since it proposes frequency changes in consecutive small steps as described in Section 2.3 and 4.1. With regard to the instant power draw shown in Figure 5.18, the **conservative** governor exhibits large fluctuations, again because of the gradual adjustments in CPU frequency. In contrast, the power consumption of the **ondemand** governor forms narrow lanes due to its drastic changes in CPU frequency. As suggested in Figure 5.18, the behaviours of their instant power draw exhibit similar patterns to that of the CPU usage shown in Figure 5.14.

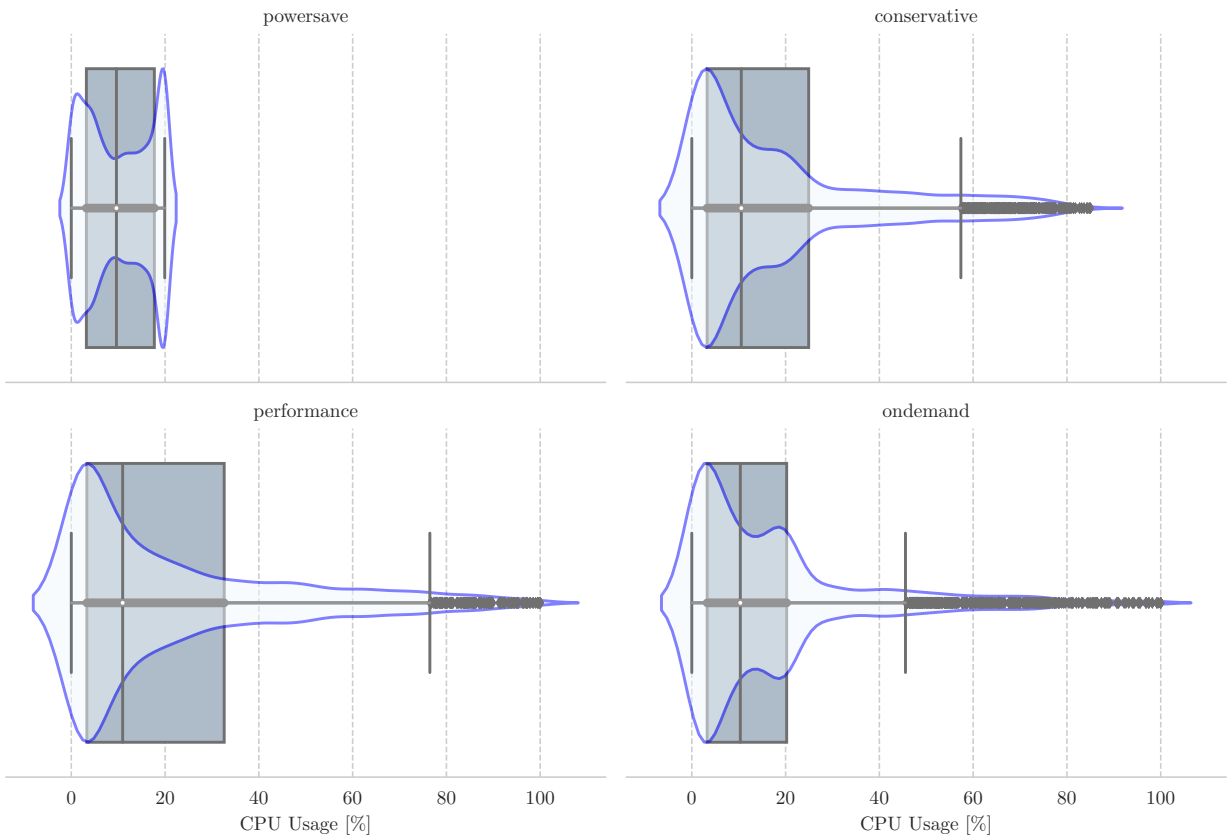


Figure 5.16: Detailed distributions of CPU usage of different scaling governors.

When it comes to CPU over-commission, however, the behaviour pattern is exactly the opposite of that of the CPU usage. As shown in Figure 5.19, the `powersave` governor has a much higher CPU over-commission compared to the others. In contrast, the `performance` governor has the lowest level since it does not stress the CPU speed whatsoever. Moreover, to further understand the behaviour of CPU over-commission, we examine its relationship with other factors.

Firstly, Figure 5.20 shows the relationship between CPU usage and over-commissioned CPU cycles. The `performance` governor barely exhibits any trace of over-commission across different levels of CPU usage, whilst other governors, especially, the `powersave` governor, demonstrate significant CPU over-commission. Due to the speed

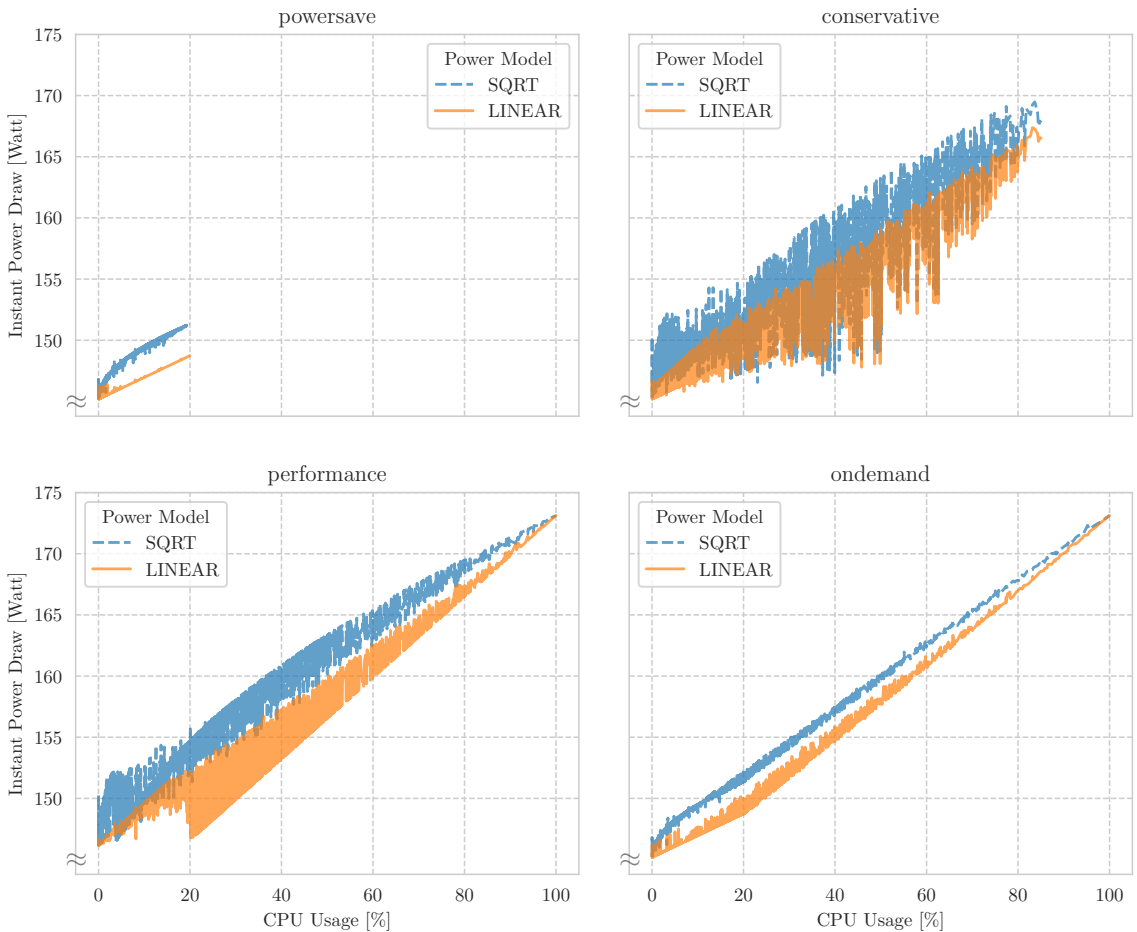


Figure 5.17: Comparison of power estimation of the two models for the four scaling governors.

cap applied, the over-commission level of the `powersave` governor reaches a dramatic level at around 20% of CPU usage. Furthermore, referring back to Figure 5.16, the `ondemand` governor has narrower interquartile range (IQR) than that of the `conservative`. Consequently, even though the overall range of CPU usage is larger in the case of `ondemand` governor than the `conservative`, `ondemand` exhibits higher levels of over-commission. Thus, the more concentrated the CPU usage (*i.e.*, the narrower the IQR), the higher level of CPU over-commission.

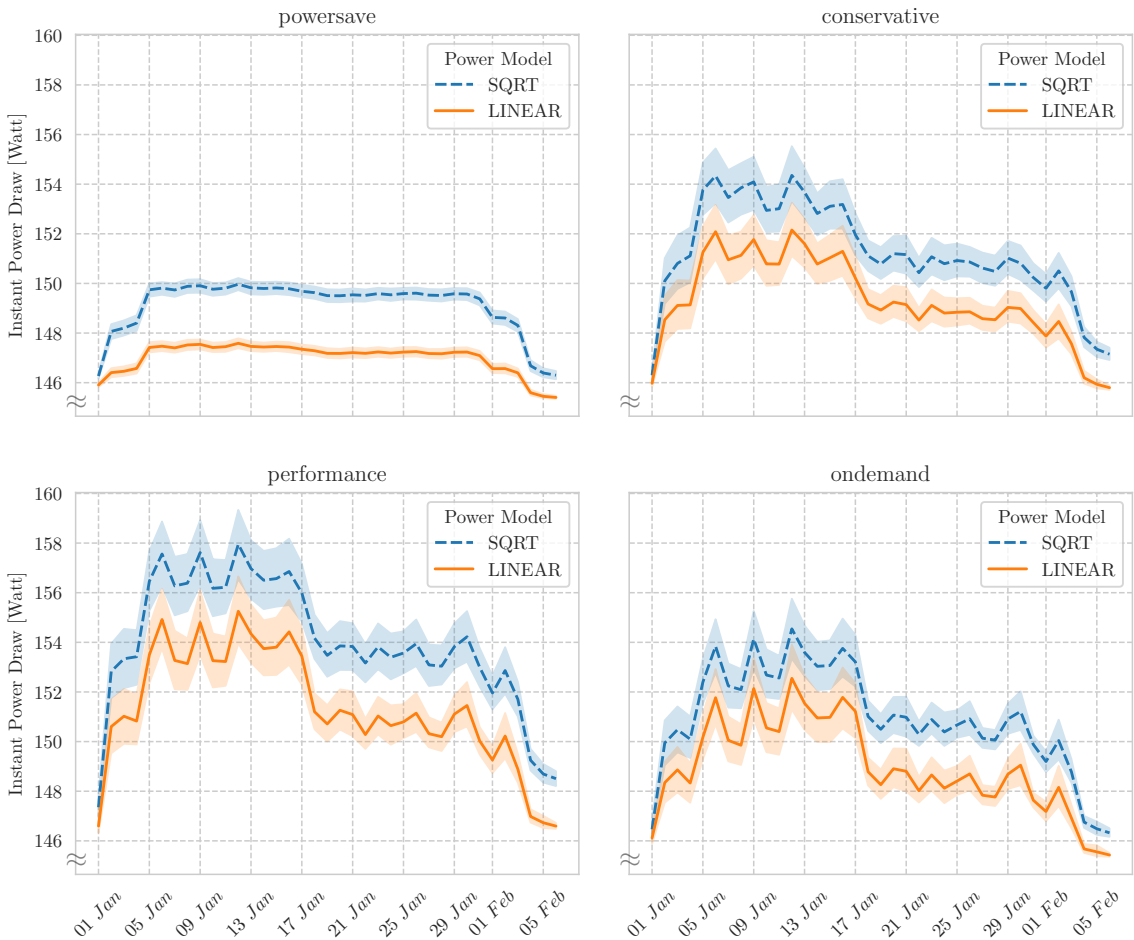


Figure 5.18: Comparison of instant power draw of different scaling governors estimated by two power models.

Secondly, Figure 5.21 illustrates the relationship between the over-commission and the instant power draw. Although the `powersave` governor has a much higher over-commission level, its power is capped because of the frequency limit. Since both the `performance` and the `ondemand` governors exploit the full range of CPU usage (0–100%), their instant power draw both reaches about 175 W. In contrast, as the CPU usage of the `conservative` is mostly below 80% (Figure 5.16), its instant power is also restrained by about 170 W. Therefore, the level of power draw is *not*

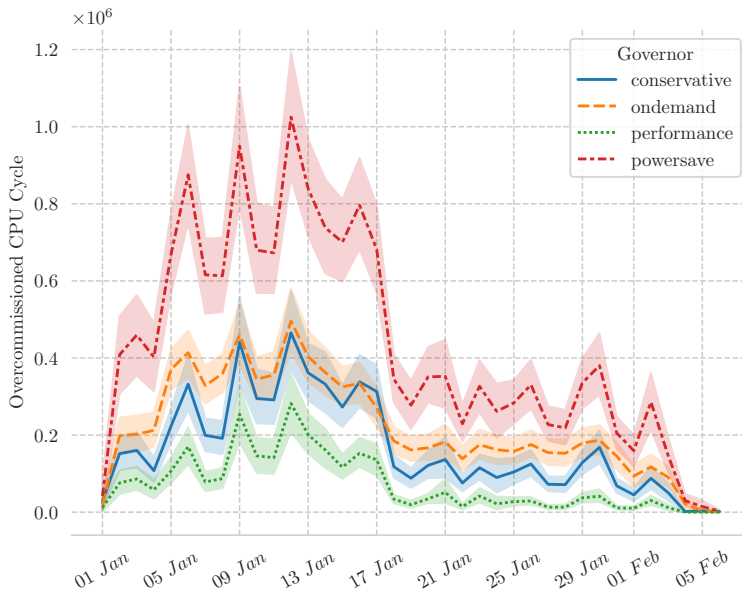


Figure 5.19: Instant CPU over-commission level of the four scaling governors.

directly associated with the CPU over-commission but depends upon the actual CPU usage. Additionally, due to the higher level of over-commission, governors other than the `performance` squeeze out the power steps caused by various discrete P-states.

Lastly, rather than looking at over-commission, we move on to the actual work committed by the CPU in Figure 5.22. Different from over-commission, there is a clear connection between the amount of committed work and the level of instant power draw. In other words, the more work the CPU commits, the higher power the system requires.

Figure 5.20: Over-commission at different CPU usage levels of the four scaling governors.

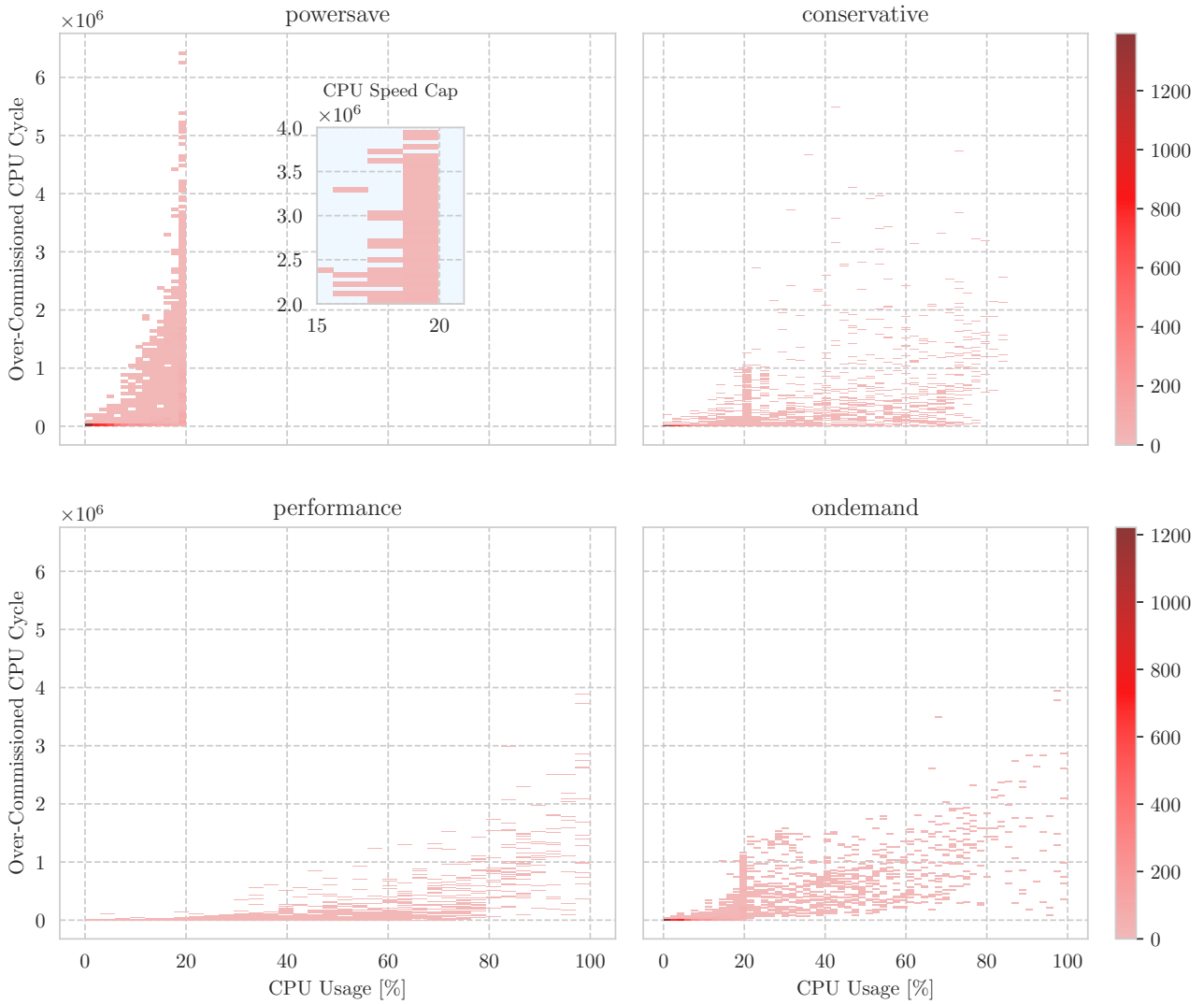


Figure 5.21: Instant power draw at different over-commission levels of the four scaling governors.

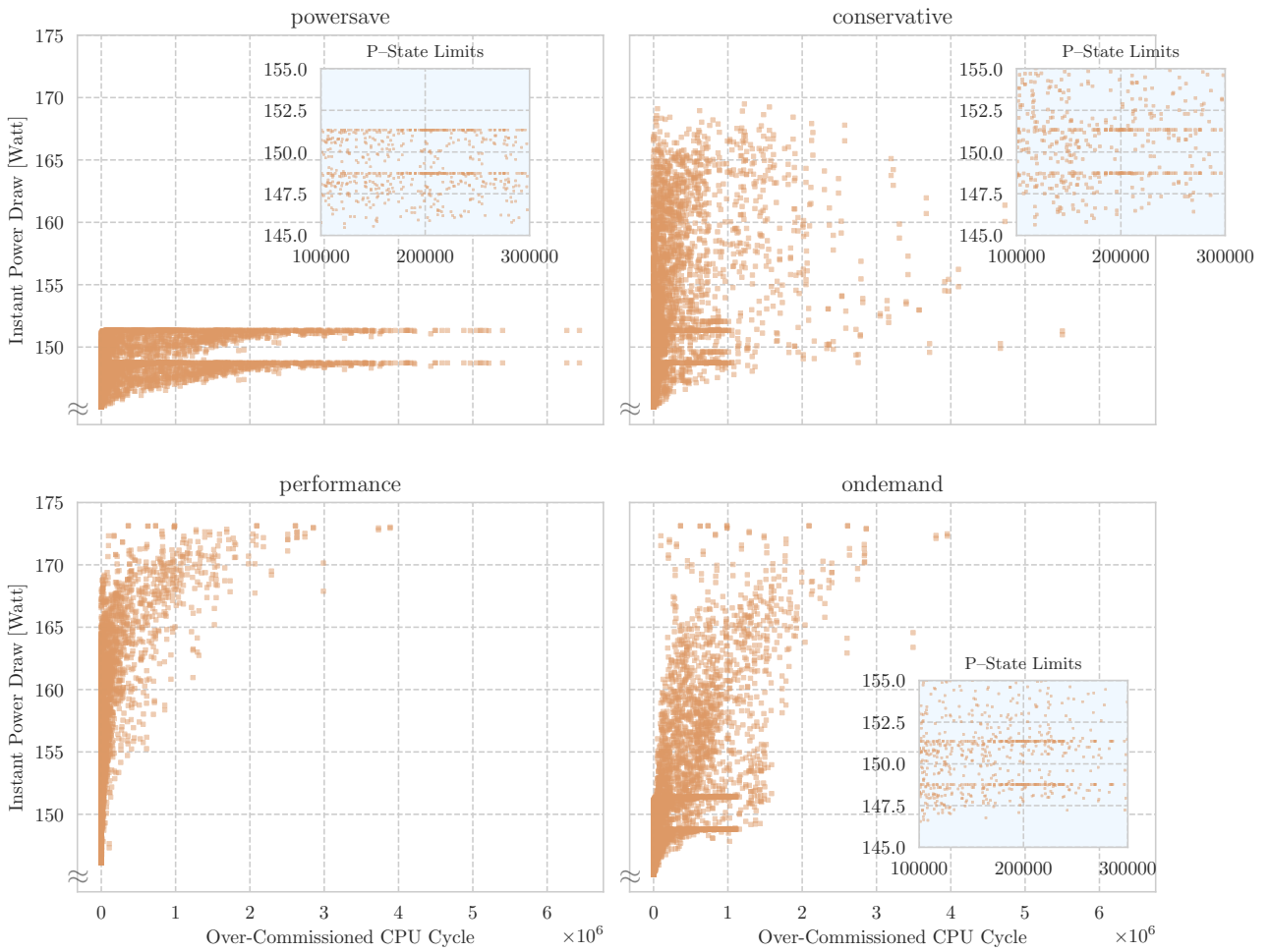
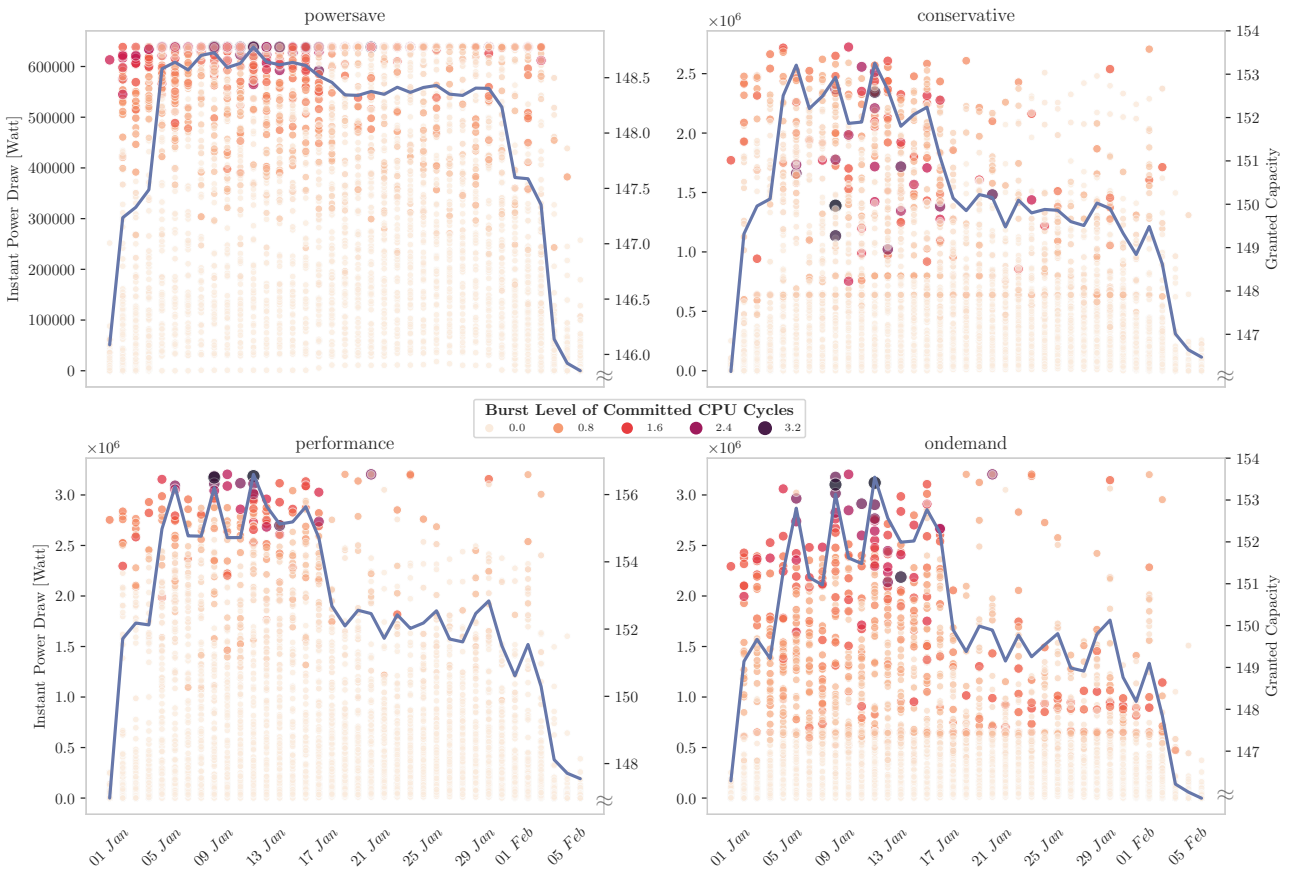


Figure 5.22: Comparison of granted work and instant power draw of the four scaling governors.



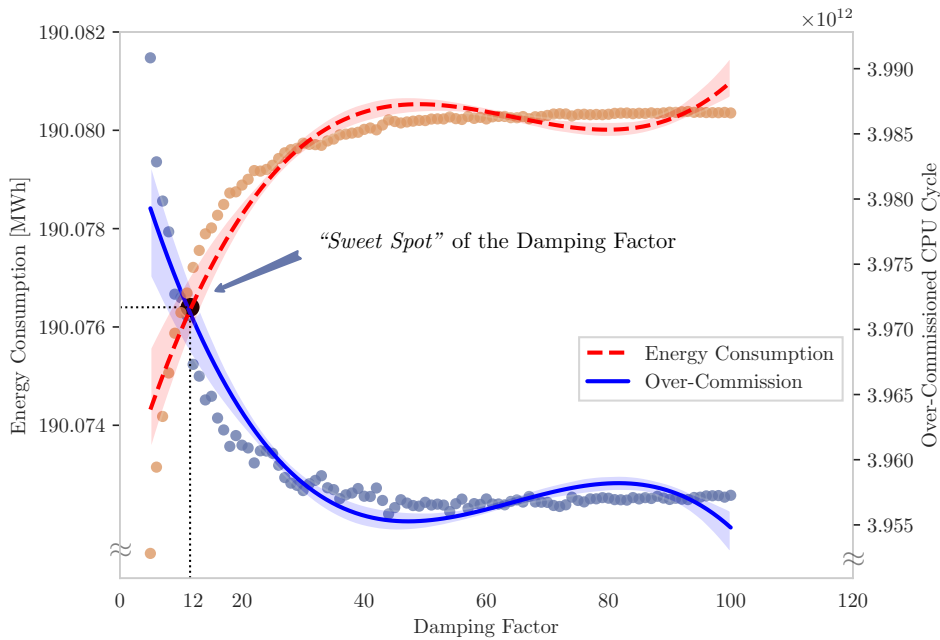


Figure 5.23: Effect of the damping factor on energy use and CPU over-commission with third order regression.

5.3.2 Damping Factor

Having studied the behaviours of the DVFS, we now focus on its proactive scheduler implemented in the EEMM extension. As described in Section 4.2, the damping factor is the hyperparameter responsible for ameliorating the stress imposed on the CPU. To this end, users can use the damping-factor plot shown in Figure 5.23 to strike a balance between the energy-saving and the level of CPU over-commission. Specifically, users determine the damping factor according to their desired level of energy saving and the acceptable level of CPU over-commission. In our case, the damping factors with which we have experimented run in the range of 0 to 110. As suggested in Figure 5.23, both the energy consumption and the level of over-commission have converged when the damping factor is greater than ~ 80 . Then, we conduct a third-order regression on both the energy consumption (the red, dotted curve) and the CPU over-commission (the blue, solid curve). Since we do not as-

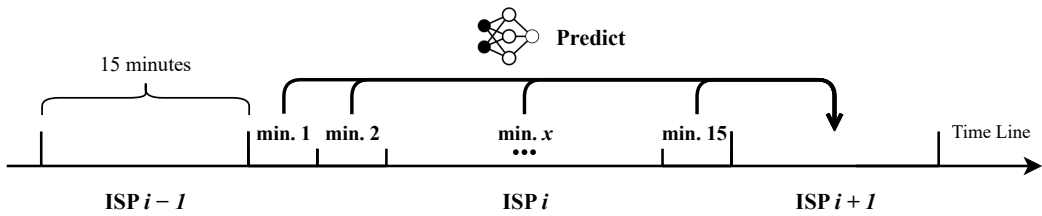


Figure 5.24: Timeline of ML prediction (abbreviation: min. = minute).

sume a specific target for either of the two, we set the damping factor to the “Sweet Spot”, the intersection between the two regression lines marked by the black dot in the figure, which is about 12. To reiterate, finding the “Sweet Spot” is not the only way, certainly not necessarily the best way, for determining the damping factor. Instead, users ought to customize it according to their needs regarding energy saving and the CPU over-commission.

5.3.3 Bounded Comparison

In this section, we describe the use of ML methods to further optimize the profit for datacenters when participating in the day-ahead and the balancing market. First and foremost, we elaborate on how the ML inferences are obtained. Referring back to Section 2.4.1, unlike the day-ahead market, trading happens in the balancing market every 15 minutes. These trading intervals are often referred to as imbalance settlement periods (ISPs) and, in turn, there are 96 consecutive ISPs per trading day. As shown in Figure 5.24, we predict the energy price of the next ISP (**ISP** $i + 1$) during the current ISP (**ISP** i) on a minute-by-minute basis. In the following sections, the predictions produced in the first minute are referred to as the first ML inferences and the last ones are referred to as the last ML inferences. Also, the average predictions of every ISP are referred to as the average ML inferences. Furthermore, it is worth noting that the closer it gets to the next ISP, the more accurate the ML inferences are, but the shorter the time that allows datacenter operators to adjust their optional schedule accordingly. Thus, there is a trade-off between the performance of the ML methods and the operational leeway that datacenter operators have.

5.3.3.1 Defining Metric

N_{ISP}	Number of ISPs
p^S	Spot price in the day-ahead market
p^B	Shortage price in the balancing market
p^F	Forecasted imbalance price

Table 5.8: Symbols used in defining the AA score.

To measure the performance of the ML methods, we define a metric base upon the needs of the DVFS scheduler. Referring back to Algorithm 5, there are two decision points at line 11 and line 15 respectively. The first decision point is to check if the imbalance price is directly profitable. If it is not, in the second decision point, we compare the price level of the balancing market with that of the day-ahead market to determine whether or not to further suppress the CPU frequency. Based upon these two decision points, we define the agreement accuracy (AA) score by Equation 5.13 as a measure of the performance of the ML methods. The \mathbb{S} function (Equation 5.12) checks the sign of the input value, and the $\mathbb{1}$ is the indicator function. Table 5.8 present the meaning of the used symbols.

$$\mathbb{S}(x) = \begin{cases} +1 & x > 0 \\ 0 & x = 0 \\ -1 & x < 0 \end{cases} \quad (5.12)$$

$$AA = \frac{\sum_i^{N_{\text{ISP}}} \mathbb{1} \left\{ \mathbb{1} \left[\mathbb{S}(p_i^B) = \mathbb{S}(p_i^F) \right] = \mathbb{1} \left[\mathbb{S}(p_i^B - p_i^S) = \mathbb{S}(p_i^F - p_i^S) \right] \right\}}{N_{\text{ISP}}} \quad (5.13)$$

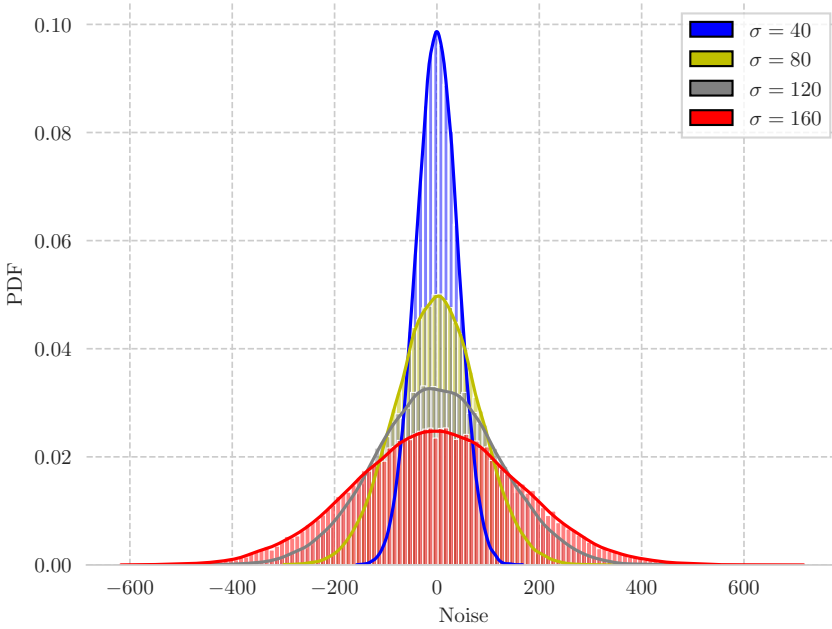


Figure 5.25: Distributions of Gaussian noises with different σ values.

5.3.3.2 Synthetic Predictor

To achieve bounded evaluation (**NFR3**), we first construct a set of synthetic predictors by adding Gaussian noises to the *actual* imbalance prices as follows:

$$p^f = p^B + E, \quad E \sim \mathcal{N}(0, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{x}{\sigma}\right)^2}, \quad (5.14)$$

where E is an additive random variable representing the error term that follows the Gaussian distribution $\mathcal{N}(0, \sigma)$. Referring back to Figure 2.6, the imbalance prices are capped by 800 €. Thus, we set $\sigma \in [0, 1200]$ in the following experiments to ensure that variation of errors is sufficiently large.

Figure 5.26 shows the trending of the AA scores of the ML methods in percentage and the synthetic predictors as the σ value increases. As previously described, the later the ML inferences are produced, the better their performance. Indeed, the

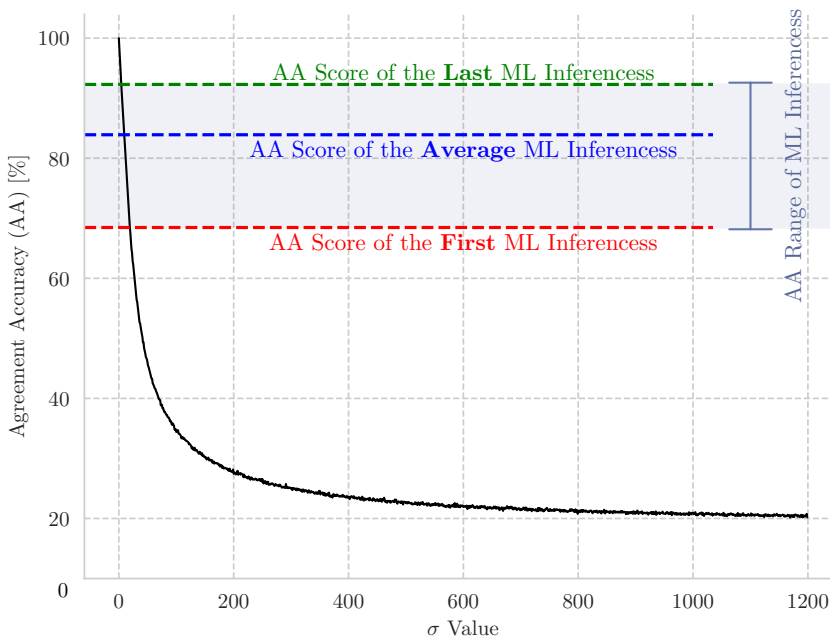


Figure 5.26: Comparison of AA scores of ML methods and synthetic predictors at different σ levels.

first ML inferences achieve the best AA score, which is about 0.92. The AA score of average ML inferences comes next (0.84), followed by the score of the first ML inferences, which is about 0.68. In respect of the synthetic predictors, when $\sigma = 0$, the AA score is a perfect 1.0. It exhibits a plunge when $\sigma \in [0, 200]$. Then, as σ continues to decrease, the AA scores of the synthetic predictors level off and converge to 0.20 at last.

Next, we turn our attention to energy consumption and the level of CPU over-commission. Figure 5.27 shows the comparison of energy consumption and over-commission between ML methods and synthetic predictors at different σ levels. Regarding the ML methods, the last ML inferences optimize energy consumption over over-commission, which, in turn, leads to the lowest energy use but also the highest over-commission level amongst the ML methods. On the contrary, the first ML inferences leverage over-commission level over energy consumption, which results in the highest energy consumption and the lowest over-commission level.

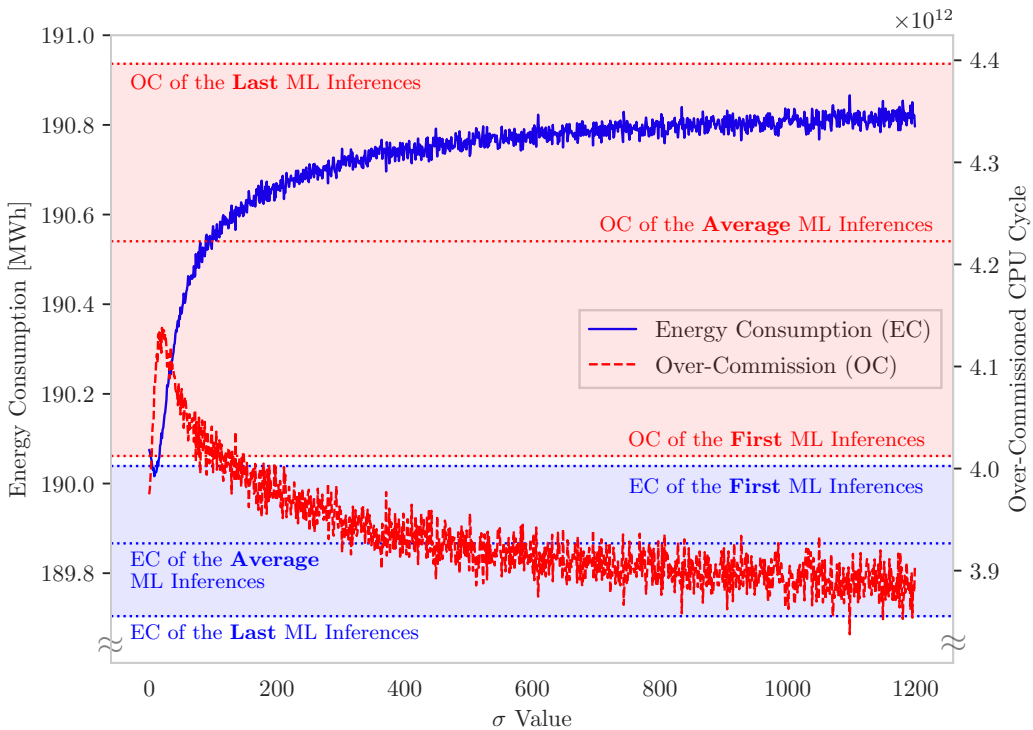


Figure 5.27: Comparison of energy consumption and over-commission between ML methods and synthetic predictors at different σ levels.

The ML method of the average inferences achieves average values in both metrics compared to the ML methods that use the first and the last inferences. Regarding the synthetic predictors, almost all of them, even the best synthetic predictor ($\sigma = 0$), is not able to achieve the same level of energy-saving as that of the worst performed ML methods, the one using the first ML inferences. When it comes to the level of over-commission, however, these synthetic predictors perform particularly well, in the sense that when $\sigma > \sim 100$, none of the ML methods can reduce the over-commission to the same level as that of the synthetic predictors. Hence, we conclude that the scheduler powered by the ML methods is good at reducing the energy consumption but perform poorly in reducing the over-commission level, compared to using the synthetic predictors.

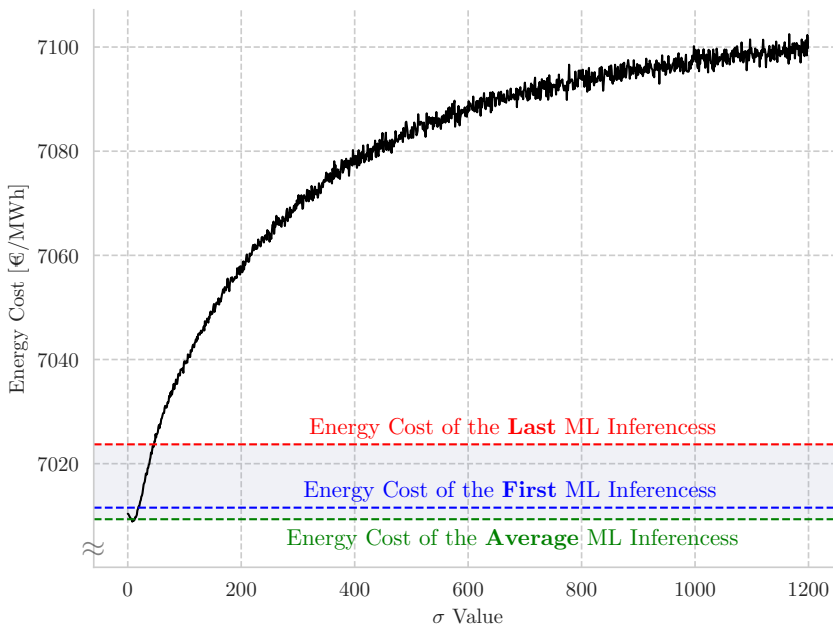


Figure 5.28: Comparison of energy costs between ML methods and synthetic predictors at different σ levels.

Now, we focus on the total energy cost. Figure 5.28 demonstrates the comparison of energy costs between ML methods and synthetic predictors at different σ levels. Overall, compared to the synthetic predictors, the ML methods are excellent at leveraging the energy cost because even the best synthetic predictor ($\sigma = 0$) is bounded by the best ML method. Furthermore, what counts most is that the best-performed ML method in saving cost is *not* the one that saves the most energy, the one using the last ML inferences, but instead, the ML method that employs the average values achieves the lowest energy cost. This observation further hammers home the paramount importance of consuming the right amount of energy at the right time.

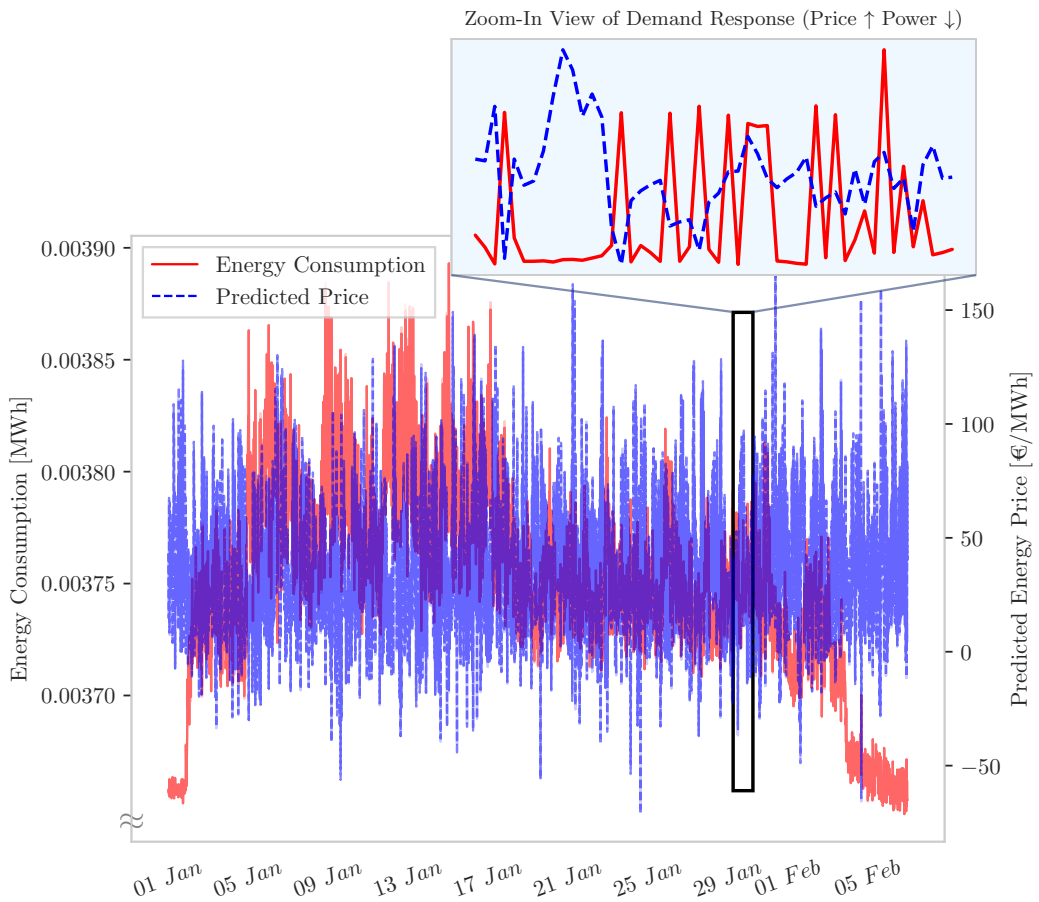


Figure 5.29: Demonstration of indirect demand response, where the energy-consumption level is adjusted in response to the predicted energy cost from a synthetic predictor of $\sigma = 50$.

5.3.3.3 Indirect Demand Response

Last, but not certainly not least, we look into the indirect DR resulted from the proactive DVFS scheduling. In Figure 5.29, we visualize the energy consumption using the solid, red line together with the predicted energy price produced by a synthetic predictor of arbitrary $\sigma = 50$ using a blue, dotted line. As we zoom in to a short timeframe, we can see the gaps between the two lines: when the energy is

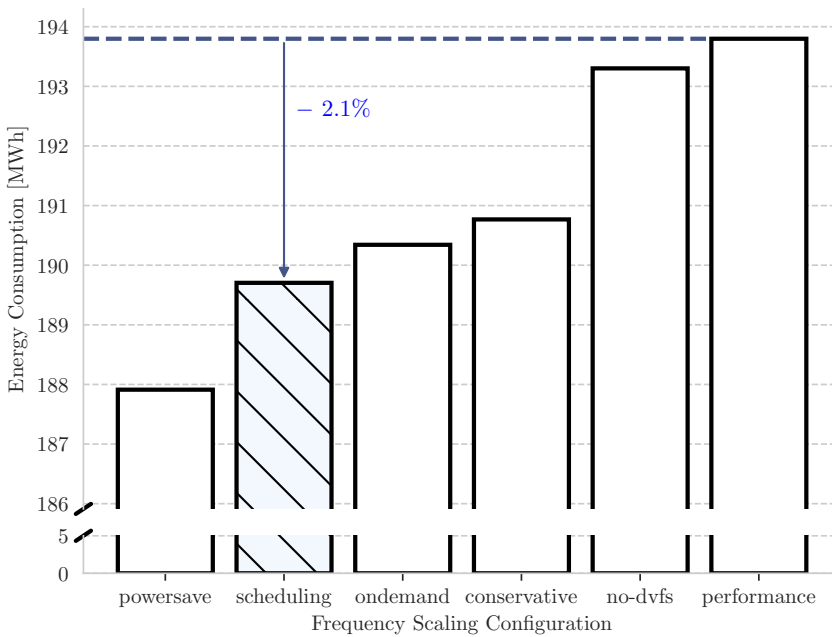


Figure 5.30: Comparison of total energy consumption for DVFS scheduling.

high, the scheduler tries to reduce the amount of energy consumed and vice versa. Note that, because of the actual demand and the effect of the damping factor, such gaps may not always be as obvious.

5.3.4 Summary

In this section, we demonstrate the behaviours of the implemented DVFS in Section 5.3.1 with a focus on CPU over-commission. Then, we determine the hyperparameter, the damping factor, of the scheduler using the damping factor plot (§5.3.2). After that, by conducting bounded comparisons between the ML methods and the synthetic estimators, we show that the ML methods are good at optimizing the energy consumption and the cost, whereas they are not as performant in curbing the over-commission level. Lastly, we visualize the energy consumption level of the resulting DVFS schedule together with a randomly chosen synthetic predictor,

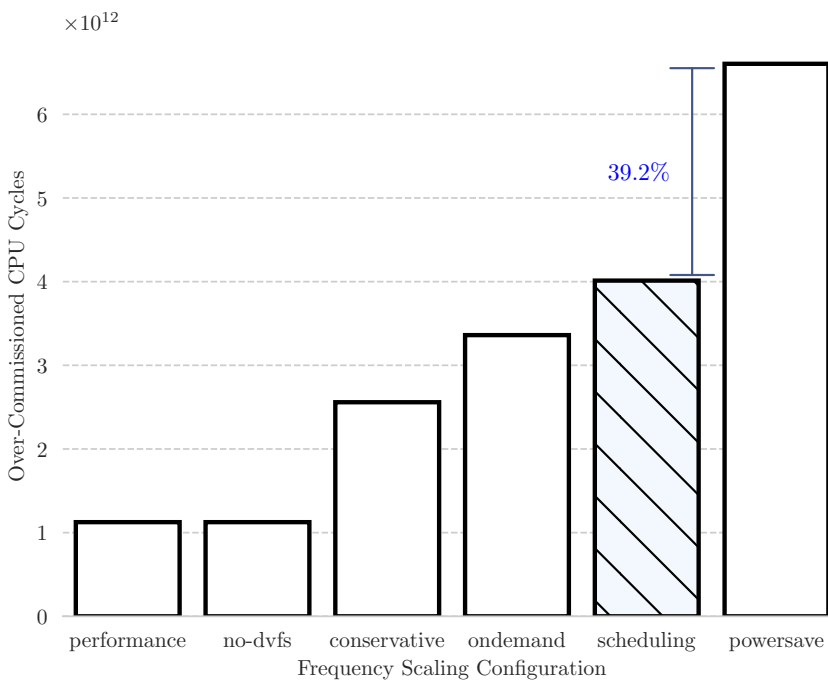


Figure 5.31: Comparison of total CPU over-commission for DVFS scheduling.

illustrating the indirect DR.

As shown in Figure 5.30, if we take the best results amongst the ML methods, we reduce about 2.1% energy consumption compared to that of the **performance** governor. Such an energy saving on this single workload is roughly the equivalent of the annual energy consumption of two Dutch households [165]. With regard to CPU over-commission, DVFS scheduling powered by the ML methods results in about 39% improvement compared to that of the **powersave** governor (Figure 5.31).

Finally, regarding the energy cost, together with the base-load procurement strategy (*i.e.*, schedule the bare-minimum energy in the day-ahead market and settle the peak load in the balancing market), the proactive DVFS scheduler powered by ML methods achieves a 2.6% of reduction in energy cost compared to the case where it is disabled. Moreover, the cost reduction is about 56.4% compared to the

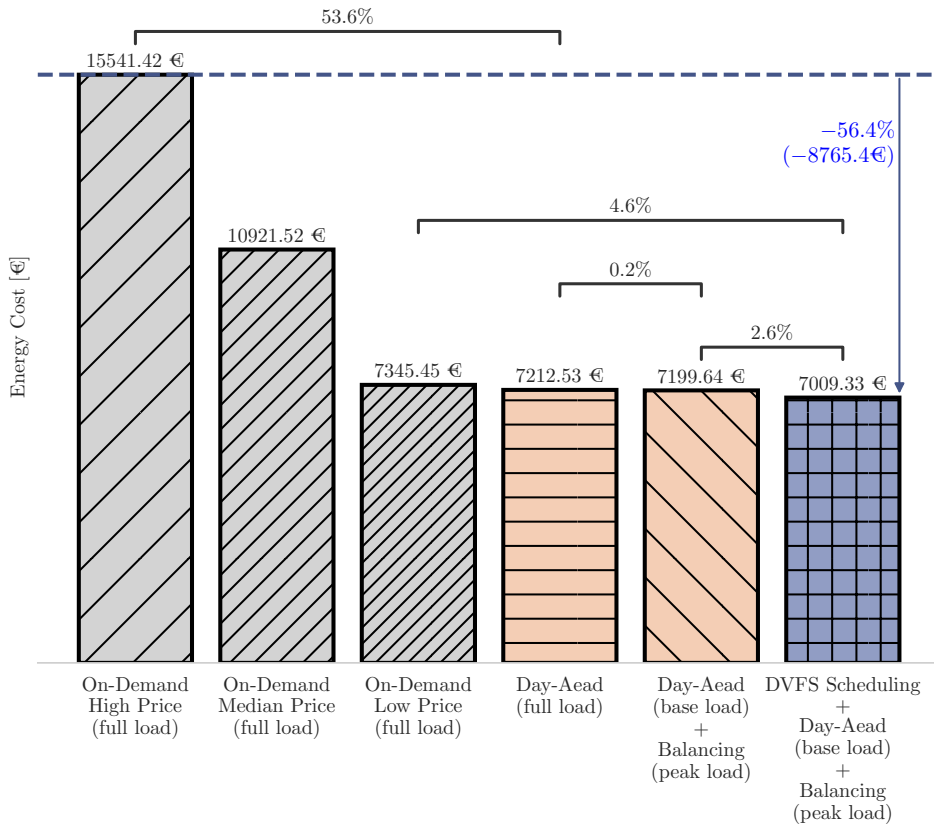


Figure 5.32: Comparison of total energy costs for DVFS scheduling.

on-demand scheme of high price.

In fact, we believe the results are rather *conservative* because they are produced by the *old* machine model whose peak load only takes up a tiny portion of its overall power load (Figure 5.6). Furthermore, the variations in power load between energy states of the old machine model are much smaller than that of the new machine model, as illustrated in Figure 5.5. In other words, the fraction of energy cost that can be leveraged is minuscule (5.3). Therefore, the advantage of the proactive DVFS scheduling has not been fully exploited. Hence, if a newer machine model were available, the above results would be expected to be much more significant.

Conclusion

In this last chapter, we first summarize all conclusions drawn from previous chapters to provide answers to research questions **RQ1** to **RQ5** in Section 6.1. Then, we reflect on the limitations of this work (§6.2) as well as lay out potential directions of future work (§6.3).

6.1 Answers to Research Questions

Answering RQ1. In this work, we extend the OpenDC simulator with advanced energy models so that it is able to model the whole power system of a typical datacenter in a flexible and highly customizable way. In Section 3.3, we present the design of the energy management and modelling system, and in Section 4.1, we describe the detailed implementations of the system and its subsystems.

Answering RQ2. To answer this and the following **RQs**, we first develop a market extension of the energy management and modelling system. Then, by using this tool, we show that there is a strong financial incentive for datacenters to participate in both the day-ahead market and the balancing market, as substantial profit can be derived by doing so (§5.2.2).

Answering RQ3. By simulating different load-forecast-based procurement strategies in Section 5.2.4, we demonstrate that the most economical strategy is the base-load strategy, *i.e.*, scheduling only the bare-minimum amount of energy, the base load, in the day-ahead market and resolve the peak load in the balancing market.

Answering RQ4. Firstly, we develop a proactive DVFS scheduler powered by ML methods that can save energy whilst curbing the level of CPU over-commission (§4.2). In Section 5.3, we employ the DVFS scheduler to make fine-grained decisions in order to leverage the profit when participating in the energy market. Further, we carry out bounded comparisons between the ML methods and synthetic predictors (§5.3.3.2), showing that the ML methods are excellent at reducing both the energy consumption and the cost, whilst not as performant in curbing the CPU over-commission. Lastly, in Section 5.3.4, we demonstrate that the DVFS scheduler is able to save about 2.1% of energy compared to the performance governor and to improve about 39.2% of CPU over-commission compared to the powersave governor. Such an energy saving on this single workload is roughly the equivalent of the annual energy consumption of two Dutch households (§5.3.4). Moreover, together with the base-load procurement strategy, the scheduler is able to save about 56.4 % energy cost compared to the on-demand scheme of high price. Furthermore, we conclude in Section 5.3.4 that these results are *conservative* and are expected to be more significant on newer machine models.

Answering RQ5. To meet the requirements of this RQ, we follow a carefully designed development pipeline from the onset (§3.1) and conduct rigorous requirement engineering in Section 3.2. The software artefacts are created and maintained using strict software engineering methods. As a result, both the datacenter simulator and its extension EEMM can be run with a few clicks and/or simple shell commands, without requiring excessive prerequisite knowledge or any manual data preprocessing from users.

6.2 Limitations

Whilst acknowledging the promising results, we recognize the limitations thereof as well. In this section, we elaborate on the limitations and their corresponding mitigations.

Internal Limitation. In our experiments, we do not assume any specific machine types or computing platforms except for the new machine model. Consequently, the estimation of energy consumption may not be particularly representative. To mitigate this limitation, we use the knowledge gained from our previous studies, employing the linear model and the square-root model as the lower and upper bound respectively. In this way, the variation of the total energy consumption is bounded to a certain range.

Construct Limitation. As introduced in Section 2.4, the energy demand is inelastic. Thus, the market prices, albeit a good indicator and predictor, may not be able to satisfactorily capture the demand-supply balance in the power grid at all times. Nevertheless, as introduced in Chapter 1, this is the inherent disadvantage of the indirect DR approach. As the functions of the smart grid and the design of the energy market advance, we believe the energy prices will become increasingly responsive.

External Limitations. Although we strive to make our scientific tools more user-friendly and easy to use, users can still use them incorrectly and, in turn, produce invalid results. For example, the market data from ENTSO-E is in CET, whilst data from TenT is in GMT, so users may well feed the software with inputs that have mismatched timestamps. To mitigate this limitation, we write detailed documentation and make tutorials with examples ^{1 2}.

¹<http://opendc-eemm.rtfid.io>

²<https://github.com/atlarge-research/opendc#documentation>

6.3 Future Work

According to our industry partners, energy planning in datacenters using simulation and modelling has been quickly gaining popularity as higher degrees of volatility caused by renewable energy sources are introduced to the market. Yet, the intersection between the energy market and datacenter simulation is still far from been fully explored. In this section, we identify potential directions of future research by means of questions; the following is a non-exhaustive list of such questions in no particular order.

1. How feasible and beneficial is it for individual datacenters to serve as BSPs instead of BRPs?
2. How can we develop a convenient and liable tool to measure P-state consumption levels, which will enable flexible experiments on more machine models?
3. What is the impact of employing not only the energy price but also the frequency level of the power grid (with less frequent grid monitoring and communication) in datacenters' decision-making on their market participation?
4. How to improve the algorithms of resource allocation (*e.g.*, power distribution, VM scaling/placement, etc.) in response to market signals?
5. What is the effect of core-level P-state frequency scaling on datacenters' participation in DR programmes?
6. How can we improve the design of the energy market to incentivize datacenters' active participation?
7. How can we optimize the direct participation in multiple markets across several, geographically distributed datacenters?
8. How to orchestrate redundancies of datacenters (*e.g.*, PSU, UPS, etc.) to provide indirect DR?
9. How can we tune the proactive scheduler base upon specifications in the SLA?

6.4 Summary

In this work, we first model the entire power system of a quintessential datacenter. By conducting simulation on real-world traces, we then demonstrate the substantial financial incentive for individual datacenters to directly participate in the energy market, specifically, the day-ahead and the balancing markets. In turn, we suggest a new short-term, direct scheme of market participation for individual datacenters in place of the current long-term, inactive participation. Furthermore, we develop a novel proactive DVFS scheduling algorithm that is able to both reduce energy consumption and save the energy cost for datacenters when participating in the energy market. Also, in developing this scheduler, we propose an innovative combination of ML methods and the DVFS technology that is able to provide the power grid with indirect DR in an effort to combat the challenges brought by renewable energy sources.

Besides the aforementioned potential societal and economic impacts, we develop and open source scientific tools, bridging the gap between domain knowledge, and sharing our data and experimental results with the community without reservation. Last, but certainly not least, with all these efforts, we believe that we have opened a new research line; as such, we call for collaborations from both the industry and academia to help datacenters actively join the smart grid to tackle the climate crisis together.

References

- [1] **ACPICA: ACPI Component Architecture.** Available at <https://www.acpica.org/>. 22
- [2] H Aalami, GR Yousefi, and M Parsa Moghadam. **Demand response model considering EDRP and TOU programs.** In *2008 IEEE/PES Transmission and Distribution Conference and Exposition*, pages 1–6. IEEE, 2008. 2
- [3] Zahra Abbasi, Georgios Varsamopoulos, and Sandeep K. S. Gupta. **Thermal aware server provisioning and workload distribution for internet data centers.** In *HPDC '10*, 2010. 19
- [4] George Andreadis, F. Mastenbroek, Vincent van Beek, and A. Iosup. **Capelin: Data-Driven Compute Capacity Procurement for Cloud Data-centers Using Portfolios of Scenarios.** *IEEE Transactions on Parallel and Distributed Systems*, 33:26–39, 2021. 57
- [5] Anirban Bag and Mostafa A. Bassiouni. **Energy Efficient Thermal Aware Routing Algorithms for Embedded Biomedical Sensor Networks.** *2006 IEEE International Conference on Mobile Ad Hoc and Sensor Systems*, pages 604–609, 2006. 19
- [6] Abhilasha Bajracharya, Md Riaz Ahmed Khan, Semhar Michael, and R. Tonkoski. **Forecasting Data Center Load Using Hidden Markov Model.** *2018 North American Power Symposium (NAPS)*, pages 1–5, 2018. 9

- [7] R. Baldick. **Incentive properties of coincident peak pricing.** *Journal of Regulatory Economics*, 54:165–194, 2018. 5, 36
- [8] Ilyas Bambrik. **A Survey on Cloud Computing Simulation and Modeling.** *SN Comput. Sci.*, 1:249, 2020. 6
- [9] L. Barroso, Urs Hölzle, and P. Ranganathan. **The Datacenter as a Computer: Designing Warehouse-Scale Machines, Third Edition.** In *The Datacenter as a Computer*, 2018. 38, 42
- [10] Luiz André Barroso and Urs Hölzle. **The Case for Energy-Proportional Computing.** *Computer*, 40, 2007. 18
- [11] Luiz André Barroso and Urs Hölzle. **The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines.** In *The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines*, 2008. 41, 42
- [12] N. Bates, G. Ghatikar, G. Abdulla, et al. **Electrical Grid and Supercomputing Centers: An Investigative Analysis of Emerging Opportunities and Challenges.** *Informatik-Spektrum*, 38:111–127, 2014. 36, 40
- [13] Frank Bellosa. **The benefits of event: driven energy accounting in power-sensitive systems.** In *EW 9*, 2000. 41
- [14] A. Beloglazov, J. Abawajy, and R. Buyya. **Energy-aware resource allocation heuristics for efficient management of data centers for Cloud computing.** *Future Gener. Comput. Syst.*, 28:755–768, 2012. 43
- [15] Anton Beloglazov, Rajkumar Buyya, Young Choon Lee, and Albert Y. Zomaya. **A Taxonomy and Survey of Energy-Efficient Data Centers and Cloud Computing Systems.** *Advances in Computers*, 82:47–111, 2010. 19
- [16] Sonja Bezjak, April Clyburne-Sherin, Philipp Conzett, et al. **The open science training handbook.** Technical report, [sn], 2018. 11
- [17] Dejene Boru, Dzmitry Kliazovich, Fabrizio Granelli, et al. **Energy-efficient data replication in cloud computing datacenters.** In *GLOBECOM Workshops*, 2013. 44

- [18] D. Brooks, V. Tiwari, and M. Martonosi. **Wattch: a framework for architectural-level power analysis and optimizations.** *Proceedings of 27th International Symposium on Computer Architecture (IEEE Cat. No.RS00201)*, pages 83–94, 2000. 41
- [19] R. Brown, Icf Incorporated, and Erg Incorporated. **Report to Congress on Server and Data Center Energy Efficiency: Public Law 109-431.** *Lawrence Berkeley National Laboratory*, 2008. 4, 37
- [20] Richard E. Brown, Eric R. Masanet, Bruce Nordman, et al. **Report to Congress on Server and Data Center Energy Efficiency: Public Law 109-431.** Technical report, Berkeley, CA, 06/2008 2008. 18
- [21] James Bucek, K. Lange, and J. Kistowski. **SPEC CPU2017: Next-Generation Compute Benchmark.** *Companion of the 2018 ACM/SPEC International Conference on Performance Engineering*, 2018. 77
- [22] C. Butler. **Climate Change, Health and Existential Risks to Civilization: A Comprehensive Review (1989–2013).** *International Journal of Environmental Research and Public Health*, 15, 2018. 1
- [23] NieuweStroom B.V. **Tarieven zakelijk grootverbruik.** Available at <https://www.nieuwestroom.nl/producten/fullflex-stroom-en-gas-grootverbruik/tarieven-zakelijk-grootverbruik/>. 76
- [24] Rodrigo N. Calheiros, Rajiv Ranjan, Anton Beloglazov, et al. **CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms.** *Softw., Pract. Exper.*, 41:23–50, 2011. 44
- [25] H. Casanova, A. Giersch, Arnaud Legrand, et al. **Versatile, scalable, and accurate simulation of distributed applications and platforms.** *J. Parallel Distributed Comput.*, 74:2899–2917, 2014. 27, 44
- [26] G. G. Castañé, Alberto Núñez, P. Llopis, and J. Carretero. **E-mc2: A formal framework for energy modelling in cloud computing.** *Simul. Model. Pract. Theory*, 39:56–75, 2013. 27, 44, 79

- [27] Fay Chang, K. Farkas, and P. Ranganathan. **Energy-Driven Statistical Sampling: Detecting Software Hotspots**. In *PACS*, 2002. 6, 41
- [28] H. Chen, A. Coskun, and M. Caramanis. **Real-time power control of data centers for providing Regulation Service**. *52nd IEEE Conference on Decision and Control*, pages 4314–4321, 2013. 40
- [29] H. Chen, Can Hankendi, M. Caramanis, and A. Coskun. **Dynamic server power capping for enabling data center participation in power markets**. *2013 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 122–129, 2013. 40
- [30] H. Chen, M. Caramanis, and A. Coskun. **The data center as a grid load stabilizer**. *2014 19th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pages 105–112, 2014. 40
- [31] H. Chen, Yijia Zhang, M. Caramanis, and A. Coskun. **EnergyQARE: QoS-Aware Data Center Participation in Smart Grid Regulation Service Reserve Provision**. *ACM Trans. Model. Perform. Evaluation Comput. Syst.*, 4: 2:1–2:31, 2019. 5
- [32] Hao Chen, Can Hankendi, Michael C Caramanis, and Ayse K Coskun. **Dynamic server power capping for enabling data center participation in power markets**. In *2013 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 122–129. IEEE, 2013. 4
- [33] Lijun Chen and N. Li. **On the Interaction Between Load Balancing and Speed Scaling**. *IEEE Journal on Selected Areas in Communications*, 33: 2567–2578, 2015. 40
- [34] Y. Chen, D. Gmach, C. Hyser, et al. **Integrated management of application performance, power and cooling in data centers**. *2010 IEEE Network Operations and Management Symposium - NOMS 2010*, pages 615–622, 2010. 37
- [35] Andrew A. Chien, Richard Wolski, and Fan Yang. **Zero-Carbon Cloud: A Volatile Resource for High-Performance Computing**. *2015 IEEE International Conference on Computer and Information Technology; Ubiquitous*

- Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing*, pages 1997–2001, 2015. 19
- [36] Todd L. Cignetti, K. Komarov, and C. Ellis. **Energy estimation tools for the Palm**. In *MSWIM '00*, 2000. 41
- [37] Pavel Machek Con Kolivas. **Linux CPU Load**. Available at <https://www.kernel.org/doc/html/latest/admin-guide/cpu-load.html>. 17
- [38] A. Conejo, J. Morales, and L. Baringo. **Real-Time Demand Response Model**. *IEEE Transactions on Smart Grid*, 1:236–242, 2010. 2, 36
- [39] Gilberto Contreras and M. Martonosi. **Power prediction for Intel XScale/spl reg/ processors using performance monitoring unit events**. *ISLPED '05. Proceedings of the 2005 International Symposium on Low Power Electronics and Design, 2005.*, pages 221–226, 2005. 6, 41
- [40] Jonathan Corbet. **Linux Per-entity load tracking**. Available at <https://lwn.net/Articles/531853/>. 17
- [41] International Business Machines Corporation. **IBM CPU utilization**. Available at <https://www.ibm.com/docs/en/informix-servers/12.10?topic=performance-cpu-utilization>. 16
- [42] Rafael Ferreira da Silva, Anne-Cécile Orgerie, H. Casanova, et al. **Accurately Simulating Energy Consumption of I/O-Intensive Scientific Workflows**. In *ICCS*, 2019. 27, 44
- [43] Dutch Data Center Association (DDA). **Data Center Factsheet**. Available at <https://www.dutchdatacenters.nl/en/factsheet/>. 4
- [44] Edward De Bono. *Six thinking hats*. Penguin uk, 2017. 47, 48
- [45] A. Demirbas. **Potential applications of renewable energy sources, biomass combustion problems in boiler power systems and combustion related environmental issues**. *Progress in Energy and Combustion Science*, 31:171–192, 2005. 3

- [46] Eaton. **Eaton Launches Industry First UPS-as-a-Reserve Service to Support the Power Grid in Frequency Containment Reserve**. Available at <http://powerquality.eaton.com/emea/about-us/news-events/2017/pr031017.asp> (2017/10/03). 40
- [47] D. Economou, S. Rivoire, C. Kozyrakis, and P. Ranganathan. **Full-System Power Analysis and Modeling for Server Environments**. 2006. 6, 41, 42
- [48] Xiaobo Fan, W. Weber, and L. Barroso. **Power provisioning for a warehouse-sized computer**. In *ISCA '07*, 2007. 41, 42, 43
- [49] X. Fang, S. Misra, G. Xue, and D. Yang. **Smart Grid — The New and Improved Power Grid: A Survey**. *IEEE Communications Surveys & Tutorials*, 14:944–980, 2012. 1
- [50] ASHRAE (Firm). *Thermal guidelines for Data processing environments*. ASHRAE, 2015. 39, 40
- [51] J. Flinn and M. Satyanarayanan. **PowerScope: a tool for profiling the energy usage of mobile applications**. *Proceedings WMCSA'99. Second IEEE Workshop on Mobile Computing Systems and Applications*, pages 2–10, 1999. 6, 41
- [52] Y. Fu, X. Han, K. Baker, and W. Zuo. **Assessments of data centers for provision of frequency regulation**. *Applied Energy*, 277:115621, 2020. 36
- [53] Yang-Yang Fu, W. Zuo, and K. Baker. **Multi-market optimization of a data center without storage systems**. 2020. 36
- [54] Yogesh Fulpagare, Y. Joshi, and A. Bhargav. **Rack level forecasting model of data center**. *2017 16th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, pages 824–829, 2017. 9
- [55] Anshul Gandhi, Mor Harchol-Balter, Rajarshi Das, and C. Lefurgy. **Optimal power allocation in server farms**. In *SIGMETRICS '09*, 2009. 19, 43
- [56] Anshul Gandhi, Y. Chen, D. Gmach, et al. **Minimizing data center SLA violations and power consumption via hybrid resource provisioning**.

- 2011 International Green Computing Conference and Workshops*, pages 1–8, 2011. 37
- [57] Mahdi Ghamkhari and H. Rad. **Energy and Performance Management of Green Data Centers: A Profit Maximization Approach.** *IEEE Transactions on Smart Grid*, 4:1017–1025, 2013. 6, 37
- [58] G. Ghatikar. **Demand Response Opportunities and Enabling Technologies for Data Centers: Findings From Field Studies.** 2012. 4, 5, 37
- [59] Sanjay Ghemawat, H. Gobiuff, and Shun-Tak Leung. **The Google file system.** In *SOSP '03*, 2003. 18
- [60] James Glanz. **Power, pollution and the internet.** *The New York Times*, 22, 2012. 5, 37
- [61] T. Guérout, T. Monteil, Georges Da Costa, et al. **Energy-aware simulation with DVFS.** *Simul. Model. Pract. Theory*, 39:76–91, 2013. 44, 77
- [62] Sandeep K. S. Gupta, Rose Robin Gilbert, Ayan Banerjee, et al. **GDCSim: A tool for analyzing Green Data Center design and resource management techniques.** *2011 International Green Computing Conference and Workshops*, pages 1–8, 2011. 19, 44
- [63] S. Gurumurthi, A. Sivasubramaniam, M. J. Irwin, et al. **Using complete machine simulation for software power estimation: the SoftWatt approach.** *Proceedings Eighth International Symposium on High Performance Computer Architecture*, pages 141–150, 2002. 41
- [64] David Guyon, Anne-Cécile Orgerie, Christine Morin, and Deb Agarwal. **How Much Energy can Green HPC Cloud Users Save?** In *2017 25th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP)*, pages 416–420. IEEE, 2017. 4
- [65] Daniel Hackenberg, R. Schöne, T. Ilsche, et al. **An Energy Efficiency Feature Survey of the Intel Haswell Processor.** *2015 IEEE International Parallel and Distributed Processing Symposium Workshop*, pages 896–904, 2015. 26

- [66] Richard R. Hamming. **The Art of Doing Science and Engineering: Learning to Learn**. 1997. 10
- [67] Hongyu He. **Modelling Energy Consumption in the OpenDC Datacenter Simulator for Analysing Energy-Aware Cloud Infrastructure**. Technical report, 2020. 65, 78
- [68] F. Heinrich, Tom Cornebize, A. Degomme, et al. **Predicting the Energy-Consumption of MPI Applications at Scale Using Only a Single Node**. *2017 IEEE International Conference on Cluster Computing (CLUSTER)*, pages 92–102, 2017. 20, 27, 44
- [69] Gernot Heiser. **Systems Benchmarking Crimes**. Available at <http://www.cse.unsw.edu.au/~Gernot/benchmarking-crimes.html>, 2020. 11
- [70] R. Hemmati, H. Saboori, and M. A. Jirdehi. **Stochastic planning and scheduling of energy storage systems for congestion management in electric power systems including renewable energy resources**. *Energy*, 133:380–387, 2017. 3
- [71] Jin-Man Heo, P. Jayachandran, I. Shin, et al. **OptiTuner: On Performance Composition and Server Farm Energy Minimization Application**. *IEEE Transactions on Parallel and Distributed Systems*, 22:1871–1878, 2011. 37, 40
- [72] Hongxun Hui, Yi Ding, Qingxin Shi, et al. **5G network-based Internet of Things for demand response in smart grid: A survey on application potential**. *Applied Energy*, 257:113972, 2020. 2
- [73] Amazon Inc. **Amazon EC2 Spot Instances**. Available at <https://aws.amazon.com/ec2/spot/?cards.sort-by=item.additionalFields.startDateTime&cards.sort-order=asc>. 4
- [74] Intel .Inc. **Intel® 64 and IA-32 Architectures Optimization Reference Manual**. Available at <https://www.intel.com/content/dam/www/public/us/en/documents/manuals/64-ia-32-architectures-optimization-manual.pdf> (2016/06), . 26

- [75] Intel .Inc. **Power Management in Intel® Architecture Servers.** Available at https://www.intel.com/content/dam/support/us/en/documents/motherboards/server/sb/power_management_of_intel_architecture_servers.pdf (2009/04), . 23, 42
- [76] VMware Inc. **VMware vSphere Documentation.** Available at https://docs.vmware.com/en/VMware-vSphere/index.html?topic=%2Fcom.vmware.wssdk.apiref.doc%2Fcpu_counters.html. 16
- [77] Uptime Institute. **Data center PUEs flat since 2013.** Available at <https://journal.uptimeinstitute.com/data-center-pues-flat-since-2013/>. 15
- [78] Intel. **What exactly is a P-state?** Available at <https://software.intel.com/content/www/us/en/develop/blogs/what-exactly-is-a-p-state-pt-1.html?language=en>. 21
- [79] A. Iosup, L. Versluis, A. Trivedi, et al. **The AtLarge Vision on the Design of Distributed Systems and Ecosystems.** *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, pages 1765–1776, 2019. 10
- [80] Alexandru Iosup, Georgios Andreadis, Vincent Van Beek, et al. **The OpenDC vision: Towards collaborative datacenter simulation and exploration for everybody.** In *2017 16th International Symposium on Parallel and Distributed Computing (ISPDC)*, pages 85–94. IEEE, 2017. 6, 19, 44
- [81] California ISO. **California ISO website.** Available at <http://www.caiso.com/Pages/default.aspx>. 34, 35
- [82] R. Jain. **The art of computer systems performance analysis - techniques for experimental design, measurement, simulation, and modeling.** In *Wiley professional computing*, 1991. 16
- [83] Ray Jain. **The Art of Computer Systems Performance Analysis.** In *Int. CMG Conference*, 1990. 10

- [84] Anish Jindal, G. Aujla, N. Kumar, and M. Villari. **GUARDIAN: Blockchain-Based Secure Demand Response Management in Smart Grid System.** *IEEE Transactions on Services Computing*, 13:613–624, 2020. 2
- [85] R. Johari and J. Tsitsiklis. **Parameterized Supply Function Bidding: Equilibrium and Efficiency.** *Oper. Res.*, 59:1079–1089, 2011. 5, 36
- [86] Murat Salim Karabinaoglu and T. Gözel. **Load forecasting modelling of data centers and IT systems by using artificial neural networks.** *2017 10th International Conference on Electrical and Electronics Engineering (ELECO)*, pages 62–66, 2017. 9
- [87] Gabor Kecskemeti. **DISSECT-CF: A simulator to foster energy-aware scheduling in infrastructure clouds.** *Simul. Model. Pract. Theory*, 58:188–218, 2015. 44
- [88] Gabor Kecskemeti, Wajdi Hajji, and Fung Po Tso. **Modelling Low Power Compute Clusters for Cloud Simulation.** *2017 25th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP)*, pages 39–45, 2017. 44
- [89] W. Kempton and J. Tomic. **Vehicle-to-grid power fundamentals: Calculating capacity and net revenue.** *Journal of Power Sources*, 144:268–279, 2005. 3
- [90] The kernel development community. **intel_pstate CPU Performance Scaling Driver.** Available at https://www.kernel.org/doc/html/latest/admin-guide/pm/intel_pstate.html, . 26
- [91] The kernel development community. **The Linux Kernel Documentations.** Available at <https://www.kernel.org/doc/html/latest/index.html>, . 23
- [92] N. Kheir. **Systems Modeling and Computer Simulation.** 1995. 10
- [93] Siwar Khemakhem, M. Rekik, and L. Krichen. **Double layer home energy supervision strategies based on demand response and plug-in electric**

- vehicle control for flattening power load curves in a smart grid.** *Energy*, 167:312–324, 2019. 2
- [94] Dzmitry Kliazovich, Pascal Bouvry, and Samee Ullah Khan. **GreenCloud: a packet-level simulator of energy-aware cloud computing data centers.** *The Journal of Supercomputing*, 62:1263–1283, 2010. 44
- [95] Jonathan Koomey et al. **Growth in data center electricity use 2005 to 2010.** *A report by Analytical Press, completed at the request of The New York Times*, 9(2011):161, 2011. 4
- [96] Jonathan G Koomey et al. **Estimating total power consumption by servers in the US and the world,** 2007. 18
- [97] K. Kurowski, A. Oleksiak, W. Piatek, et al. **DCworms - A tool for simulation of energy efficiency in distributed computing infrastructures.** *Simul. Model. Pract. Theory*, 39:135–151, 2013. 44
- [98] D. Kusic, J. Kephart, James E. Hanson, et al. **Power and performance management of virtualized computing environments via lookahead control.** *Cluster Computing*, 12:1–15, 2008. 19, 43
- [99] Etienne Le Sueur and Gernot Heiser. **Dynamic voltage and frequency scaling: The laws of diminishing returns.** In *Proceedings of the 2010 international conference on Power aware computing and systems*, pages 1–8, 2010. 4
- [100] Hongjian Li, Yuyan Zhao, and Shuyong Fang. **CSL-driven and energy-efficient resource scheduling in cloud data center.** *The Journal of Supercomputing*, 76:481 – 498, 2019. 4
- [101] N. Li, Lijun Chen, and S. Low. **Optimal demand response based on utility maximization in power networks.** *2011 IEEE Power and Energy Society General Meeting*, pages 1–8, 2011. 36
- [102] Yang Li, David Chiu, Changbin Liu, et al. **Towards dynamic pricing-based collaborative optimizations for green data centers.** *2013 IEEE 29th International Conference on Data Engineering Workshops (ICDEW)*, pages 272–278, 2013. 36

- [103] X. Liang. **Emerging Power Quality Challenges Due to Integration of Renewable Energy Sources**. *IEEE Transactions on Industry Applications*, 53:855–866, 2017. 1
- [104] Seung-Hwan Lim, Bikash Sharma, Gunwoo Nam, et al. **MDCSim: A multi-tier data center simulation, platform**. *2009 IEEE International Conference on Cluster Computing and Workshops*, pages 1–9, 2009. 19
- [105] Minghong Lin, A. Wierman, L. Andrew, and Eno Thereska. **Dynamic right-sizing for power-proportional data centers**. *2011 Proceedings IEEE IN-FOCOM*, pages 1098–1106, 2011. 37, 40
- [106] Minghong Lin, Zhenhua Liu, A. Wierman, and L. Andrew. **Online algorithms for geographical load balancing**. *2012 International Green Computing Conference (IGCC)*, pages 1–10, 2012. 6, 37
- [107] Changbin Liu, Lu Ren, B. T. Loo, et al. **Cologne: A Declarative Distributed Constraint Optimization Platform**. *Proc. VLDB Endow.*, 5:752–763, 2012. 36
- [108] Zhenhua Liu, A. Wierman, Y. Chen, et al. **Data center demand response: avoiding the coincident peak via workload shifting and local generation**. In *SIGMETRICS '13*, 2013. 5, 37
- [109] Zhenhua Liu, Iris Liu, S. Low, and A. Wierman. **Pricing data center demand response**. In *SIGMETRICS '14*, 2014. 5, 6, 36, 37
- [110] Zhenhua Liu, Minghong Lin, A. Wierman, et al. **Greening Geographical Load Balancing**. *IEEE/ACM Transactions on Networking*, 23:657–671, 2015. 6, 37
- [111] B. Louis, K. Mitra, S. Saguna, and C. Åhlund. **CloudSimDisk: Energy-Aware Storage Simulation in CloudSim**. *2015 IEEE/ACM 8th International Conference on Utility and Cloud Computing (UCC)*, pages 11–15, 2015. 44
- [112] A. Malik, Kashif Bilal, S. Malik, et al. **CloudNetSim++: A GUI Based Framework for Modeling and Simulation of Data Centers in OM-**

- NeT++**. *IEEE Transactions on Services Computing*, 10:506–519, 2017. 44
- [113] Christopher Malone and Christian Belady. **Metrics to characterize data center & IT equipment energy use**. In *Proceedings of the Digital Power Forum, Richardson, TX*, volume 68. sn, 2006. 14, 15, 18, 19
- [114] K. Mare. **Demand Response and Open Automated Demand Response Opportunities for Data Centers**. *Lawrence Berkeley National Laboratory*, 2010. 5, 38, 40
- [115] András Márkus, A. Kertész, and Gabor Kecskemeti. **Cost-Aware IoT Extension of DISSECT-CF**. *Future Internet*, 9:47, 2017. 44
- [116] Deborah T Marr, Frank Binns, David L Hill, et al. **Hyper-Threading Technology Architecture and Microarchitecture**. *Intel Technology Journal*, 6 (1), 2002. 26
- [117] George Marsh. **From intermittent to variable: can we manage the wind?** *Renewable energy focus*, 10(5):42–47, 2009. 3
- [118] Eric Masanet, Arman Shehabi, Nuo Lei, et al. **Recalibrating global data center energy-use estimates**. *Science*, 367(6481):984–986, 2020. 4
- [119] Fabian Mastenbroek, Georgios Andreadis, Soufiane Jounaid, et al. **OpenDC 2.0: Convenient modeling and simulation of emerging technologies in cloud datacenters**. *CCGRID*, 2021. 44, 57
- [120] Lykomidis Mastroleon, Nicholas Bambos, Christoforos E. Kozyrakis, and Dimitris Economou. **Automatic power management schemes for Internet servers and data centers**. *GLOBECOM '05. IEEE Global Telecommunications Conference, 2005.*, 2:5 pp.–, 2005. 19
- [121] David Meisner, Christopher M. Sadler, L. Barroso, et al. **Power management of online data-intensive services**. *2011 38th Annual International Symposium on Computer Architecture (ISCA)*, pages 319–330, 2011. 37

- [122] Tridib Mukherjee, Ayan Banerjee, Georgios Varsamopoulos, et al. **Spatio-temporal thermal-aware job scheduling to minimize energy consumption in virtualized heterogeneous data centers.** *Comput. Networks*, 53: 2888–2904, 2009. 19
- [123] UNITED NATIONS. **Paris Agreement.** Available at https://unfccc.int/sites/default/files/english_paris_agreement.pdf. 1
- [124] Alberto Núñez, J. L. Vázquez-Poletti, A. Caminero, et al. **iCanCloud: A Flexible and Scalable Cloud Infrastructure Simulator.** *Journal of Grid Computing*, 10:185–209, 2012. 44, 79
- [125] Essent N.V. **ENERGIE BESPAREN VIA ESSENT.** Available at <https://www.essent.nl/content/particulier/energie-besparen/index.html>. 76
- [126] U.S. Department of State. **Leaders Summit on Climate.** Available at <https://www.state.gov/leaders-summit-on-climate/>. 1
- [127] J. Ousterhout. **Always measure one level deeper.** *Communications of the ACM*, 61:74 – 83, 2018. 11
- [128] Fingrid Oyj. **Two-price and one-price system.** Available at <https://www.fingrid.fi/en/electricity-market/balance-service/description-of-balance-model/two-price-and-one-price-system/>. 91
- [129] N. Paterakis, O. Erdinç, A. Bakirtzis, and J. Catalão. **Optimal Household Appliances Scheduling Under Day-Ahead Pricing and Load-Shaping Demand Response Strategies.** *IEEE Transactions on Industrial Informatics*, 11:1509–1519, 2015. 3
- [130] Michael K Patterson, Stephen W Poole, Chung-Hsing Hsu, et al. **TUE, a new energy-efficiency metric applied at ORNL’s Jaguar.** In *International Supercomputing Conference*, pages 372–382. Springer, 2013. 15
- [131] Massoud Pedram. **Energy-efficient datacenters.** *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 31(10):1465–1484, 2012. 4

- [132] K. Peffers, T. Tuunanen, M. Rothenberger, and S. Chatterjee. **A Design Science Research Methodology for Information Systems Research**. *Journal of Management Information Systems*, 24:45 – 77, 2008. 10
- [133] M. Pipattanasomporn, M. Kuzlu, S. Rahman, and Y. Teklu. **Load Profiles of Selected Major Household Appliances and Their Demand Response Opportunities**. *IEEE Transactions on Smart Grid*, 5:742–750, 2014. 3
- [134] PricewaterhouseCoopers. **Vergelijking van gas-en elektriciteitsprijzen 2017**. Available at <https://zoek.officielebekendmakingen.nl/blg-850506.pdf/>. 76
- [135] Asfandyar Qureshi, Rick Weber, H. Balakrishnan, et al. **Cutting the electric bill for internet-scale systems**. In *SIGCOMM '09*, 2009. 6, 37
- [136] H. Rad and A. Leon-Garcia. **Optimal Residential Load Control With Price Prediction in Real-Time Electricity Pricing Environments**. *IEEE Transactions on Smart Grid*, 1:120–133, 2010. 36
- [137] R. Raghavendra, P. Ranganathan, V. Talwar, et al. **No "power" struggles: coordinated multi-level power management for the data center**. In *ASPLOS*, 2008. 19, 43
- [138] Lei Rao, X. Liu, Le Xie, and Wenyu Liu. **Minimizing Electricity Cost: Optimization of Distributed Internet Data Centers in a Multi-Electricity-Market Environment**. *2010 Proceedings IEEE INFOCOM*, pages 1–9, 2010. 6, 37
- [139] N. Rasmussen. **Electrical Efficiency Modeling for Data Centers**. 2007. 80
- [140] F. Rawson. **MEMPOWER: A Simple Memory Power Analysis Tool Set**. 2004. 41
- [141] Ian Rogers. **The Google Pagerank algorithm and how it works**. 2002. 73
- [142] C. Root, H. Presume, D. Proudfoot, et al. **Using battery energy storage to reduce renewable resource curtailment**. *2017 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT)*, pages 1–5, 2017. 8

- [143] S. Ruepp, Artur Pilimon, Jakob Thrane, et al. **Combining hardware and simulation for datacenter scaling studies**. *2017 International Conference on Optical Network Design and Modeling (ONDM)*, pages 1–6, 2017. 6
- [144] Frederick Ryckbosch, Stijn Polfliet, and Lieven Eeckhout. **Trends in server energy proportionality**. *Computer*, 44(9):69–72, 2011. 18
- [145] H. Sæle and O. S. Grande. **Demand Response From Household Customers: Experiences From a Pilot Study in Norway**. *IEEE Transactions on Smart Grid*, 2:102–109, 2011. 3, 5, 36
- [146] Ted Samson. **Power capping yields savings and floor space**, 2009. 40
- [147] K. Seta, Hideki Hara, Tomonori Kuroda, et al. **50% active-power saving without speed degradation using standby power reduction (SPR) circuit**. *Proceedings ISSCC '95 - International Solid-State Circuits Conference*, pages 318–319, 1995. 19
- [148] Subhadra Bose Shaw, C. Kumar, and Anil Kumar Singh. **Use of time-series based forecasting technique for balancing load and reducing consumption of energy in a cloud data center**. *2017 International Conference on Intelligent Computing and Control (I2C2)*, pages 1–6, 2017. 9
- [149] S. Shen, V. V. Beek, and A. Iosup. **Statistical Characterization of Business-Critical Workloads Hosted in Cloud Datacenters**. *2015 15th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, pages 465–474, 2015. 78
- [150] Suresh Siddha. **Process Scheduling Challenges in the Era of Multi-core Processors**. 2007. 26
- [151] Alexandru Iosup, Siqi Shen, Vincent van Beek. **The Grid Workloads Archive**. Available at <http://gwa.ewi.tudelft.nl/datasets/gwa-t-12-bitbrains>. 77
- [152] Hendrawan Soeleman and Kairshik Roy. **Ultra-low power digital sub-threshold logic circuits**. *Proceedings. 1999 International Symposium on Low Power Electronics and Design (Cat. No.99TH8477)*, pages 94–96, 1999. 19

- [153] Ingo Stadler. **Power grid balancing of energy systems with high renewable energy penetration by demand response.** *Utilities Policy*, 16(2):90–98, 2008. 2
- [154] Wenhong Tian, Yong Zhao, Minxian Xu, et al. **A Toolkit for Modeling and Simulation of Real-Time Virtual Machine Allocation in a Cloud Data Center.** *IEEE Transactions on Automation Science and Engineering*, 12: 153–161, 2015. 44
- [155] Michael Tighe, Gaston Keller, Michael Bauer, and Hanan Lutfiyya. **DC-Sim: A data centre simulation tool for evaluating dynamic virtualized resource management.** In *2012 8th international conference on network and service management (cnsm) and 2012 workshop on systems virtualization management (svm)*, pages 385–392. IEEE, 2012. 44
- [156] Adel Nadjaran Toosi, Rodrigo N. Calheiros, Ruppa K. Thulasiram, and Rajkumar Buyya. **Resource Provisioning Policies to Increase IaaS Provider’s Profit in a Federated Cloud Environment.** *2011 IEEE International Conference on High Performance Computing and Communications*, pages 279–287, 2011. 44
- [157] TenneT TSO. **Imbalance Pricing System.** Available at https://www.tennet.eu/fileadmin/user_upload/S0_NL/Imbalance_pricing_system.pdf. 93
- [158] Fort Collins Utilities. **Coincident Peak.** Available at <https://www.fcgov.com/utilities/business/manage-your-account/rates/electric/coincident-peak>. 36
- [159] N. Uv and Kishore Kumar G Pillai. **Energy Management of Cloud Data Center Using Neural Networks.** *2018 IEEE International Conference on Cloud Computing in Emerging Markets (CCEM)*, pages 85–89, 2018. 9
- [160] Thiago Lara Vasques, Pedro S. Moura, and A. D. Almeida. **A review on energy efficiency and demand response with focus on small and medium data centers.** *Energy Efficiency*, 12:1399–1428, 2019. 4

- [161] Akshat Verma, Puneet Ahuja, and A. Neogi. **pMapper: Power and Migration Cost Aware Application Placement in Virtualized Systems**. In *Middleware*, 2008. 19, 43
- [162] Andreea Valeria Vesa, Tudor Cioara, Ionut Anghel, et al. **Energy Flexibility Prediction for Data Center Engagement in Demand Response Programs**. *Sustainability*, 12(4):1417, 2020. 4
- [163] N. Vijaykrishnan, M. Kandemir, M. J. Irwin, et al. **Energy-driven integrated hardware-software optimizations using SimplePower**. *Proceedings of 27th International Symposium on Computer Architecture (IEEE Cat. No.RS00201)*, pages 95–106, 2000. 41
- [164] Inc. VMware. **Difference between cpu.usage and cpu.utilization counters for HostSystem object**. Available at <https://kb.vmware.com/s/article/2055995>. 16
- [165] Centraal Bureau voor de Statistiek. **Energy consumption private dwellings; type of dwelling and regions**. Available at <https://www.cbs.nl/en-gb/figures/detail/81528ENG?q=parts%20of%20the%20country>. 116
- [166] E. Vrettos, F. Oldewurtel, Fengtian Zhu, and G. Andersson. **Robust Provision of Frequency Reserves by Office Building Aggregations**. *IFAC Proceedings Volumes*, 47:12068–12073, 2014. 3
- [167] Hangsheng Wang, X. Zhu, L. Peh, and S. Malik. **Orion: a power-performance simulator for interconnection networks**. *35th Annual IEEE/ACM International Symposium on Microarchitecture, 2002. (MICRO-35)*. *Proceedings.*, pages 294–305, 2002. 41
- [168] Hao Wang, Jianwei Huang, X. Lin, and H. Rad. **Proactive Demand Response for Data Centers: A Win-Win Solution**. *IEEE Transactions on Smart Grid*, 7:1584–1596, 2016. 2, 4, 6, 37
- [169] K. Wang, Liqiu Gu, Xiaoming He, et al. **Distributed Energy Management for Vehicle-to-Grid Networks**. *IEEE Network*, 31:22–28, 2017. 3

- [170] Peijian Wang, Lei Rao, X. Liu, and Yong Qi. **D-Pro: Dynamic Data Center Operations With Demand-Responsive Electricity Prices in Smart Grid.** *IEEE Transactions on Smart Grid*, 3:1743–1754, 2012. 6, 37
- [171] Weiming Wang, AmirAli Abdolrashidi, N. Yu, and Daniel Wong. **Frequency regulation service provision in data center with computational flexibility.** *Applied Energy*, 251:113304–113304, 2019. 40
- [172] A. Wierman, Zhenhua Liu, Iris Liu, and H. Rad. **Opportunities and challenges for data center demand response.** *International Green Computing Conference*, pages 1–10, 2014. 3, 4, 5, 38
- [173] Wikipedia. **CPU utilization.** Available at https://en.wikipedia.org/wiki/Load_%28computing%29. 16
- [174] M. Wilkinson, M. Dumontier, I. J. Aalbersberg, et al. **The FAIR Guiding Principles for scientific data management and stewardship.** *Scientific Data*, 3, 2016. 11
- [175] H. Xu and B. Li. **Cost efficient datacenter selection for cloud services.** *2012 1st IEEE International Conference on Communications in China (ICCC)*, pages 51–56, 2012. 37
- [176] Yunjian Xu, N. Li, and S. Low. **Demand Response With Capacity Constrained Supply Function Bidding.** *IEEE Transactions on Power Systems*, 31:1377–1394, 2016. 5, 36
- [177] Zichuan Xu and Weifa Liang. **Minimizing the Operational Cost of Data Centers via Geographical Electricity Price Diversity.** *2013 IEEE Sixth International Conference on Cloud Computing*, pages 99–106, 2013. 15
- [178] Rahul Yadav, Weizhe Zhang, K. Li, et al. **An adaptive heuristic for managing energy consumption and overloaded hosts in a cloud data center.** *Wireless Networks*, 26:1905–1919, 2020. 4
- [179] Y. Yao, Longbo Huang, A. Sharma, et al. **Data centers power reduction: A two time scale approach for delay tolerant workloads.** *2012 Proceedings IEEE INFOCOM*, pages 1431–1439, 2012. 37

- [180] Y. Yu, H. Wang, Gang Yin, and C. Ling. **Reviewer Recommender of Pull-Requests in GitHub**. *2014 IEEE International Conference on Software Maintenance and Evolution*, pages 609–612, 2014. 11
- [181] John Zedlewski, Sumeet Sobti, Nitin Garg, et al. **Modeling Hard-Disk Power Consumption**. In *FAST*, 2003. 41
- [182] Linqun Zhang, Shaolei Ren, C. Wu, and Z. Li. **A truthful incentive mechanism for emergency demand response in colocation data centers**. *2015 IEEE Conference on Computer Communications (INFOCOM)*, pages 2632–2640, 2015. 6, 37
- [183] Q. Zhang, M. Zhani, S. Zhang, et al. **Dynamic energy-aware capacity provisioning for cloud computing environments**. In *ICAC '12*, 2012. 37, 40
- [184] Qin Zhang, Xifan Wang, Min Fu, and Jianxue Wang. **Smart grid from the perspective of demand response**. *Automation of Electric Power Systems*, 17:49–55, 2009. 2
- [185] Bo Zheng, Guannan Geng, Philippe Ciais, et al. **Satellite-based estimates of decline and rebound in China’s CO2 emissions during COVID-19 pandemic**. *Science Advances*, 6(49):eabd4998, 2020. 1
- [186] Y. Zhou, Y. Yi, Gaoying Cui, et al. **Demand response control strategy of groups of central air-conditionings for power grid energy saving**. *2016 IEEE International Conference on Power and Renewable Energy (ICPRE)*, pages 323–327, 2016. 2