

Energy Efficient Data Center Management Strategies

A thesis submitted to University of Hyderabad in partial fulfillment

for the degree of

Doctor of Philosophy

by

Vemula Dinesh Reddy

Reg.No. 14MCPC19



SCHOOL OF COMPUTER AND INFORMATION SCIENCES
UNIVERSITY OF HYDERABAD
HYDERABAD -500046

Telangana

India

May, 2018



CERTIFICATE

This is to certify that the thesis entitled “**Energy Efficient Data Center Management Strategies**” submitted by **Vemula Dinesh Reddy** bearing **Reg. No. 14MCPC19** in partial fulfillment of the requirements for the award of **Doctor of Philosophy in Computer Science** is a bonafide work carried out by him under my supervision and guidance at IDRBT, Hyderabad.

This thesis is free from plagiarism and has not been submitted previously in part or in full to this or any other University or Institution for award of any degree or diploma.

Parts of this thesis have been published online in the following publications:

1. Soft Computing 2017 (Chapter 3).
2. IEEE Transactions on Sustainable Computing 2017 (Chapter 5).
3. IEEE IT Professional 2017 (Chapter 6).

Further, the student has passed the following courses towards fulfillment of coursework requirement for Ph.D:

	Course Code	Name	Credits	Pass/Fail
1	BT701	Data Structures and Algorithms	4	Pass
2	BT702	Operating System and Programming	4	Pass
3	BT709	Advanced Software Engineering	4	Pass
4	BT718	Enterprise Architecture	4	Pass

Supervisor	Director	Dean
Dr. G. R. Gangadharan	Dr. A.S. Ramasastrri	Prof. Arun Agrawal
Associate Professor	IDRBT	School of Computer and
IDRBT	Hyderabad-500 057, India	Information Sciences, UOH
Hyderabad-500 057, India		Hyderabad-500 046, India

DECLARATION

I, **Vemula Dinesh Reddy**, hereby declare that this thesis entitled “**Energy Efficient Data Center Management Strategies**” submitted by me under the guidance and supervision of **Dr. G. R. Gangadharan**, is a bonafide research work and is free from plagiarism. I also declare that it has not been submitted previously in part or in full to this University or any other University or Institution for the award of any degree or diploma. I hereby agree that my thesis can be deposited in Shodganga/INFLIBNET.

A report on plagiarism statistics from the University Librarian is enclosed.

Date:

Signature of the Student

(Vemula Dinesh Reddy)

Reg. No.: 14MCPC19

//Countersigned//

Signature of the Supervisor:

(Dr. G. R. Gangadharan)

Dedicated To My Family & Teachers

Acknowledgements

I would like to thank my parents **Sri. Prabhakar Reddy** and **Smt. Rajyalaxmi** for their best support and encouragement all times in pursuing my dreams. I would also like to thank my wife **Shruthi** and daughter **Rigveditha** for their infallible love. I would also like to thank my grand parents **Sri. Venkat Reddy** and **Smt. Venkatamma** for their best support.

I owe my greatest gratitude to my supervisor **Dr. G. R. Gangadharan**, Associate Professor, IDRBT, Hyderabad for his precious regular encouragement and timely support from the introductory level to the concluding level that enabled me to understand and implement the concepts learned.

I would like thank **Prof. Marco Aiello**, University of stuttgart, Germany for his inspiring thoughts during my research. I would like to specially thank **Dr. G. Subramanya V.R.K Rao**, AVP-Technology, Cognizant Technology Solutions for his suggestions during my research.

It is my previlage to thank **Prof. Arun Agrawal**, Dean, School of Computer and Information Sceinces (SCIS), UoH, Hyderabad for his academic support throughout research work. It is an honor for me to thank **Dr. A.S. Ramasastrri**, Director, IDRBT for reviewing the work and timely suggestions throughout my research. I also thank **Mr. B. Sambamurthy**, Former Director, IDRBT for extending his cooperation at the preliminary stage of my research. I would like to extend my sincere thanks to **Prof. V. Ravi** and **Dr. M.V.N.K. Prasad** for their regular reviews of the progress of research and kind inputs as the members of my Doctoral Research Committee.

I would like to express my thanks to my friends including J. Chandrashekar, M. Sandhya, D. Pradeep Kumar, K. Ilaiyah for their invaluable support. I would like to express my special thanks to my research colleagues for their technical support in my research work, their daily interactions with me and kind support. I would like to express my thanks to research associates Srujana, Vamshi Krishna, and Jyothi.

.....*Vemula Dinesh Reddy*

Abstract

The proliferation of Cloud computing has resulted in the establishment of large-scale data centers around the world. These data centers are most energy-intensive building types and consume large amounts of electrical energy resulting in high operating costs and carbon dioxide emissions. The high energy consumption of data centers is drawing more and more attention due to economic, social, and environmental concerns. Today, it is no surprise that reducing energy costs is one of the top priorities for many energy-related businesses.

Determining the optimal placement of virtual machines is an essential aspect of a data center to improve physical resource utilization and to reduce the energy consumption while satisfying the service level agreement. Energy efficient resource provisioning aims to find the near-optimal solution that improves the resource utilization and decreases the energy consumption of the data center in an acceptable time. Determining a set of virtual machines that can be migrated from an over-utilized or under-utilized host has a significant impact on the energy consumption of the data center. So, designing a virtual machine selection policy, considering different resources along with CPU utilization plays an important role in improving the energy efficiency of the data centers. In this thesis, we develop energy efficient virtual machine placement and selection algorithms based on soft computing approaches that minimize the energy consumption while fulfilling the service level agreements.

Forecasting data center electrical energy demand is very challenging due to dynamic nature and complexity of workloads. Developing forecasting models with accurate predictions give operators enough time to

avoid the risk of over-provisioning during non-peak period, and reduces the risk of under-provisioning in peak period. We develop two machine learning approaches for forecasting energy demand of the chillers in a data center.

In order to predict growth or set effective goals, it is important to choose the correct metric and being aware of their expressivity and potential limitations. Understanding and analyzing data center metrics allows the operators to have a better view on possible inefficiencies by focusing on the core parameters. It is necessary to compare the current approaches in a data center with industry standards and assess whether the practices are still valid and/or optimal. Determining and implementing the best practices for data center operations is needed to optimize the workflows and to decrease the operating costs in the long term. We analyze the potential metrics for the next-generation data centers and their inter dependencies. This thesis further provides a set of best practices to improve the energy efficiency of data centers.

This dissertation contributes to energy efficient data center management strategies by developing novel algorithms in the perspective of resource management and energy forecasting. Further, the thesis analyzes the best practices and metrics for sustainable data centers. More specifically, the results presented in this thesis outline opportunities to control the growing energy demand of data centers, so that IT can evolve into an energy efficient utility with the potential to facilitate a more sustainable expansion of services.

Contents

Acknowledgments	v
Abstract	vii
List of Figures	x
List of Tables	xi
1 Introduction	1
1.1 Data Centers	1
1.2 Typical Power Flow in a Data Center	4
1.3 Energy Consumption and Carbon Footprint of Data Centers	5
1.4 Problem Statement, Objectives, and Contributions	8
1.5 Thesis Organization	11
2 Literature Review	13
2.1 Virtual Machine Placement and Selection using Soft Computing Approaches	13
2.1.1 Non-heuristic Approaches	15
2.1.2 Heuristic Approaches	18
2.1.3 Meta-heuristic Approaches	20
2.2 Machine Learning Approaches for Data Center Monitoring	23
2.3 Analysis and Discussion	27
2.4 Justification for Present Work	29

3	Energy aware Virtual Machine Placement and Selection Approaches in Cloud Data Centers	32
3.1	Modelling Resource Allocation and Power Consumption in Cloud Data Centers	33
3.1.1	Problem Definition	33
3.1.2	Fitness Evaluation	35
3.2	System Architecture for Energy Efficient VM Placement and Selection (EE-VMPS)	36
3.3	Energy Efficient VM Placement Approaches in Cloud Data Centers	37
3.3.1	Approach 1 : VM Placement Using Modified Discrete Particle Swarm Optimization (MDPSO)	37
3.3.2	Approach 2 : VM Allocation Using Interactive PSO-GA (IP-SOGA)	41
3.3.3	Approach 3 : VM Allocation Using Imitation Based Optimization (IBO)	45
3.4	VM Selection and Migration using MBS-VM	48
3.5	Performance Evaluation	52
3.5.1	Experimental Environment	52
3.5.2	Analysis of VM Allocation Approaches	54
3.5.3	Performance Analysis in terms of Migrations	55
3.5.4	Performance Analysis in terms of SLA Violations	56
3.5.5	Performance Analysis in terms of Convergence	57
3.5.6	Speedup and Parallel Efficiency of IPSOGA	58
3.5.7	Analysis of VM Selection Algorithms	59
3.6	Summary	60
 4	 Machine Learning Approaches for Forecasting Data Center Energy Demand	 63
4.1	Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction	63
4.1.1	Approach 1: Multi-layer Feed Forward Neural Networks (MFNN)	64
4.1.2	Approach 2: Deep Learning with Parallel Stochastic Gradient Descent (DPSGD)	67

CONTENTS

4.2	Description of Data Set and Preprocessing	70
4.3	Experimental Analysis	71
4.3.1	Performance Evaluation	73
4.4	Summary	77
5	Metrics for Sustainable Data Centers	80
5.1	A Taxonomy of Data Center Metrics	80
5.1.1	Energy Efficiency Metrics	85
5.1.2	Cooling Metrics	87
5.1.3	Greenness Metrics	89
5.1.4	Performance / Productivity Metrics	91
5.1.5	Thermal and Air Management Metrics	93
5.1.6	Network Metrics	95
5.1.7	Storage Metrics	95
5.1.8	Security Metrics	98
5.1.9	Financial Impact Metrics	98
5.2	Analysis of Metrics	101
5.3	Summary	104
6	Best Practices for Sustainable Data Centers	107
6.1	Research Methodology	107
6.2	Data Center Management Best Practices	108
6.2.1	Energy Efficiency Practices	109
6.2.2	Cooling, Thermal, and Air Management Practices	111
6.2.3	Green Practices	114
6.2.4	Storage and Network Practices	115
6.2.5	Security Practices	116
6.3	Recommendations for Data Center Operators and IT Professionals .	118
6.4	Summary	121
7	Conclusions and Future Directions	123
	References	127
	List of Publications	169

APPENDICES	170
A Overview of Techniques Used	171
A.1 Particle Swarm Optimization (PSO)	171
A.2 Discrete Binary Version of PSO	172
A.3 Genetic Algorithm (GA)	172
A.4 Coordination of the particles for Modified Discrete Particle Swarm Optimization (MDPSO) approach	174
B Data Center Metrics Definitions	175
B.1 Energy Efficiency Metrics	175
B.2 Cooling Metrics	187
B.3 Greenness Metrics	190
B.4 Performance / Productivity Metrics	195
B.5 Thermal and Air Management Metrics	202
B.6 Network Metrics	209
B.7 Storage Metrics	212
B.8 Security Metrics	214
B.9 Financial Impact Metrics	219
C Data Center Assessment Checklist	223
D Fact Sheet for Publications	235

List of Figures

1.1	Basic Data Center Topology (Based on [4])	2
1.2	Data Center Power Flow (Based on [4])	5
2.1	A Typical Cloud Resource Provisioning Scenario	15
2.2	Classification of Approaches	28
2.3	Meta-heuristic Approaches	28
2.4	Machine learning applications for data center monitoring	29
2.5	Various machine learning approaches for data center monitoring	30
3.1	System Architecture of EE-VMPS	36
3.2	Flow diagram of IPSOGA	42
3.3	Comparison of active hosts in different algorithms	54
3.4	Number of Migrations vs. Number of Virtual Machines	56
3.5	Convergence Analysis	58
3.6	Comparison of VM selection algorithms	61
4.1	Structure of an Artificial Neuron	65
4.2	Basic structure of a four-layer feed forward network.	65
4.3	Experimental data vs predictions for test data	76
4.4	Experimental data vs predictions for training data	76
5.1	Categories of Components of Data Centers	81
5.2	Relationships between Energy Efficiency Metrics	89
5.3	Relationship between Green Metrics	91
5.4	Relationship between Thermal and Air Management Metrics	93
5.5	Relationship between Financial Metrics	101

LIST OF FIGURES

6.1	Hot Aisle / Cold Aisle Containment in Data Centers	113
A.1	Flow of Particle Swarm Optimization	172
A.2	Flow of Genetic Algorithm	173
A.3	Particle coordination	174
B.1	Energy Proportional System, adopted from [332]	198
D.1	Publication [1]	235
D.2	Publication [2]	236
D.3	Publication [3]	237
D.4	Publication [4]	238
D.5	Publication [5]	239

List of Tables

1.1	Tier Classification of Data Centers (Based on [5])	3
2.1	Non-heuristic Approaches	17
2.2	Heuristic approaches	19
2.3	Meta-heuristic approaches	22
2.4	Machine Learning approaches for data center monitoring	25
2.5	Machine Learning approaches for data center monitoring	26
3.1	Physical Machines Configurations	53
3.2	Virtual Machines Requirements	53
3.3	Performance comparison in terms of energy consumption	55
3.4	Overall SLA Violations	57
3.5	Speedup and Parallel Efficiency of IPSOGA	60
4.1	Comparison of different normalization functions	71
4.2	Comparison of different activation functions for MFNN	72
4.3	Comparison of different activation functions for DPSGD	72
4.4	Performance comparison of MFNN with two hidden layers	74
4.5	Performance comparison of DPSGD with two hidden layers	75
4.6	Performance comparison of various prediction models	77
4.7	Values for T and P for significance level of 0.05	77
4.8	Statistical analysis using Wilcoxon Signed Rank test (Z and P values)	78
5.1	Taxonomy of Data Center Metrics	84
5.2	Energy Efficiency Metrics	86
5.3	Cooling Metrics	88

LIST OF TABLES

5.4	Green Metrics	90
5.5	Performance Metrics	92
5.6	Thermal and Air Management Metrics	94
5.7	Network Metrics	96
5.8	Storage Metrics	97
5.9	Security Metrics	99
5.10	Financial Impact Metrics	100
6.1	Data centers configurations	109
6.2	Energy efficiency practices for data centers	112
6.3	Cooling, Thermal, and Air management practices for data centers	114
6.4	Green practices for data centers	116
6.5	Storage and Network practices for data centers	117
6.6	Security practices for data centers	118
6.7	Implementation issues and challenges in data centers	119
B.1	Efficiency level of PUE and DCiE	179
B.2	Standard Values of Airflow Efficiency	203
B.3	AHSRAE thermal recommendations for Class 1 data centers	203
B.4	Compliance Of RCI	207
B.5	SHI and RHI for different infrastructures	208
C.1	Key Questions and Metrics for Energy Efficiency	225
C.2	Key Questions and Metrics for Thermal and Air Management	226
C.3	Key Questions and Metrics for Cooling Plant	228
C.4	Key Questions and Metrics for Overall Performance and Distribu- tion Chain	229
C.5	Key Questions and Metrics Greenness	231
C.6	Key Questions and Metrics for Network	232
C.7	Key Questions and Metrics for Storage	233
C.8	Key Questions and Metrics for Security	234
D.1	Fact sheet for publication [1]	235
D.2	Fact sheet for publication [2]	236

LIST OF TABLES

D.3	Fact sheet for publication [3]	237
D.4	Fact sheet for publication [4]	238
D.5	Fact sheet for publication [5]	239

Chapter 1

Introduction

In today's data-centric world, people expect the right data to be available everywhere, anytime, cheaply, in growing quantities and 'processed-right'. Increasing data storage, processing and continuity requires a globally growing and more efficient data center management strategies to lower the environmental footprint. The increase in the number of online services users, leading to an ever-increasing demand for computing resources, has resulted in an increase in the requirement of power and space to host these computer resources and IT infrastructure in a data center. Thus, energy efficiency is becoming an important concern in designing and managing data centers.

1.1 Data Centers

Data centers are structures or groups of structures, dedicated to a centralized accommodation, operation and interconnection of Information and Communications Technology (ICT) equipment providing data storage, processing, and transportation services [1]. The data center encompasses all of the facilities and infrastructures for power distribution, Heating, Ventilation and Air Conditioning (HVAC) control, together with the necessary levels of resilience and security that are required to provide the desired service availability. Data centers are the critical infrastructures for any business nowadays [2]. Therefore, data centers are designed in such a way that their long term uninterrupted operation is guaranteed.

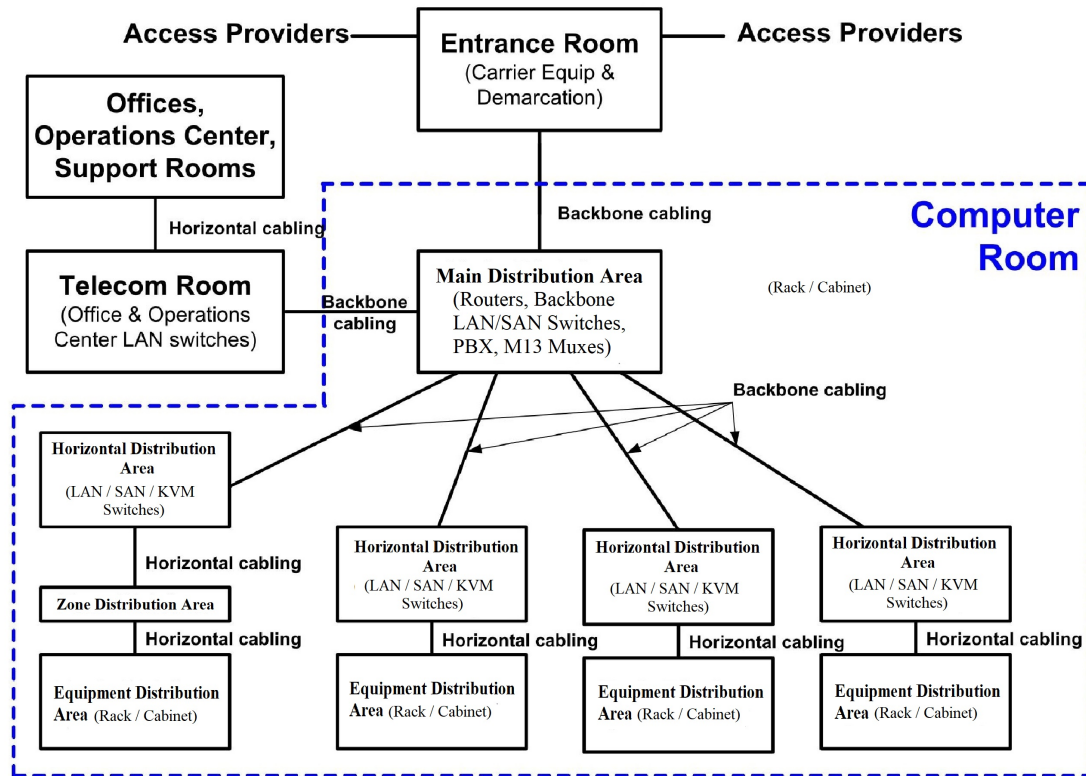


Figure 1.1: Basic Data Center Topology (Based on [4])

They employ various redundant or backup techniques in both software and hardware level to ensure their reliability. Further, they employ several air-conditioning controls and security solutions to ensure their thermal safety and security respectively. Data centers come in different sizes depending on their design objectives and functionalities. They range from small facilities hosting a few computers without sophisticated power and cooling system infrastructure to massive facilities hosting hundreds of thousands of servers and offering a variety of services, with the advent of virtualization and cloud computing [3].

A typical data center includes a single entrance room, possibly one or more telecommunications rooms, one main distribution area, and several horizontal distribution areas [4]. The topology of a typical data center is given in Figure 1.1. Entrance Room (ER) is the interface between the access provider and the data center structured cabling. Main Distribution Area (MDA) is the hub of the cabling system

and powerful switches. MDA may be located in the computer room. Equipment Distribution Area (EDA) is the space allocated for cabinets, racks, end equipment, and other communications hardware. Horizontal Distribution Area (HDA) is the space that supports cabling to the equipment distribution areas. The HDA houses cross-connects and active equipment (LAN, SAN, KVM switches) for connecting to the EDA(s). There is an optional interconnection area within HDA, called as Zone Distribution Area (ZDA), that provides flexibility for future expansion.

Tier Level	Requirements
1	<ul style="list-style-type: none"> • Consists of non-redundant capacity components and a single non-redundant distribution path serving the critical environment. • Susceptible to disruption from both planned and unplanned activities. • Expected availability of 99.671%.
2	<ul style="list-style-type: none"> • Meets or exceeds all Tier 1 requirements. • Consists of redundant capacity components and a single, non-redundant distribution path. • Susceptible to disruption from both planned and unplanned activities. • Expected availability is 99.741%.
3	<ul style="list-style-type: none"> • Meets or exceeds all Tier 1 and Tier 2 requirements. • N+1 Redundant capacity components and multiple independent distribution paths. • Susceptible to disruption from unplanned activities. • Expected availability of 99.982%.
4	<ul style="list-style-type: none"> • Meets or exceeds all Tier 1, Tier 2 and Tier 3 requirements. • 2N+1 redundant capacity components and multiple independent distribution paths. • Not susceptible to disruption from both planned and unplanned activities. • Expected availability of 99.995%.

Table 1.1: Tier Classification of Data Centers (Based on [5])

Data center reliability and resilience are often referred to as uptime and is rated by “tier”. The Tier classification, as defined by the Uptime Institute, presents the data center industry “a consistent mechanism for comparing typical facilities based on their up-time and facility performance” [6]. The Tier Standard levels describe the availability of a particular site infrastructure design. The higher the Tier

1.2 Typical Power Flow in a Data Center

level, the greater the expected availability. The requirements of each Tier Level are shown in Table 1.1 [5].

1.2 Typical Power Flow in a Data Center

The typical power flow in a data center consists of two key parts: facility infrastructure and IT equipment. The power is delivered to a data center by a local utility company. The utility power enters to the Automatic Transfer Switch (ATS) in the facility infrastructure. In case of emergency, the power comes from power generators. While the utility power is available, the power flows into subsidiary circuits, often called “switchgear”. The switchgear passes power to the uninterruptible power supply (UPS) units and key facility systems such as elevators, lighting, HVAC equipment. Figure 1.2 shows the typical power flow among the power, cooling, and IT systems of a data center [4]. The switchgear passes power to chillers, cooling towers, and Computer Room Air Conditioners or Handlers (CRACs/CRAHs). If the normal power source from the utility company is not available, the ATS triggers the power generator. Once the power generator starts up, the ATS switches the load from the normal power to the emergency power. The power enters the UPS that provides the emergency power to a load if the normal power source fails. The UPS protects sensitive IT equipment in a data center from power fluctuations and outage. The UPS is connected in-line with the battery backup system. In case of power outage, ATS turns on the power generators and the power is supplied to the IT load. After passing the UPS, the power flows to power distribution units (PDUs). The PDUs convert high voltage to a more usable voltage for IT equipment in a data center. This power is distributed to the downstream loads via a common circuit breaker. At this point, the power leaves the facility infrastructure boundary. The PDUs’ power flows to each power supply in the rack. The next step in the power flow in IT equipment is fans, which are one of the crucial factors to make data centers more energy efficient.

Figure 1.2 also shows the locations where power can be measured. Moving from the broad measurement to the detailed, the first is the data center power consumption at the meter (Measurement Point-A). This is used as the numerator

1.3 Energy Consumption and Carbon Footprint of Data Centers

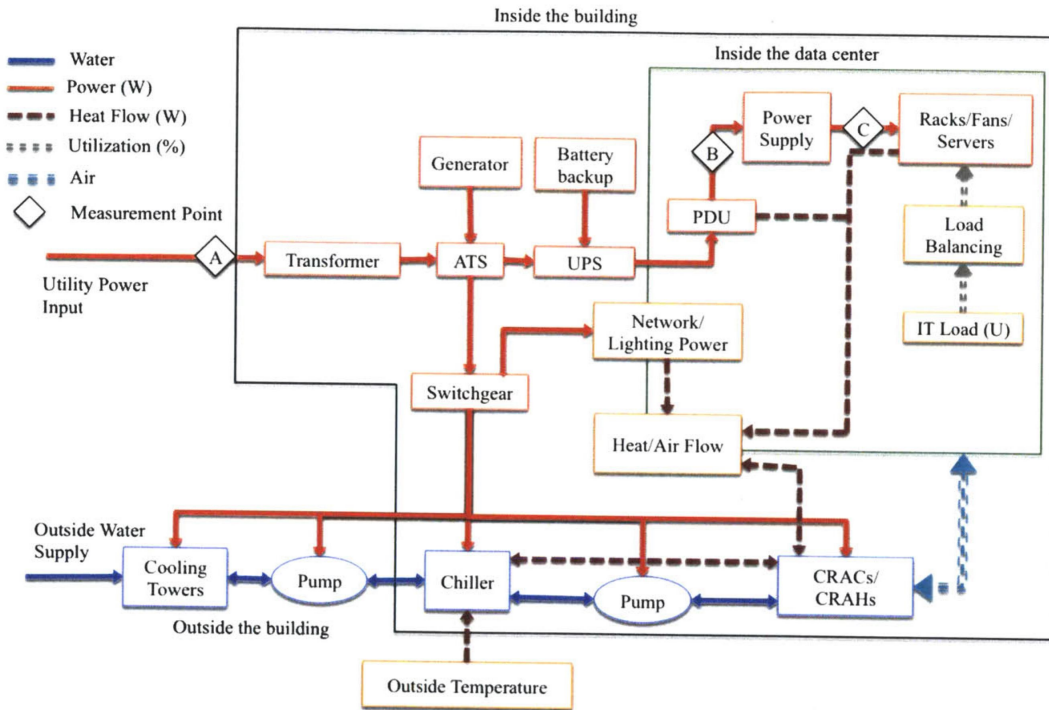


Figure 1.2: Data Center Power Flow (Based on [4])

in Power Usage Effectiveness (PUE) [7] calculation. The next place is Measurement Point-B that measures power consumed at a rack with intelligent rack PDUs (IRPs). These measurements denote the IT equipment power consumption, aggregated to a rack. A third place to measure power is Measurement Point-C, placed at the individual outlets of a rack PDU to get the total rack power consumption. By measuring power consumption at the component level, it is possible to take specific actions to improve energy efficiency of data centers.

1.3 Energy Consumption and Carbon Footprint of Data Centers

The use of cloud computing services and applications continues to increase at a rapid rate, leading to the rise of vast ‘hyperscale’ cloud data centers. These data centers are most energy-intensive building types and consume large amounts

1.3 Energy Consumption and Carbon Footprint of Data Centers

of electrical energy resulting in high operating costs and carbon dioxide (CO₂) emissions to the environment. The world's Information and communications technology (ICT) infrastructure is estimated to consume 1,500 TWh of electricity, roughly 10% of global usage. U.S Data centers consumed 1.4% and 1.8% of all the electricity used in U.S. in 2010 and 2014 respectively [8, 9]. Energy demand is set to grow more quickly, with demand growing more than 70% in next 20 years with data centers being the key contributors for this expansion [10]. A survey by IBM shows that the average resource utilization rate is lower than 20% in several data centers. Further, 70% of power is consumed by idle servers in several data centers [11, 12, 13]. Thus, we observe that a source of high energy consumption is not only the amount of computing resources used and power inefficiency of the hardware but also lies in the inefficient usage and dynamic power ranges of servers.

Energy consumption in data centers can be reduced by efficient use of servers in data centers, that can be achieved by efficiently placing virtual machines (VM) on servers and decommissioning of unused servers. With rapid development of virtualization technology, migration and dynamic placement of virtual machines becomes significant for efficient data center management [14, 15]. During the provisioning of resources to virtual machines, the resource utilization should be maximized by reducing the number of active hosts [16].

Live migration enables to respond to peak load periods by moving virtual machines from one physical machine to another physical machine. This technique is widely used to reduce underutilization in data centers by dynamically migrating VMs to few host considering different thresholds and Service Level Agreements (SLA). Then, idle physical machines can be turned to sleep mode to improve the energy efficiency. Migration takes place between two physical machines. Once a new configuration file is created on the target physical machine, the selected VM memory is copied to the target physical machine. Meanwhile if any memory page is changed on the source, those pages will be tagged and copied. Migration helps performing maintenance without disrupting operations, optimizing resource pools and avoids failing.

The operational efficiency of data centers assumes central importance. Even small gains in efficiency translate into end-user perceivable cost reductions, providing key competitive advantage. Data center consumers would like to provision

1.3 Energy Consumption and Carbon Footprint of Data Centers

resources transparently with minimal latency. From the perspective of consumers, the performance of a data center is measured in terms of response time, virtual machine provisioning time, etc. To serve the customers in a better way, data center providers should follow an optimal VM provisioning within the short time. The objective of energy efficient resource provisioning is to find the near-optimal solution that improves the resource utilization and decreases the energy consumption of the data center in an acceptable time. To reduce the energy consumption, in this thesis, we propose energy efficient virtual machine placement and selection algorithms using soft computing approaches that minimize the energy consumption while fulfilling the service-level-agreement.

The operations of a data center are quickly transforming from individual and disconnected tactical activities with a primary historical goal of “high service levels at any cost” to “service at what cost” using predictable approaches. Energy efficiency of a data center is influenced by many factors, such as data center layout design and characteristics, ambient weather conditions, rack density, the operation of HAVC systems and their behavior. This complicated connection makes it hard to predict data center energy consumption. With sensor data and information about data center operations, forecasting energy consumption helps in planning and operations of data centers. Well-planned resource provisioning makes good Return on Investment (ROI) and elasticity of computing infrastructures. To take advantage of well-planned resource provisioning, in this thesis, we present energy demand prediction approaches using machine learning.

To predict growth or to set effective goals, it is essential to choose the correct metrics and to be aware of their expressivity and potential limitations. But there are multitude of metrics available to analyze energy efficiency of the data centers. This makes it difficult to select the right metrics to measure the efficiency of a data center. In this thesis, we present an analysis of metrics that are commonly used in data centers, starting from the power grid and going all the way up to the service delivery. Further, data center operators should compare their current approaches with industry standards and assess whether their practices are still valid and/or optimal. Most of the data center operators are not familiar with their existing practices, which hinders on incorporating streamlined new practices. It is important to include best practices into data center operations because baselines are

1.4 Problem Statement, Objectives, and Contributions

created for existing conditions. In this thesis, we partially fill this voids by evaluating different data centers and proposing a set of best practices for sustainable data centers.

This research work is supported by NextGenSmart-DC, a project sponsored by the Ministry of Electronics and Information Technology (MeitY), Government of India and The Netherlands Organization for Scientific Research (NWO), Government of Netherlands under Indo Dutch Science Industry Collaboration (629.002.102). The project's premise is that the application of state-of-the-art ICT tools and techniques can improve the environmental footprint of energy-intensive facilities and deliver Smarter and Greener Data Centers. This project intends to achieve (i) Energy aware job scheduling and load balancing (ii) Optimal trade-offs between energy savings and application performance (iii) Predictive operations and automated decision making (iv) Best practices to identify and to implement measures looking to improve energy efficiency at different levels of data center. By keeping a cyclic approach to research of testing, evaluating and redesigning, this project ensures that the solutions are practically feasible and deliver an actual improvement on the data center efficiency.

1.4 Problem Statement, Objectives, and Contributions

This thesis aims to develop few techniques for optimizing and forecasting energy consumption in data centers for better planning and operations. Further, this thesis aims to develop metrics and a set of guidelines of best practices for sustainable data centers. In this thesis, we investigate the following research problems related to energy efficient data centers:

- **Energy efficient virtual machines placement and selection.** To reduce the number of active servers in a data center, it is necessary to have an efficient virtual machine (VM) placement strategies in place. Determining the optimal placement (allocation) of VMs is an essential aspect of the data center to improve physical resource utilization and to reduce the energy consumption while satisfying the service level agreement (SLA). Determining

1.4 Problem Statement, Objectives, and Contributions

a set of VMs from an over-utilized or under-utilized host has a significant impact on the virtual machine migration time and energy consumption of the data center, and can cause the SLA violation. So, designing a VM selection policy, considering different resources along with CPU utilization plays an important role in improving the energy efficiency of the data centers. The problem consists in determining the best subset of VMs to migrate which would provide the most beneficial system reconfiguration.

- **Forecasting data center energy demand.** Forecasting data center electrical energy demand is very challenging due to highly dynamic nature and complexity of workloads. Developing forecasting models with accurate predictions enable operators to avoid the risk of over-provisioning during non-peak periods, and reduces the risk of under-provisioning in peak periods. It is necessary to have an efficient forecast model for data centers to predict and estimate proper energy demand in real-world situations.
- **Analysis of metrics and practices of data centers.** In order to predict growth or set effective goals, it is important to choose the correct metric and being aware of their expressivity and potential limitations. Understanding and analyzing data center metrics allows the operators to have a better view on possible inefficiencies by focusing on the core parameters. Determining and implementing the best practices for data center operations is required to optimize the workflows and to decrease operating cost in the long term in a data center.

To address these problems, we formulate the following objectives.

- Explore the research on energy-efficient resource management strategies for data centers.
- Develop novel methods for dynamic VM placement and selection.
- Develop efficient methods for forecasting data center energy consumption.
- Explore and analyze the relationship between diverse metrics to measure the efficiency of various data center components.

1.4 Problem Statement, Objectives, and Contributions

- Develop the best practices for sustainable data centers.

The salient contributions of this thesis are as follows:

- Solutions for energy aware virtual machine allocation approaches (see Chapter 3) using
 - i Modified Discrete Particle Swarm Optimization (MDPSO) approach that minimizes the power consumption of the physical machines by estimating the increase in the power consumption before a VM is placed.
 - ii Interactive Particle Swarm Optimization and Genetic Algorithm (IP-SOGA) that performs parallel processing of particle swarm optimization (PSO) and genetic algorithm (GA) using multi-threading and shared memory for information exchange to enhance convergence time and global exploration.
 - iii Imitation Based Optimization (IBO), a swarm based approach for virtual machine placement.
- A novel virtual machine selection method considering Memory, Bandwidth and Size of the Virtual Machines (MBS-VM) (see Chapter 3).
- Solutions for forecasting energy demands of data centers (see Chapter 4) using
 - i Multi-Layer Feed Forward Neural Networks (MFNN).
 - ii Deep learning approach with Parallel Stochastic Gradient Descent training (DPSGD).
- A taxonomy of metrics for sustainable data centers (see Chapter 5).
- Best practices for sustainable data centers (see Chapter 6).

1.5 Thesis Organization

This thesis is structured as follows:

Chapter 2 presents a literature review on virtual machine placement and selection using soft computing techniques and machine learning based algorithms for monitoring data center operations.

Chapter 3 presents three novel approaches (MDPSO, IPSOGA, and IBO) for energy efficient virtual machine placement and a novel virtual machine selection mechanism for cloud data centers. Further, we present a novel virtual machine selection method considering the factors such as memory, bandwidth and size of the VMs (MBS-VM).

The first part of this chapter is published in *Soft Computing, Springer*.

Chapter 4 presents Multi-Layer Feed Forward Neural Networks and Deep learning approach with Parallel Stochastic Gradient Descent for forecasting data center energy demand.

Chapter 5 presents an analysis of metrics that are commonly used in data centers, starting from the power grid and going all the way up to the service delivery. This chapter presents a classification based on the different core dimensions of data center operations such as energy efficiency, cooling, greenness, performance, thermal and air management, network, security, storage, and financial impact. Our work on analysis of metrics is published in *IEEE Transactions on Sustainable Computing*.

Chapter 6 describes a set of best practices to improve the energy efficiency of the data centers which spans the categories of Energy Efficiency, Cooling, Air and Thermal management, Greenness, Storage, and Networks. This chapter is broadly based on the contents of our paper that is accepted in *IEEE IT Professional*.

Chapter 7 summarizes the contributions of the thesis and outlines the future directions.

Chapter 2

Literature Review

Cloud computing provides on-demand access to distributed computing resources on a pay-as-you-go basis. Energy efficiency is becoming an important concern in designing and managing data centers in the era of cloud computing. Therefore, developing energy-efficient resource management techniques for large scale cloud data centers is inevitable. In this chapter, we primarily focus on the state-of-the-art research in virtual machine allocation and selection using soft computing approaches in cloud data centers. Also, we analyze the related works on monitoring data centers using various machine learning approaches.

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

Soft Computing approaches are adaptive mechanisms to facilitate intelligent behaviour in complex and real world problems. The guiding principles of soft computing approaches is to “exploit the tolerance for imprecision, uncertainty, partial truth, and approximation to achieve traceability, robustness, and low solution cost” [17]. Soft Computing approaches are used for solving complex problems such as NP hard for which there are no effective algorithms [18, 19].

In order to reduce the amount of workload in traditional data centers, some specific applications are tied to physical servers and thus making data centers very expensive to maintain. With the advent of virtualization, cloud data centers are

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

providing better services on demand and have become more secure and flexible. In a cloud data center, a physical machine (can also be called as host or server) can host multiple VMs with dynamic work load types and different resource specifications. The servers that host heterogeneous VMs with dynamic and unpredictable workloads can cause an imbalance in violating service level agreements and proper usage of resources. To improve overall system performance and to ensure proper resource utilization in cloud environment, load balancing mechanism distributes and balances the excessive workload ideally across all the nodes. Here, VMs are assigned to suitable hosts and resource utilization is balanced within the hosts. Thus, selecting appropriate algorithms aids in optimal utilization of available resources, enables scalability and there by reducing the resource consumption and response time.

Fig. 2.1 illustrates a typical cloud resource provisioning scenario having user requests, VM and physical machine relationship in a cloud data center. The physical machines represent the available resources like CPU, memory and storage. All the user requests are executed on VMs and may have interdependencies between them. The server virtualization platform makes the physical resource be virtualized and manages the VMs hosted by the physical machines. Each physical machine is allocated to multiple VMs. The VM manager is responsible for placement and consolidation of these virtual machines.

As cloud resource provisioning in data centers is considered as a NP hard problem [20, 21], soft computing approaches are applied by researchers to solve the problem of virtual machine placement and selection, to find the optimal solution with reduced cost and computational time, satisfying the service level agreement. In this section, we classify the virtual machine placement and selection approaches based on the types of soft computing approaches (non-heuristic / heuristic / meta-heuristic) used by the researchers. Non-heuristic approaches guarantee the optimal solution for the problem. But as the size of the instances grows these approaches takes large amount of time to find the optimal solution [22, 23]. Heuristic approaches are generally created by “experience” for specific optimization problems and they intend to find a good solution to the problem by “trail-and-error” in a acceptable amount of time. The solutions may not be the best or optimal solution but they might be better than an educated guess [24]. Meta-heuristic approaches are

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

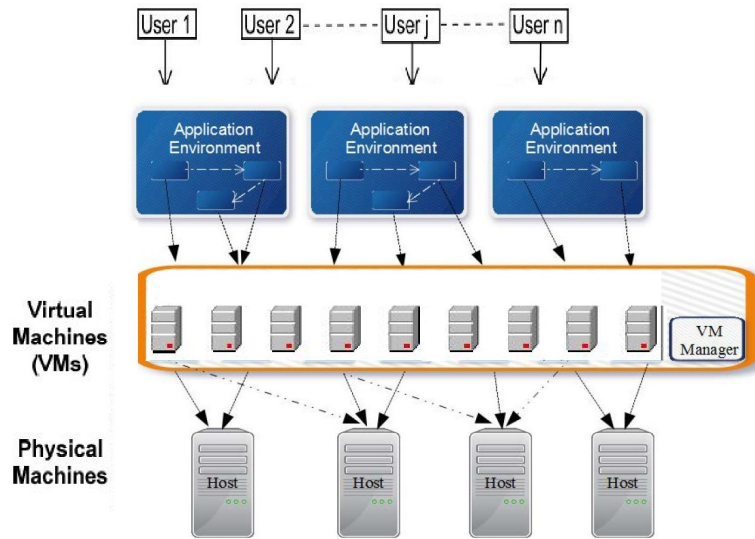


Figure 2.1: A Typical Cloud Resource Provisioning Scenario

higher level heuristic designed to find, generate, or select a heuristic (partial search algorithm) that may provide a sufficiently good solution, with incomplete information. Meta-heuristics sample a set of solutions that is too large to be completely sampled. These approaches may make few assumptions about the optimization problem being solved, and so they may be usable for a variety of problems.

2.1.1 Non-heuristic Approaches

Portaluri et al. [25] proposed energy-aware dynamic allocations of virtual machines considering computational and network requirements using multi resource Best-fit algorithm. Zhang et al. [26] developed energy efficient VM selection algorithms for overloaded hosts based on dynamic programming and greedy algorithm. Garg et al. [27] considered different Quality of Service (QoS) requirements of workloads and proposed a virtual machine scheduling algorithm to maximize the resource utilization and profit. Shen et al. [28] developed a novel application-aware bandwidth guarantee framework and a VM migration algorithm considering the variations in network demands.

Tseng et al. [29] solved the network-aware VM placement optimization problem using integer linear programming, to minimize communication time for VMs

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

of the same type. Goudarzi et al. [30] created multiple copies of VMs and put them on different servers. They used dynamic programming to minimize the energy consumption by switching off the underutilized servers. However, in this approach, SLA violation is not considered and network overhead is created by multiple copies of VMs. Addya et al. [31] proposed a coalition-based cooperative structure to compute the pricing that users pay for their requested VMs and used Integer Linear Programming for energy-aware virtual machine placement. Zhang et al. [32] proposed a heuristic design for virtual machine placement using the least resource wastage and the least power consumption approaches. Quang-Hung et al. [33] proposed energy-aware and performance-per-watt oriented Best-fit algorithm to choose the physical machine that has maximum performance (in terms of performance-per-watt) to assign a VM. Li et al. [34] proposed an approach based on multi-dimensional space partition model for efficiently placing virtual machines to improve the resource utilization of data centers. However, this model has not considered SLA and VM migration cost. Wang et al. [35] presented a decentralized double threshold VM selection policy considering the utilization of the physical hosts. However, they have not considered the energy consumption and this policy may not give the optimal solution always. Mansouri et al. [36] proposed a system to optimize cost using object replication and migration across cloud service providers while maintaining the latency threshold for the application. They used an optimal offline algorithm and a Receding Horizon Control (RHC) technique to make a trade-off between residential and migration costs.

Table 2.1 presents a set of non-heuristic approaches for virtual machine placement and selection in data centers, considering energy efficiency and / or utilization of single / multiple resources.

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

Reference	Methodology	Considered	
		Energy ?	Multiple Resources ?
Mansuri et al. (2017) [36]	Randomized online algorithm	yes	no
Portaluri et al.(2017) [25]	Multi Resource Best Fit	no	yes
Tseng et al.(2017) [29]	Integer Linear Programming	no	yes
Addya et al. (2017) [31]	Queuing Model	yes	no
Zhang et al. (2017) [26]	Dynamic Programming	yes	no
Xia et al. (2017) [37]	Mixed Integer Programming	yes	yes
Cui et al. (2017) [38]	Progressive-decompose-rounding	no	no
Rimal et al. (2016) [39]	Easy Backfilling, Minimum Completion Time	no	no
Imlai et al.(2016) [40]	Elastic Scheduling	yes	no
Shen et al. (2016) [28]	Graphcut Algorithm	no	no
Verma et al. (2016) [41]	Best-Fit Algorithm	no	no
Wang et al. (2016) [42]	Mixed integer programing (MIP)	no	no
Zhou et al. (2016) [43]	Adaptive Three-Threshold Algorithm	yes	no
Agrawal et al. (2015)[44]	Linear programming	yes	no
Zhang et al. (2015) [32]	Evolution approximation algorithm	yes	no
Quang-Haung et al. (2014)[45]	MinDFT-ST and MinDFT-FT	no	yes
Quang-Hung et al. (2014) [33]	Energy-aware and Performance per watt oriented Best fit	no	yes
Wang et al. (2013) [35]	Two-threshold decentralized migration	no	no
Singh et al. (2013) [46]	Bankers algorithm	no	no
Li et al. (2013) [34]	EAGLE	yes	no
Beloglazov et al. (2012) [47]	Optimal Online Deterministic Algorithm, Random Choice, Maximum correlation,	yes	no
Goudarzi et al. (2012)[30]	Dynamic programing	yes	no
Bose et al. (2011) [48]	Deterministic approach	no	no
Le et al. (2011) [49]	Binary Integer Programming	yes	no
Beloglazov et al. (2010) [50]	Minimum migration time (MMT)	yes	no
Buyya et al. (2010) [51]	Single Threshold, highest potential growth and random choice policy	yes	no
Meng et al. (2010) [52]	Two tier Cluster and Cut algorithm	no	no
Dhiman et al. (2009)[53]	vGreen	yes	no
Cardosa et al. (2009) [54]	PowerExpand MinMax	yes	no
Verma et al. (2008) [55]	pMapper	no	no
Nathuji et al . (2007) [56]	Virtual power management	yes	yes

Table 2.1: Non-heuristic Approaches

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

2.1.2 Heuristic Approaches

Lee et al. [57] proposed a virtual machine consolidation approach considering the competing resource demands in multiple dimensions. They developed two heuristic approaches namely Dot-Product and Norm-based Greedy to select the next VM. Xiao et al. [58] measured the unevenness in the multi-dimensional resource utilization of a server using skewness. They developed an algorithm to capture the rising trend of resource usage patterns to minimize the skewness and to increase the overall resource utilization. Wang et al. [59] presented a heuristic min-cost algorithm to solve energy consumption problems in data centers, resolving difficulties with integer decision variables and non-linearity of the power model.

Grange et al. [60] proposed an Attractiveness-Based Blind Scheduling algorithm using as a greedy heuristic approach for choosing a definitive placement for a task at the time of its submission. This algorithm compares multiple possible placements. Wang et al. [59] solved VM placement problem using a mixed integer programming approach. This approach takes more time for large number of virtual machines and physical machines.

Zhang et al. [66] presented an approximate approach based on bin packing algorithm to migrate virtual machines. They considered the resource utilization and the migration cost to get the optimal solution. Cao et al. [84] proposed an energy-aware heuristic framework for VM consolidation to achieve a better energy-performance trade-off. Most of the works apply greedy heuristics to model VM consolidation as variants of the bin packing problem such as First Fit Decreasing (FFD) [85], Best Fit [86], Best Fit Decreasing [77], and so on [87, 88].

For reducing power consumption and SLA violations, Mostafa et al. [62] considered maximum absolute deviation for VM placement by grouping both physical and virtual machines. An extended Best Fit Decreasing algorithm was proposed for reducing the number of virtual machines and a learning automata was used as a trade-off between power consumption and SLA violations. Wood et al. [85] presented a system called sandpiper that automates monitoring tasks and helps in detecting hotspots, determines a mapping between physical and virtual resources and initiates the necessary migrations in a virtualized environment. Soomro et al. [61] presented two novel heuristic algorithms based on First Fit Decreasing (FFD) for VM placement. FFD preprocesses all VMs and sorts them in descending order

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

Reference	Methodology	Considered	
		Energy ?	Multiple Resources ?
Soomro et al. (2017) [61]	First-fit Decreasing	yes	yes
Ghobaei et al. (2017) [62]	Best-fit Decreasing	yes	no
Zhou et al. (2017) [63]	Network topology aware redundant VM placement approach	no	no
Zhou et al. (2016) [43]	Energy-Aware Best Fit Decreasing	no	no
Xiao et al. (2016) [64]	Game Theory	no	no
Verma et al. (2016) [41]	Best-fit Heuristic	no	no
Su et al. (2015)[65]	Iterative heuristic search	yes	no
Zhang et al. (2015) [66]	Bin packing	no	yes
Liang et al. (2014) [67]	App_VM_Reconfiguration	yes	yes
Dai et al. (2014) [68]	Minimum power VM placement	yes	no
Fang et al. (2013) [69]	Greedy Bin-Packing	yes	yes
Huang et al. (2013) [70]	Opportunity cost based heuristic approach	yes	no
Dong et al. (2013) [71]	Best Fit with hierarchical clustering	yes	no
Wang et al. (2013) [72]	Max-Min Multidimensional Stochastic Bin Packing	yes	no
Gupta et al. (2013) [73]	Heuristic bound approach	no	no
Dong et al. (2013) [74]	Power aware best fit decreasing	yes	no
Beloglazov et al. (2012) [47]	Power Aware Best Fit Decreasing Algorithm	yes	no
Chen et al. (2012) [75]	Cost-aware two-phase heuristic algorithm	yes	no
Somani et al. (2012) [76]	Bin Packing	yes	yes
Beloglazov et al. (2012) [77]	Modified Best Fit Decreasing Algorithm	yes	no
Jin et al. (2012) [78]	Max-Min multidimensional Stochastic bin packing	no	yes
Le et al. (2011) [49]	Approximation algorithms based on greedy formulations	yes	no
Simarro et al. (2011) [79]	Dynamic cost-aware load distribution	yes	no
Machida et al. (2010) [80]	Multiple k-redundancy	no	yes
Younge et al. (2010) [81]	Greedy Algorithm	yes	no
Li et al. (2009) [82]	EnaCloud	yes	yes
Bobroff et al. (2007) [83]	Binpacking and Timeseries forecasting	yes	no

Table 2.2: Heuristic approaches

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

based on their sizes. Zhou et al. [63] presented a VM placement optimization approach to enhance the reliability of services in a cloud environment. The best set of VM and hosting servers were selected using network topology. Further, a heuristic was used to solve the problem of VM reassignment optimization. Fang et al. [89] presented a method to reduce job delay by assigning virtual machines to a few hypervisors and by moving the communicating parties to closer locations. Su et al. [90] addressed the problem of VM placement by considering affinity and conflict between virtual machines. They analyzed the impact of affinity, and conflict between virtual machines.

Table 2.2 presents a list of heuristic approaches for virtual machine placement and selection in data centers, considering energy efficiency and / or utilization of single / multiple resources.

2.1.3 Meta-heuristic Approaches

Zhao et al. [91] explored the balance between server power savings and virtual machine performance. They proposed an algorithm based on ant colony optimization, to minimize server power consumption and guarantee VM performance. Duan et al. [92] proposed a PreAntPolicy that consists of a prediction model based on fractal mathematics and a scheduler based on an improved ant colony algorithm. Based on load trend prediction, the model executes the scheduler while minimizing energy consumption. Tang et al. [93] considered energy consumption of networks and servers and proposed a genetic algorithm and a hybrid genetic algorithm for virtual machine placement. Zheng et al. [94] proposed a Biogeography-Based Optimization (BBO) technique that studies the geographical distribution of species migration and extinction of existing species and rise of new species for VM placement. They consider the resource wastage and the power consumption while placing the virtual machine on to physical machines. Wang et al. [95] proposed an improved particle swarm optimization (PSO) to minimize energy consumption during the provision of data-intensive services with a global QoS guarantee in a data center.

Liu et al. [96] proposed an Ant Colony Optimization (ACO) algorithm for virtual machine placement. They combined ACO with a local search technique namely Order Exchange and Migration (OEM). This method effectively minimizes

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

the energy consumption by reducing the active physical machines. Sawant et al. [97] proposed a genetic algorithm based scheduling for placing virtual machines in a cloud environment. They used the average load of each virtual machine on a host in a time cycle to find the mapping solution satisfying the total load variation (σ). Goiri et al. [98] considered the problem of multifaceted resource management in data centers, considering the cost of energy consumption, SLAs, outsourcing capabilities, heterogeneity management, and economic modeling. They proposed an algorithm that tries to find those combinations from the scoring matrix that maximizes the overall system benefit. However, they confined the number of movements per round because this algorithm has a chance to enter into a periodic cycle without converging.

Dashti et al. [109] presented a modified PSO approach for allocation of migrated virtual machines from overloaded hosts to improve energy efficiency in a cloud computing environment. But placing migrated VMs again in other hosts may lead to aggressive migration which causes high energy consumption and leads to SLA violations in a cloud computing environment. Wu et al. [122] proposed a genetic algorithm for virtual machine placement considering the energy consumption of servers and communication network in a data center. Kumar et al. [110] focused on minimizing the total resource wastage with efficient VM allocation using PSO in cloud operations.

Virtual machine migration is a major way for reducing unnecessary consumptions in a data center. Palmieri et al. [127] proposed a swarm based meta-heuristic to the resource scheduling and balancing problem, minimizing runtime and fairly balancing load in next generation grids. However, they have not considered the energy consumption of computing resources of local grid nodes. Ferdous et al. [128] proposed Ant Colony Optimization technique to consolidate the virtual machines focusing on balancing different computing resources with the goal of minimizing resource wastage and power consumption. Kansal et al. [105] proposed a virtual machine migration technique using Firefly algorithm. This approach has the capability to handle multiple modes and migrates the virtual machines based on their utilization levels. A virtual machine having high utilization is migrated to the under-utilized host. Efficient re-optimization strategy for big data access in multi-tenant cloud infrastructures. Palmieri et al. [129] developed a meta-heuristic based on Greedy Randomized Adaptive Search Procedure (GRASP) for path rerouting

2.1 Virtual Machine Placement and Selection using Soft Computing Approaches

Reference	Methodology	Considered	
		Energy ?	Multiple Resources ?
Ding et al. (2018) [99]	Discrete firefly algorithm	yes	no
Zhao et al. (2018) [91]	ACO	no	no
Aryania et al. (2018) [100]	ACO	yes	yes
Duan et al. (2017)[92]	Improved ACO	no	yes
Kaur et al. (2017) [101]	Gravitation algorithm	yes	no
Chen et al. (2017) [102]	Constrained immune memory and immunodominance clone optimization	yes	no
Zheng et al. (2016) [94]	Bi-geography based optimization	no	yes
Wang et al. (2016) [95]	Improved PSO	yes	no
Portaluri et al. (2016) [103]	Fuzzy Logic based optimization	yes	yes
Mu et al. (2016) [104]	PSO with Gauss Strategy	no	no
Kansal et al. (2016) [105]	Firefly optimization	yes	no
Liu et al. (2016) [96]	ACO with OEM	yes	no
Gao et al. (2016) [106]	ACO	no	no
Wang et al. (2016) [107]	Integer bi-level genetic Algorithm	yes	yes
Xu et al. (2015) [108]	Improved PSO	yes	yes
Dashti et al. (2015) [109]	Modified PSO	yes	no
Kumar et al. (2015) [110]	PSO	yes	no
Joshi et al. (2015) [111]	Cuckoo search	yes	no
Tang et al. (2015) [112]	Hybrid Genetic Algorithm	no	no
Luo et al. (2014) [113]	Hybrid shuffled frog leaping	yes	no
Kruekaew et al. (2014) [114]	Artificial Bee Colony (ABC)	no	no
Dong et al.(2014) [115]	ACO	no	no
Yang et al. (2014) [116]	NSGA	yes	no
Gao et al. (2013) [117]	Mean variance optimization	yes	yes
Wang et al. (2013) [118]	PSO based local fitness first	yes	no
Ma et al. (2012) [119]	ACO	yes	yes
Wu et al. (2012) [120]	Simulated annealing	yes	no
Goudarzi et al. (2012) [121]	Semi-static optimization, Dynamic optimization	yes	no
Wu et al. (2012) [122]	Genetic Algorithm	yes	no
Jeyarani et al. (2011) [123]	Self Adaptive PSO	yes	no
Mark et al. (2011) [124]	GA and PSO	yes	yes
Xu et al. (2010) [125]	Improved GA	yes	yes
Piao et al. (2010) [126]	Greedy randomized adaptive search (GRASP)	no	no

Table 2.3: Meta-heuristic approaches

2.2 Machine Learning Approaches for Data Center Monitoring

and VM migration in a federated cloud. Several researchers used variants of ACO to address VM consolidation and have shown better results [130, 131, 132, 133].

Table 2.3 presents a list of meta-heuristic approaches for virtual machine placement and selection in data centers, considering energy efficiency and/ or utilization of single / multiple resources.

2.2 Machine Learning Approaches for Data Center Monitoring

Energy consumption modeling of the data centers quantifies the energy consumption as a function of input parameters. These models aim at forecasting the energy consumption of the individual data center to determine the energy supply requirements at various levels. Accurate modeling and predictions of data center energy consumptions enables various energy management opportunities such as: estimating energy supply, early stage design decisions, estimating improvements to data center energy performance, and energy infrastructure planning [134, 135]. Machine learning has the distinct advantage that distilled expertise from other disciplines. Machine learning (ML) algorithms have achieved remarkable successes in solving many complex tasks by learning from raw data. Therefore, the integration of machine learning in data centers has sparked great interest in recent years [136].

Verma et al. [41] proposed a dynamic resource demand prediction and allocation framework by classifying service tenants according to their changing resource requirements in multi-tenant service clouds. They applied exponential moving average for short term predictions and polynomial regression, ARX, and ARMAX for long term predictions. This model is useful for preparing the correct type of VMs in advance but does not consider energy consumption. Dong et al. [74] proposed a forecast based power aware best fit decreasing model that uses the estimated resource demand and places the virtual machines using a heuristic method. Broff et al. [83] proposed a dynamic server migration and consolidation algorithm. They used time series forecasting for forecasting resource demands and proposed a method for classifying workload signatures to identify efficient servers.

Krioukov et al. [137] proposed various forecasting models to predict the future user requests and adjusted the servers accordingly in a heterogeneous environment.

2.2 Machine Learning Approaches for Data Center Monitoring

Abdelzaher et al. [138] use control theory to design a feedback control framework augmented by elements of scheduling and queuing theory for resource provisioning. They proposed a closed-loop approach to manage the trade-off between the number of applications hosted on a machine and the amount of resources (CPU or memory) allocated to each application.

Liang et al. [67] proposed a modified Index Curve Model to predict the application requests. To avoid the time delay for VM allocation, they separated reconfiguration and real allocation. Kousiouris et al. [202] studied the effect of different critical parameters on the performance of VMs. Further, they used the optimized artificial neural networks to predict the performance degradation a priori to the execution. This helps in pro-actively provisioning the resources as needed by an application. Xu et al. [185] presented a unified reinforcement learning approach for predicting the demand and then automated the configuration processes of VMs. However, this work focuses on CPU and memory resources and does not consider network-I/O and disk-I/O bandwidth. Bankole et al. [203] developed a prediction model for TPC-W benchmark web application. They applied neural network (NN), linear regression (LR), and support vector regression (SVR) to forecast the future resource usage in multi-tier web applications.

Do et al. [204] presented a novel application profiling technique using canonical correlation analysis (CCA). CCA is used to identify the resources that affect the workload behavior and to find a suitable host. Further, a performance prediction model is developed based on the application profiles generated using CCA. Gong et al. [205] presented a model to predict the dynamic patterns in resource demands and adjusted their resource allocations automatically. They have used Fast Fourier Transform (FFT) to identify the dominant resource usage patterns. Wood et al. [206] proposed an autonomous model to correlate the resource usage on native environment. The model achieved a higher precision in estimating resources needed to an application that is migrating from physical to virtual environment. Du et al. [207] proposed models based on artificial neural network and regression to predict the interference among virtual machines to manage the resources efficiently. Elprince et al. [208] proposed an autonomous cloud management model that predicts the workload demand and places virtual containers accordingly. For improved utilization of resources, Suhad et al. [169] proposed the grouping of tasks based on their similar resource requirements. This task grouping

2.2 Machine Learning Approaches for Data Center Monitoring

	Problem Handled	Approach
Wang et al. (2018) [139]	Classifying traffic flows	Naive Bayes Discretization
Yu et al. (2018) [140]	Traffic prediction	Neural networks
Shoukourian et al. (2017) [141]	Coefficient of Performance (COP) prediction	LSTM
Nishi et al. (2017) [142]	Power variation trend	Support vector machine
Liu et al. (2017) [143]	Early prediction of job failures	On line extreme learning machine
Kim et al. (2017) [144]	Detect inactive VMs	Linear support vector machine
Ahmed et al. (2017) [145]	Automated detection of the performance faults	Machine learning approach based on correlation stability violations
Bashar et al. (2017) [146]	Admission Control of service requests	Bayesian Networks based predictive modeling framework
Son et al. (2017) [147]	Migration	Fuzzy-logic-based learning
Jobava et al. (2017) [148]	VM clustering	Learning automata
Sidhu et al. (2016) [149]	Identify different nodes with similarities	Bayesian networks
Moskalenko et al. (2016) [150]	Moment of reduction in the functional efficiency	Information-extreme machine learning
Verma et al. (2016) [41]	Resource demand prediction	ARX, ARMX
Ukdave et al. (2016)[151]	Detecting interference between applications	Collaborative filtering
Shen et al. (2016) [152]	Identifying unproductive instances.	Decision Tree Verification
Sieber et al. (2016) [153]	Estimate hypervisor resources	Weighted orthogonal distance regression
Tarutani et al. (2016) [154]	Predicting Temperature Distribution	Machine-learning
Hieu et al. (2015) [155]	VM consolidation	Original least squares
Masoumzadeh et al. (2015) [156]	Dynamic Consolidation	Fuzzy Q-learning
Dabbagh et al. (2015) [157]	Estimation of the number of physical machines	Wiener Predictor
Fallah et al. (2015) [158]	High energy consumption	Learning automata
Liang et al. (2014) [67]	VM placement	Modified Index Curve Model
Versick et al. (2013) [159]	Power consumption estimation	Auto-Adaptive Resource Allocation
Dong et al. (2013) [74]	VM placement	Forecast based power aware best fit decreasing
Berral et al. (2013) [160]	Resource usage	M5P regression tree
Sato et al. (2013) [161]	Resource usage prediction	Autoregressive Model
Vasic et al. (2012) [162]	Classifying workloads and estimating interference index	K-means clustering
Yuan et al (2012) [163]	Uncertain task flow	reinforcement learning
Li et al. (2011) [164]	Thermal forecasting model	Autoregressive moving average
Bobroff et al. (2007) [83]	VM placement	Time series forecasting
Zhang et al. (2007) [165]	Resource requirements of applications	Regression-based model
Moore et al. (2006) [166]	Cooling and heat management	Weatherman

Table 2.4: Machine Learning approaches for data center monitoring

2.2 Machine Learning Approaches for Data Center Monitoring

Table 2.5: Machine Learning approaches for data center monitoring

	Problem Handled	Approach
Arif et al. (2017) [167]	Live virtual machine migration	Decision tree algorithm
Kim et al. (2017) [144]	Detecting inactive VMs	Support vector machine
Zhang et al. (2017) [168]	Switch failure detection	Hidden Semi-Markov Model
Yousif et al. (2017) [169]	Clustering the workload	K-means clustering and density based clustering
Li et al. (2017) [170]	Optimize the control policy and performance	Deep-Q-Network
Nakamura et al. (2016) [171]	Workload misplacement	Cooling control model
Shaw et al. (2016) [172]	Energy consumption	Reinforcement Learning
Verma et al. (2016) [41]	Resource demand prediction	Polynomial regression, auto-regressive with external input (ARX), and auto-regressive moving average with external input (ARMAX)
Dabbagh et al. (2015) [173]	VM placement and Energy Consumption Prediction	Release-Time Aware heuristic Approach
Hien et al. (2015) [155]	Resource usage prediction	Original least square and Linear Regression
Tseng et al. (2015) et al. [174]	VM migration	Support vector machine
Tang et al. (2014) [175]	Future Memory prediction	Auto regression Model
Shoukourian et al. (2014) [176]	Power consumption prediction	Adaptive Energy and Power Consumption Prediction
Farahnakian et al. (2014) [177]	Dynamic Consolidation	Reinforcement Learning
Mijumbi et al. (2014) [178]	Virtual network resource allocation	reinforcement learning algorithm
Liu et al. (2014) [179]	workload prediction	Time Series Analysis
Dong et al. (2013) [180]	CPU utilization prediction	Auto Regressive Integrated Moving Average model (ARIMA), ExponentialSmoothing State Space model (ETS), Random Walk model (RW), and Structural Time Series model (STS)
Ramezani et al. (2013) [181]	Predicting VM workload patterns	neural networks and fuzzy expert systems.
Sato et al. (2013) [161]	Resource usage prediction.	Autoregression Model
Canali et al. (2013) [182]	Clustering VMs	Bhattacharya distance
Zhao et al. (2013) [183]	Workload classification	Support Vector Machine and KNN
Farahnakian et al. (2013) [184]	CPU utilization prediction	Linear Regression
Xu et al. (2012) [185]	Automate the configuration process	Unified Reinforcement Learning
Fang et al. (2012) [186]	Cloud Resource Prediction and Provisioning scheme	ARIMA model for our prediction algorithm
Moreno et al. (2012) [187]	Resource utilization patterns	Architected Neural Network.
Kundu et al. (2012) [188]	Predict the performance of a VM-hosted application	Artificial neural network, support vector machine
Wang et al. (2011) [189]	Temperature prediction in Data Centers	Artificial neural networks
Roy et al. (2011) [190]	Workload forecasting	Auto-regressive moving average
Kousouris et al. (2011) [191]	Interference prediction	Artificial Neural Network
Berral et al. (2011) [192]	CPU utilization	Reinforcement learning and Q Learning Algorithms
Hu et al. (2011) [193]	Dirty Page Detection	Time-series based precopy algorithm
Kraft et al. (2011) [194]	Degradation of disk	Trace-driven approach
Dhiman et al. (2010) [195]	Power prediction	Gaussian mixture model
Akoush et al. (2010) [196]	Live migration	Average page dirty rate Model
Berral et al. (2010) [197]	Predict the power consumption	Linear regression, M5P
Rao et al. (2009) [198]	VM configuration	Reinforcement Learning
Choi et al. (2008) [199]	CPU Utilization	History matrix based Prediction
Tang et al. (2006) [200]	Temperature Prediction	Cross interference method
Ganesh et al. (2006) [201]	Capacity predictor	Trend Analysis

enables a cloud data center to identify an optimal VM placement strategy and then tries to allocate those VMs from the complemented groups or clusters on the same physical machine. This type of placement reduces competition among resources within the same physical machine. Keller et al. [209] considered network affinity while migrating virtual machines. They proposed live migration of ensembles, an algorithm to perform incremental migration of data-plane state and other traffic sources synchronously.

Table 2.4 presents a list of machine learning approaches for data center monitoring.

2.3 Analysis and Discussion

We classified the virtual machine placement and selection approaches based on non-heuristic, heuristic, and meta-heuristic approaches. Based on our literature survey, from Figure 2.2, we observe that 40% of the research papers focus on non-heuristic approaches, 34% of the research papers focus on heuristic approaches, and 26% of the research papers focus on meta-heuristic approaches. Researchers used non-heuristic methods to solve VM placement problem with global constraints. However, these algorithms support limited workflow and do not guarantee optimal solution. Heuristic approaches like First-fit decreasing, Best-fit decreasing, etc. are deterministic in nature and do not always guarantee the optimal solution. Meta-heuristic methods support large workflow sizes with global constraints and have less computation time. However, improving the performance, exploration and exploitation capacities of the meta-heuristic algorithms is still a challenging issue. It seems promising to propose novel meta-heuristics methods to solve VM placement and selection. We also observe that in recent years, these meta-heuristic approaches are receiving considerable attention from the research community. Various meta-heuristic approaches proposed by researchers for VM placement and selection are shown in Figure 2.3.

Based on our literature survey, we identify that workload demand prediction, resource utilization prediction, and power consumption/variation prediction are the major problems related to data center monitoring, solved using machine learning approaches. Fig. 2.4 presents various data center monitoring issues solved using machine learning. Fig. 2.5 illustrates the solutions for the said problems (in

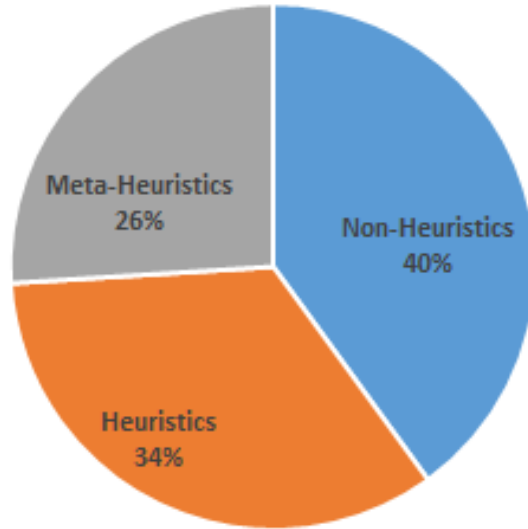


Figure 2.2: Classification of Approaches

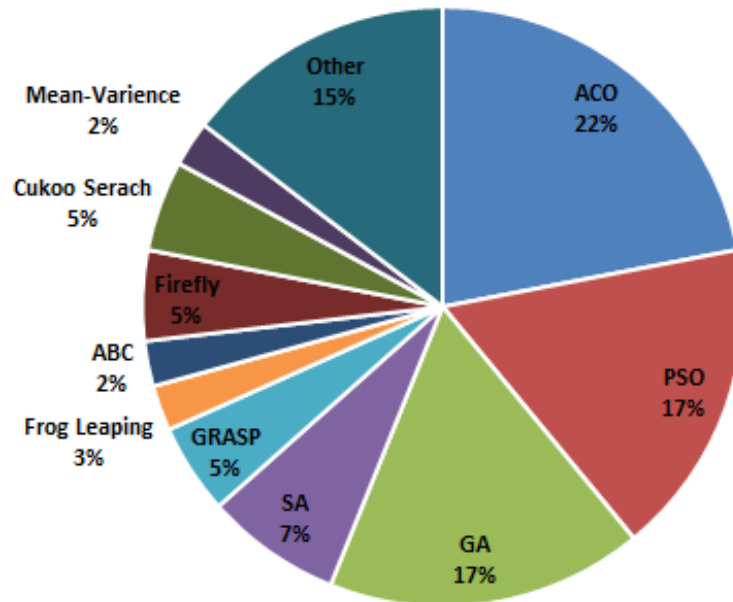


Figure 2.3: Meta-heuristic Approaches

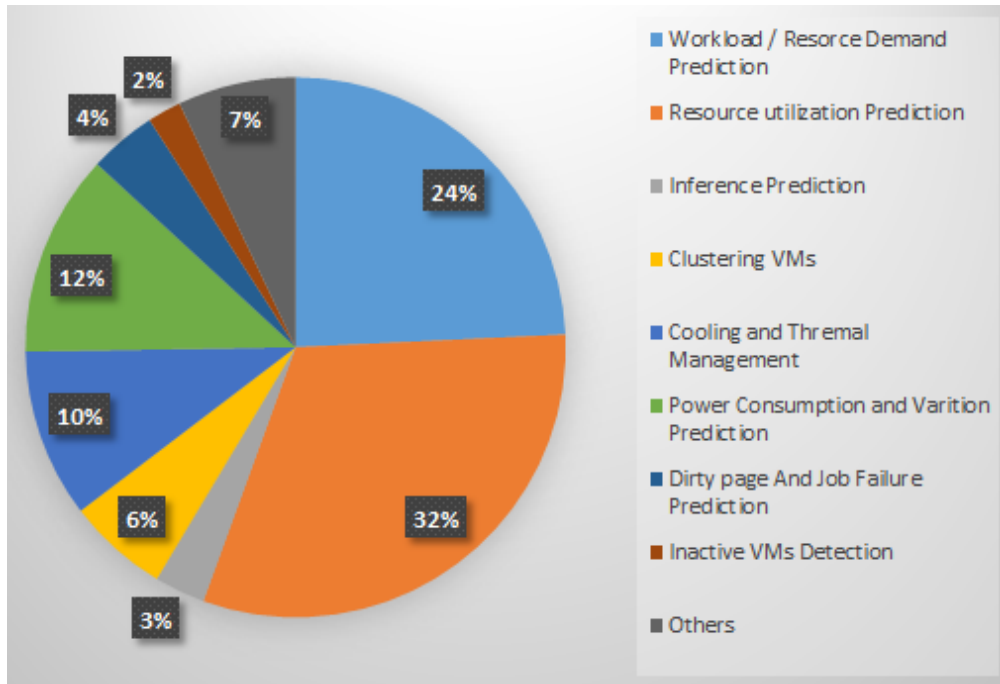


Figure 2.4: Machine learning applications for data center monitoring

Fig. 2.4) using various machine learning approaches. It can be observed that Artificial Neural Networks (ANN) (12%), Auto Regressive Moving Average (ARIMA) (10%), Reinforcement Learning (RL) (9%), and Support Vector Machines (SVM) (10%) are mostly used by the researchers.

2.4 Justification for Present Work

This chapter's contribution has been to reveal the state-of-the-art research on virtual machine placement and selection in cloud data center using soft computing approaches and on monitoring data centers using machine learning approaches. Based on literature review, we observe that VM placement and selection problems are solved generally by heuristics, non-heuristics, and meta-heuristic approaches. Analyzing from the aspect of consideration of energy and utilization of multiple resources for VM placement and selection, we observe that there is a scope to propose meta-heuristic approaches for solving VM placement and selection in cloud data centers. Further, we observe that forecasting energy consumption in data centers

2.4 Justification for Present Work

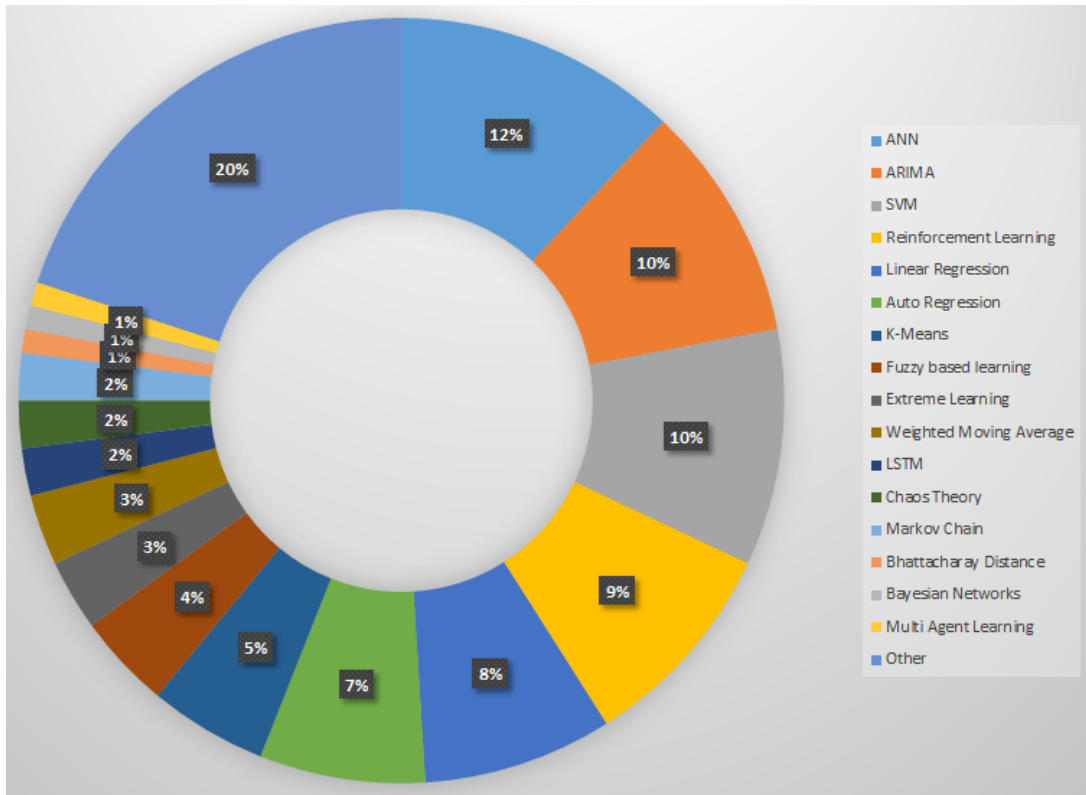


Figure 2.5: Various machine learning approaches for data center monitoring

using machine learning and deep learning techniques help data center operators to predict and estimate the accurate energy demand in real world situations.

Chapter 3

Energy aware Virtual Machine Placement and Selection Approaches in Cloud Data Centers

The placement of virtual machines over physical hosts becomes a vital component of any cloud management framework which needs to be optimized continuously due to the variability of workloads experienced by modern applications. In this chapter, we formally define the virtual machine placement problem and describe three novel optimization algorithms namely “Modified Discrete Particle Swarm Optimization (MDPSO)”, “Interactive PSO-GA”, and “Imitation Based Optimization (IBO)” for dynamic VM placement. Further, we present a VM selection algorithm for live migration to improve the energy efficiency of the data centers. To understand the implications of the proposed approaches, we present a comparative analysis of the proposed algorithms.

3.1 Modelling Resource Allocation and Power Consumption in Cloud Data Centers

Let M be the number of virtual machines and N be the number of physical hosts / physical machines (PM) respectively. V is defined as a set consisting of virtual machines, where v_i is an instance of a virtual machine and $|V| = M$. P is a collection of physical machines, where p_j is an instance of a physical machine and $|P| = N$.

$$V = \{v_1, v_2, \dots, v_k\} \quad (3.1)$$

$$P = \{p_1, p_2, \dots, p_t\} \quad (3.2)$$

In our model, each physical machine (p_j) is characterized as a 5 tuple:

$$p_j = (id_j, cpu_j, storage_j, bw_j, cores_j)$$

where id_j is the unique identity of the physical machine, cpu_j is the computing power of the physical machine generally given in million instructions per second (MIPS), $storage_j$ is the capacity of the random access memory of the physical machine, $cores_j$ is the number of cores in a physical machine and bw_j is the bandwidth that is allocated to a physical machine. Each virtual machine is characterized by a 4 tuple, given by:

$$v_i = (id_i, cpu_i, storage_i, bw_i)$$

where id_i gives the unique ID of the VM. cpu_i , $storage_i$, and bw_i are the quantity of processing power, memory, and bandwidth requested by the VM respectively. Generally, a physical machine can host one or more virtual machines but each virtual machine must be assigned to only one physical machine. We use a variable y_{ij} to indicate whether a virtual machine is allocated to a physical machine or not. y_{ij} is 1, if i^{th} virtual machine is allocated to the j^{th} physical machine and 0 otherwise.

3.1.1 Problem Definition

We aim to derive a mapping from the set of virtual machines (V), to the set of physical machines (P) that should maximize the resource utilization and minimize

3.1 Modelling Resource Allocation and Power Consumption in Cloud Data Centers

the energy consumption. Resource utilization of a physical machine can be decomposed into cpu utilization, memory utilization and bandwidth utilization by all the virtual machines in that host. v_i^{cpu}/p_j^{cpu} gives the cpu utilization of i^{th} virtual machine (v_i) in the j^{th} host (p_j). v_i^{mem}/p_j^{mem} gives the memory utilization of i^{th} virtual machine (v_i) in the j^{th} host (p_j), and v_i^{bw}/p_j^{bw} gives the bandwidth utilization of i^{th} virtual machine (v_i) in the j^{th} host (p_j). The objective of maximizing physical resource utilization is defined as follows:

$$\text{Maximize} \left(\frac{v_i^{cpu}}{p_j^{cpu}}, \frac{v_i^{mem}}{p_j^{mem}}, \frac{v_i^{bw}}{p_j^{bw}} \right)$$

The allocation must satisfy the following constraints :

$$\forall i \sum_{j=1}^m y_{ij} = 1 \quad (3.3)$$

$$\forall j \sum_{i=1}^n y_{ij} \cdot v_i^{cpu} \leq p_j^{cpu} \quad (3.4)$$

$$\forall j \sum_{i=1}^n y_{ij} \cdot v_i^{mem} \leq p_j^{mem} \quad (3.5)$$

$$\forall j \sum_{i=1}^n y_{ij} \cdot v_i^{bw} \leq p_j^{bw} \quad (3.6)$$

where y_{ij} is a boolean variable. $y_{ij} = 0$ indicates that the virtual machine (v_i) is allocated to a physical machine (p_j) or not. Equation 3.3 ensures that one virtual machine is allocated to only one physical machine though one physical machine may have more than one virtual machines. Equations 3.4, 3.5, and 3.6 check the resource requests for each physical machine and guarantee that they will not exceed the capacity of a physical machine. It is possible to define the upper and lower utilization limits for a physical machine to ensure reliability.

Generally, we find that the increasing energy consumption is primarily caused by idle or under-utilized servers. Modeling server power consumption based on CPU utilization is traditionally followed in predicting the power consumption for fixed workloads [210]. Hence, we model the energy consumption of a physical machine based on CPU utilization. The cpu utilization of a physical machine (p_j) for a given virtual machine (v_i) is represented as follows:

3.1 Modelling Resource Allocation and Power Consumption in Cloud Data Centers

$$p_u^{j(i)} = \begin{cases} \left(\frac{v_i^{cpu}}{p_j^{cpu}} \right) * 100 & \text{if } y_{ij} = 1 \\ 0 & \text{otherwise.} \end{cases} \quad (3.7)$$

The physical machine utilization is the sum of the utilization rates of all VMs in that physical machine (p_j), given by

$$p_j^u = \sum_{i=1}^m p_u^{j(i)} \quad (3.8)$$

We used SPEC proven methodologies for iteratively measuring the server performance and adopted them to calculate energy consumption [211]. The power consumption of a physical machine (p_j) with utilization ‘ u ’ is defined as $E(p_j(u))$ and calculated according to Equation 3.9.

$$E(p_j(u)) = a_i + (a_{i+1} - a_i) * (10 * u - i) \quad (3.9)$$

where $i = \lfloor 10 * u \rfloor$ and $(a_{i+1} - a_i)$ represents the increase in energy consumption of the server when the utilization is increased from i to $i + 1$.

3.1.2 Fitness Evaluation

Let ‘ u ’ be the current utilization of a physical machine (p_j). A virtual machine v_i will be placed in p_j if the physical machine satisfies the conditions in Equations 3.3, 3.4, 3.5, and 3.6. The increase in the energy consumption of the physical machine is calculated using Equation 3.10.

$$\delta P_j^i = E(p_j(u1)) - E(p_j(u)) \quad (3.10)$$

where ‘ $u1$ ’ is the utilization of the physical machine after v_i is placed in p_j and δP_j^i is the increase in the energy consumption of the physical machine p_j if the virtual machine v_i is allocated to the physical machine p_j .

We minimize the total data center energy consumption by optimizing the increase in power consumption while placing a virtual machine to a physical machine. Our objective is to minimize the δP_j^i when placing a VM into a physical machine, satisfying service level agreements. Our aim is to minimize the fitness function (f) given as follows.

$$f = \sum_{j=1}^n \sum_{i=1}^m \delta P_j^i \quad (3.11)$$

3.2 System Architecture for Energy Efficient VM Placement and Selection (EE-VMPS)

The high level architecture of EE-VMPS in cloud data centers is shown in Figure 3.1. The problem involves placing the VMs to physical machines in energy efficient manner, and optimizing the current VM allocation. EE-VMPS integrates the following two modules :

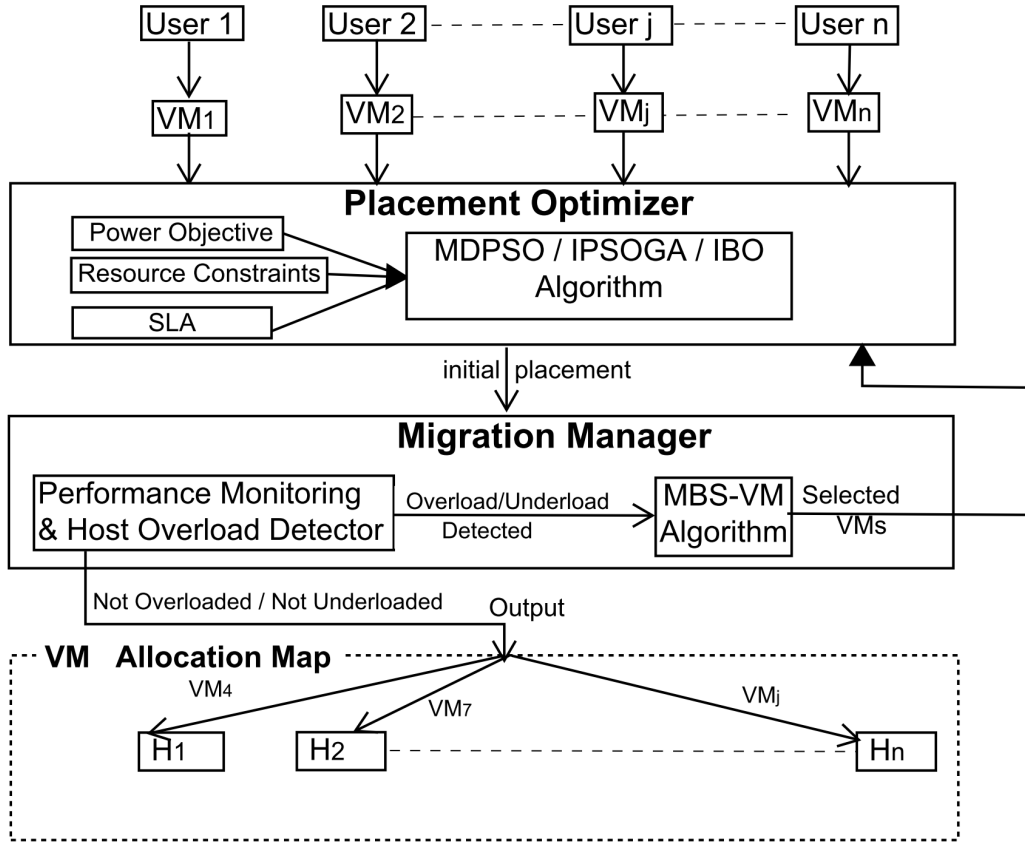


Figure 3.1: System Architecture of EE-VMPS

1. **Placement Optimizer** : This module performs the VM provisioning for new requests and places the VMs to physical machines using a Modified Discrete Particle Swarm Optimization (MDPSO) or an Interactive Particle Swarm Optimization and Genetic Algorithm (IPSOGA) or an Imitation

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

Based Optimization (IBO) (See Section 3.3) subject to the given resource constraints and SLA constraints.

2. **Migration Manager** : This module performs real time monitoring considering the utilization of VMs and physical machines and checks whether the physical machines exceed the given threshold or not. If any physical machine is over-utilized or under-utilized, then the migration manager selects some of the virtual machines using Memory, Bandwidth, and Size based Virtual Machine (MBS-VM) selection algorithm (See Section 3.4) and places them in other physical machines using the placement optimizer.

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

3.3.1 Approach 1 : VM Placement Using Modified Discrete Particle Swarm Optimization (MDPSO)

Discrete Particle Swarm Optimization (DPSO) is a version of classical Particle Swarm Optimization (PSO) that is modified to work with discrete valued search spaces consisting of binary/integer values [212]. DPSO shows discrete or subjective discrimination between variables [213]. DPSO is widely used to solve several combinatorial optimization problems when a good specialized algorithm is not available. In DPSO, a particle is composed of binary or integer variables and velocity is converted to chance of probability. This algorithm efficiently maintains the divergence of the swarm.

We propose a modified discrete particle swarm optimization approach (MDPSO) to find optimal energy aware virtual machine allocation. Particles are modeled to present a map of virtual machines to be placed in available physical machines in a data center. We define a novel particle encoding scheme in MDPSO with modified operators and velocity updation.

3.3.1.1 Particle Position and Velocity

Assume that ‘ m ’ workloads are allocated to any of the ‘ n ’ heterogeneous physical machines. We define a “position (X)” as a list of ‘ m ’ possible values and each value is chosen from the set P given in Equation 3.2. Every value in X represents a physical machine in a cloud data center chosen from the set P . Let $(v_1, v_2, v_3, \dots, v_m)$ is the list of VMs to be placed in the available physical machines in a data center. Then the particle position is a list of size M that gives the corresponding allocation of each VM in an optimization period (t), represented as follows.

$$X_i(t) = (p_1, p_5, p_1, \dots, p_j, \dots, p_n, \dots, p_j)$$

The velocity specifies the direction and tendency of a particle to move to a better position, shown in Equation 3.12. We model the velocity of a particle as a set $V_i(t)$, containing elements such as $(p_i \rightarrow p_j)$ where each $p_i, p_j \in P$ and $|V_i(t)|=M$. In MDPSO, the particles use the information of the local best and the global best to adjust their velocities.

$$V_i(t) = \{(p_i \rightarrow p_j) \mid p_i, p_j \in P \text{ and } |V_i(t)|=M\} \quad (3.12)$$

3.3.1.2 Transition (\odot) Operator

The transition operator is modeled to reflect the updation of a particle position ($X_i(t)$) based on its velocity. We use \odot to denote the transition operator. A particle $X_i(t)$ uses a new velocity $V_i(t)$ to adjust its current position and builds a new position $X_i(t+1)$. Let p_i be the value from $X_i(t)$ and $(p_j \rightarrow p_k)$ be the corresponding velocity from $V_i(t)$. Then $p_i \odot (p_j \rightarrow p_k)$ is defined as follows:

$$p_i \odot (p_j \rightarrow p_k) = \left\{ \begin{array}{ll} p_k & \text{if } p_i = p_j \\ p_i & \text{otherwise} \end{array} \right\} \quad (3.13)$$

If we apply $(p_i \rightarrow p_j)$ to evaluate a value different from p_i , then the position remains the same.

3.3.1.3 Subtraction (\ominus) Operator

The subtraction operator is defined to reflect the change between two particles, each representing a placement solution. For any two particles positions X_i and X_k , $X_i \ominus X_k$ depicts the velocity. In our modified approach, we define $X_i \ominus X_k$ for each physical machine $p_i, p_k \in P$ as follows:

$$p_i \ominus p_k = \{p_k \rightarrow p_i\} \quad (3.14)$$

where p_i is the value from X_i and p_k is the corresponding value from X_k .

3.3.1.4 Addition (\oplus) Operator

We define the addition operator (\oplus) to calculate the particle velocity updation, which depends on different factors including the local best of each particle, the global best of the swarm and the velocity inertia of a particle. Between two velocities $V_i(t)$ and $V_j(t)$, the addition defines the updated velocity. Let $(p_i \rightarrow p_x), (p_y \rightarrow p_j)$ be the velocity values from $V_i(t)$ and $V_j(t)$ respectively. $V_i(t) \oplus V_j(t)$ is defined as follows:

$$V_i(t) \oplus V_j(t) = \left\{ \begin{array}{ll} p_i \rightarrow p_j & \text{if } p_x = p_y \\ p_i \rightarrow p_x & \text{if } p_x \neq p_y \end{array} \right\} \quad (3.15)$$

3.3.1.5 Multiplication Operator (\otimes)

We modify the multiplication operator to update the particle velocity when it is multiplied by a real positive coefficient. The \otimes operator is stochastic, and defined only for a coefficient between 0 and 1. We can perform 'n' addition operations and one multiplication operation if the coefficient is greater than one. For each value $(p_i \rightarrow p_j)$ of $V_j(t)$, the multiplication operator is defined as follows:

$$c \otimes V_j(t) = \left\{ \begin{array}{ll} p_i \rightarrow p_i & \text{if } \hat{c} \leq c \\ p_i \rightarrow p_j & \text{if } \hat{c} > c \end{array} \right\} \quad (3.16)$$

where $c \in [0, 1]$ and $\hat{c} = \text{random}(0, 1)$.

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

Algorithm 3.1: VM Placement Using MDPSO

Input: List of virtual machines and physical machines

- 1 Initialize the parameters k , r , s , *population size* and learning factors (k_1, k_2) .
- 2 Initialize the population using random sequences.
- 3 **foreach** $particle \in Swarm$ **do**
 - 4 $fitness[i] = \sum_{j=1}^m E(p_j(u_1)) - E(p_j(u))$.
 - 5 $L_i = X_i(t)$. // Local Best
- 6 $\hat{L}_j(t) =$ particle with minimum $\sum_{j=1}^m E(p_j(u_1)) - E(p_j(u))$ in the swarm. // Global Best
- 7 $prevFitness = fitness$.
- 8 **repeat**
 - 9 **foreach** $particle \in Swarm$ **do**
 - 10 // Update Velocity and Position of each particle
 - $V_i(t+1) = w \otimes V_i(t) \oplus k_1 r_1(t) \otimes (L_i(t) \ominus X_i(t)) \oplus k_2 r_2(t) \otimes (\hat{L}_j(t) \ominus X_i(t))$
 - $X_i(t+1) = X_i(t) \odot V_i(t+1)$;
 - 11 $fitness[i] = \sum_{j=1}^m E(p_j(u_1)) - E(p_j(u))$.
 - 12 **if** $fitness[i] < prevFitness[i]$ **then**
 - 13 | $L_i = X_i(t+1)$.
 - 14
 - 15 $\hat{L}_j(t+1) =$ particle with minimum $\sum_{j=1}^m E(p_j(u_1)) - E(p_j(u))$ in the swarm.
 - 16 **if** $fitness(\hat{L}_j(t+1)) < fitness(\hat{L}_j(t))$ **then**
 - 17 | $\hat{L}_j(t) = \hat{L}_j(t+1)$.
 - 18 | $Flag = True$;
 - 19 | $count ++$;
 - 20 | **if** $Flag == True$ and $count > r$ **then**
 - 21 | | Return $\hat{L}_j(t)$.
 - 22
 - 23 **else**
 - 24 | **if** $fitness(\hat{L}_j(t+1)) == fitness(\hat{L}_j(t))$ **then**
 - 25 | | $Flag = false$;
 - 26 | | $count1 ++$;
 - 27 | | **if** $count1 > s$ **then**
 - 28 | | | Return $\hat{L}_j(t)$.
 - 29
 - 30 **else**
 - 31 | $Flag = false$.
 - 32 | $count = count1 = 0$.
- 33 **until** k -times ;

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

Using the said operations, we improve the DPSO by reformulating the equations of velocity and position in PSO (see Appendix A.1) as follows:

$$V_i(t+1) = w \otimes V_i(t) \oplus k_1 r_1(t) \otimes (L_i(t) \ominus X_i(t)) \oplus k_2 r_2(t) \otimes (\hat{L}_j(t) \ominus X_i(t)); \quad (3.17)$$

$$X_i(t+1) = X_i(t) \odot V_i(t+1); \quad (3.18)$$

The pseudocode for the allocation of virtual machines using our MDPSO is illustrated in Algorithm 3.1. The particle data in this case is the VM-Physical Machine map that shows the allocation of a VM to a physical machine. Initially all the VMs are placed randomly in different physical machines. We calculate the fitness of each particle in the swarm using the Equation 3.11. Then we choose the particle having a minimum fitness as the global best. After each iteration, all particles try to move closer to the coordinates of the optimal solution in the population which provides less energy consumption. This happens during the updation of the velocities of each particle using their local best and global best as shown in Equation 3.17. These velocities are applied to the respective particles to update their position as shown in Equation 3.18. Individual particles of the swarm are manipulated until the particle matches the target by updating the velocity and position. Algorithm 3.1 terminates either after k -iterations or achieves continuous improvement in the global best r -times.

3.3.2 Approach 2 : VM Allocation Using Interactive PSO-GA (IPSOGA)

Interactive PSO-GA combines the searching abilities of both genetic algorithm (GA) and particle swarm optimization (PSO). To make particles more suitable to the environment before producing offsprings, we incorporate the social interaction between the algorithms to enhance the population on each generation. The combination of GA and PSO algorithms has been investigated in many studies [214, 215, 216, 217]. But, to the best of our knowledge, parallelization of these algorithms is not considered in the literature instead authors worked on serialization of these algorithms in one or other ways [218, 219]. Parallelization helps to enhance the computational throughput and global search capability. If the algorithm is properly designed to take advantage of parallelism, it can execute faster than their sequential counterparts. So, by taking the recent advances in multi core

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

architectures, we propose to use multi-threading and shared memory to develop a parallelized optimization algorithm called “Interactive PSO-GA” (IPSOGA).

This algorithm performs parallel processing of particle swarm optimization (PSO) and genetic algorithm (GA) using multi-threading. Firstly, basic PSO and GA algorithms will run in parallel in separate threads. After each iteration, these threads share the top ranked particles information with each other to form the population for the next iteration and override the poor-performing particles. This technique helps balancing between improving convergence time and accuracy.

Our proposed IPSOGA is illustrated in Figure 3.2. Both the PSO and GA algorithms are population based and each particle data represents a candidate solution to the optimization problem. Particles are modeled to present a map of virtual machines to be placed in available physical machines in a data center. Our algorithm starts with randomly initializing the population and other parameters like inertia coefficients, learning factors, crossover rate, and mutation rate. We calculate the fitness of each particle in the swarm and then we choose the particle having a maximum fitness as the global best. With the said information, the

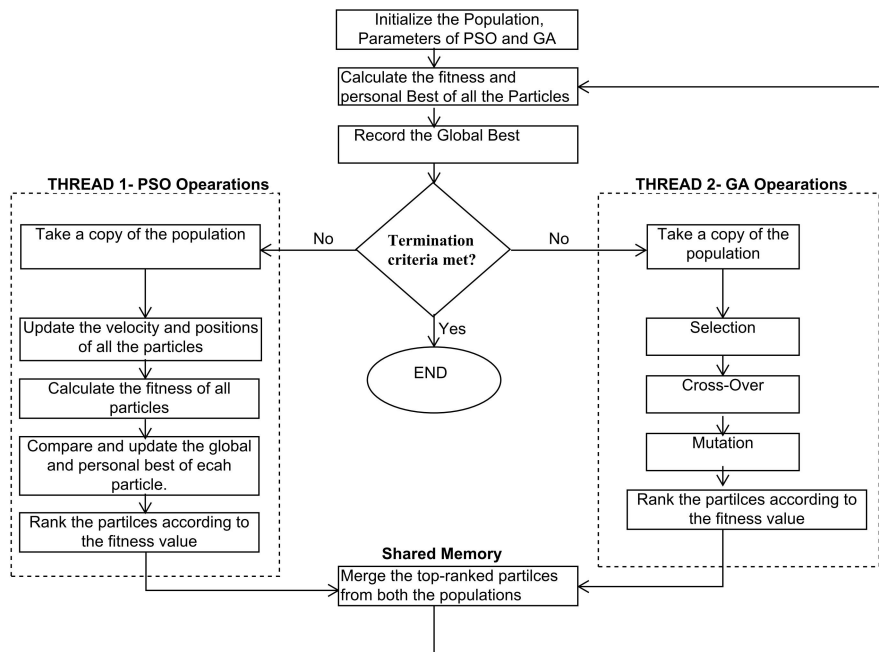


Figure 3.2: Flow diagram of IPSOGA

program is broken into two discrete parts that can run concurrently in two different

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

threads. The first thread runs the PSO algorithm and the second thread runs GA algorithm. The first thread follows the flow explained in Appendix A.1, where all the particles try to move closer to the coordinates of the optimal solution in the population. The second thread follows the flow given in Appendix A.3.

After each iteration, both the threads rank the particles of the respective population according to the particles fitness value. Finally, in step 10, we create a new population for the next iteration by merging the top particles (according to fitness) in shared memory context as shown in Algorithm 3.2, where the process of merging is done asynchronously. Our algorithm terminates either after M-iterations or achieves continuous improvement in the global best r-times as shown in Algorithm 3.1.

On a multi-core machine, we can improve the actual performance of compute bound operations, by executing multiple threads in parallel, with every core handling a separate thread simultaneously. To implement parallel version of an algorithm, the code should be easily vectorized or partitioned to run in parallel on a multi-core system. Furthermore, it should be scalable to permit efficient utilization of computational resources. Evolutionary algorithms such as PSO and GA are better suited for parallel implementation. PSO and GA are population based algorithms and the number of individuals in each population can be adjusted according to the availability and capacity of processors. However, increasing the population size allows for a higher global exploration in the search space and avoids the chance of trapping into local minima.

The basic PSO algorithm was designed in synchronous fashion for updating particle's best (pBest) and swarm's best (gBest) fitness values along with their positions. In this approach, we need to evaluate the fitness of all the particles before updating the pBest and gBest values. However, updating the pBest and gBest after each individual iteration (i.e, in an asynchronous fashion) improves convergence rates [220, 221]. Because of asynchronous update, there is a chance that other particles quickly react to the change in the global best value. In our approach, we evaluate the fitness values concurrently, which indicates the synchronous design. The parallel implementation requires updating gBest of the population after PSO and GA complete their updates. Consequently, this will reduce the convergence rate compared to the asynchronous implementation. To improve the convergence rate, we enhance the population using the ranked particles from both PSO and GA

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

Algorithm 3.1: VM Placement Using IPSOGA

```

1 Initialize the parameters population size, mutation-rate, crossover rate,
  learning factors ( $k_1, k_2$ ), maximum iterations(M), and r.
2 Initialize the population using sub-random sequences.
3 foreach particle  $\in$  Swarm do
4   | Calculate the fitness of the particle.
5   | Store the local best of the particle.
6  $\hat{L}_j(t)$  = particle with maximum fitness value in the swarm. // Global Best
7 repeat
8   | Perform the PSO operation and GA operation in different threads.
9   | Rank the individual populations of PSO and GA according to the fitness
  values.
10  newPopulation = mergePopulation(psoPopulation, gaPopulation);
11   $\hat{L}_j(t + 1)$  = particle with maximum fitness in the population.
12  if  $fitness(\hat{L}_j(t + 1)) > fitness(\hat{L}_j(t))$  then
13    |  $\hat{L}_j(t) = \hat{L}_j(t + 1)$ .
14    | Flag = True;
15    | count ++;
16    | Update the populations of the PSO and GA with the newPopulation.
17    | if  $Flag == True$  and  $count > r$  then
18      | Return  $\hat{L}_j(t)$ .
19    |
20  |
21  else
22    | if  $fitness(\hat{L}_j(t + 1)) == fitness(\hat{L}_j(t))$  then
23      | Flag = false;
24      | count1 ++;
25      | if  $count1 > s$  then
26        | Return  $\hat{L}_j(t)$ .
27      |
28    | else
29      | Flag = false.
30    | count = count1 = 0.
31 until M-times ;

```

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

Algorithm 3.2: mergePopulation()

```
input : population1, population2
output: newPopulation
1 temp1= first particle of the population1;
2 temp2= first particle of the population2;
3 p=g=0;
4 for  $i=0; i < \text{population1.length}(); i++$  do
5   if  $\text{fitness}(\text{temp1}) \geq \text{fitness}(\text{temp2})$  then
6     newPopulation[i]=population1[p];
7     temp1= population1[+p];
8   else
9     newPopulation[i]=population2[g];
10    temp2= population2[+g];
11 return newPopulation
```

algorithms. After each iteration, we rank the individuals of the populations concurrently. We then update the initial population with the top-ranked individuals chosen from both the populations, as shown in Algorithm 3.2.

3.3.3 Approach 3 : VM Allocation Using Imitation Based Optimization (IBO)

In an attempt to address few challenges like decreasing the convergence time and increasing the consistency in providing optimal solutions, we propose a new optimization algorithm named imitation based optimization (IBO) based on the human behavior of learning and imitating. Having a good reason for copying the behaviour of the leader, other particles of the swarm starts paying attention to the leader. They remember all / some behaviour of the leader and tries to imitate the behaviour with some confidence. Imitation is an important aspect of social learning. If a particle is capable of carrying out the behaviour which they have observed to achieve the desired result, then they start imitating but may or may not achieve most. This explains the differences in the states of the particles. The individual being observed (Leader) is an important factor in this approach. In

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

general, people are more likely to be influenced by the person with status and power. In our case the chance of imitating a particle increases, if it has a good fitness value than the other particles in the swarm. The leader may have a powerful influence if they have gained good position through their movement.

Algorithm 3.3: Calculating Mahalanobis Distance - MD (A, B)

Input: Two matrices A and B with same number of columns

- 1 Center the data of each matrix on arithmetic mean of each column to get \hat{A} and \hat{B} .
- 2 Calculate the covariance matrices for \hat{A} and \hat{B} as:

$$C_1 = \frac{1}{n1} * \hat{A}^T * \hat{A}$$

$$C_2 = \frac{1}{n2} * \hat{B}^T * \hat{B}$$

where $n1$ and $n2$ are the number of rows in \hat{A} and \hat{B} respectively.

- 3 Calculate the pooled covariance matrix of C_1 and C_2 as:

$$PC = \frac{1}{n1 + n2} [n1 * C_1 + n2 * C_2]$$

- 4 Calculate the mean difference matrix Diff as:

$$Diff = mean(A) - mean(B)$$

- 5 Calculate the Mahalanobis distance as:

$$MD = \left\{ \begin{array}{ll} \sqrt{Diff * PC^{-1} * Diff^T} & \text{if inverse of PC is possible} \\ 1 & \text{otherwise} \end{array} \right\}$$

The proposed algorithm is a population-based algorithm where the population consists of a group of individuals (candidate solutions). Each individual starts with a random position and has some level of capacity (fitness) for solving a problem. The individual with high fitness is announced as the best participant (Leader). This information is propagated to all the individuals in the population. The procedure of IBO is divided into two phases. The first phase is “grouping phase”

3.3 Energy Efficient VM Placement Approaches in Cloud Data Centers

where the individuals who are similar to the best participant will be grouped together. This similarity is designed using the Mahalanobis Distance as given in Algorithm 3.3, and this group is named as “pioneers”. The individuals having less Mahalanobis Distance are treated as more similar and vice versa.

We have chosen the Mahalanobis Distance (MD) for deriving a similar individual because it distinguishes the pattern of a certain group from other groups in a better way [222]. It is a distance measure based on correlations between variables by which different patterns could be identified and analyzed with respect to a reference baseline [223]. The MD measure can be used to determine the similarity of a set of values from an unknown sample to a set of values measured from a collection of known samples. One of the main applications of MD is to introduce a scale based on all characteristics to measure the degree of abnormality. It measures the distances in multi-dimensional spaces taking into account the correlations between variables or characteristics that may exist. When compared to other multivariate measurement techniques such as the Euclidean distance, MD is superior [224]. The Euclidean distance, for instance, does not indicate how well the unknown data matches with the reference data and does not take into account the distribution of the points (correlation). Calculating MD between the two candidate solutions is given in Algorithm 3.3.

The second phase is the “Imitation phase”. In this phase each person in the population except those in pioneers group enhances their knowledge by imitating the best person in the group. This imitation is done by updating the solution using Algorithm 3.4.

The flow of algorithm is explained in Algorithm 3.5. Our algorithm starts with a random population. Each individual in the population represents VM-Host mapping as discussed in Section 3.3.1.1. Then we calculate the fitness of each individual in the population and decide the best participant of the population. The best participant is the one who has the lowest fitness value (energy consumption) and this information is propagated to all other individuals. Furthermore, our approach categorizes the generated population of particles into similar (pioneers) and non-similar particle groups based on similarity with the best particle which is calculated using the Mahalanobis distance as described in Algorithm 3.3. Then each individual in the non-similar particle group tries to imitate the best participant by doing the operations as presented in Algorithm 3.4. Once all the individuals are

3.4 VM Selection and Migration using MBS-VM

Algorithm 3.4: Update Particle

Input: particle A, Best particle B

- 1 Convert the given particles into matrix form according to the problem.
 - 2 Calculate the mean of matrix A and B. //
*mean returns a row vector containing
the mean of each column.*
 - 3 Calculate the mean difference D:
 - 4 $D = \text{mean}(A) - \text{mean}(B)$
 - 5 **foreach** row in A **do**
 - 6 **foreach** column $\in A$ [r] **do**
 - 7 $A[\text{row}][\text{column}] = A[\text{row}][\text{column}] + D[\text{column}]$
 - 8 dist-A = Calculate the Mahalanobis distance between A, B.
 - 9 create a new particle (C) using sub-random sequences.
 - 10 dist-C = Calculate the Mahalanobis distance between C, B.
 - 11 **if** $\text{dist-C} < \text{dist-A}$ **then**
 - 12 Return C.
 - 13 **else**
 - 14 Return A.
-

updated, we decide the best participant again and compare it with the previous one. This process continues until we get continuous improvements in the updated particle for ‘r-times’ or will not change for ‘s-times’ or up to ‘max-iterations (k)’.

3.4 VM Selection and Migration using MBS-VM

We perform migration of the virtual machines to increase the efficiency and to ensure the compliance of SLA. We choose two extreme states of the servers for VM selection: over-utilization and under-utilization. If a server is under-utilized, then we select all the virtual machines in that physical machine. If a server is over-utilized, then we select the virtual machines for migration using the proposed selection algorithm until the physical machine utilization drops below the upper threshold.

3.4 VM Selection and Migration using MBS-VM

Algorithm 3.5: VM Placement Using IBO

Input: List of decision variables

- 1 Initialization:
- 2 grouping factor(γ) = 0.5,
- 3 population size =40, $r=2$, and $s=10$,
- 4 Initialize all the particles of swarm using sub-random sequences.
- 5 **foreach** $particle \in Swarm$ **do**
- 6 | Calculate the fitness of the particle.
- 7 | $L_{particle} = X_i(t)$. // Local Best
- 8 $\hat{L}(t)$ = particle with minimum fitness value in the swarm. // Global Best
- 9 $prevFitness = fitness$. //save a copy of the array
- 10 **repeat**
- 11 | **foreach** $particle \in swarm$ **do**
- 12 | Calculate the Mahalanobis Distance (MD) between the $\hat{L}(t)$ and particle.
- 13 | **if** $MD < \gamma$ **then**
- 14 | add the particle to poineers list.
- 15 |
- 16 | **foreach** $particle \notin poineers$ **do**
- 17 | $X_i(t+1)$ = Update the particle according to Algorithm 2
- 18 | Calculate the fitness of the updated particles
- 19 | **if** $fitness[particle] < prevFitness[particle]$ **then**
- 20 | $L_i = X_i(t+1)$.
- 21 |
- 22 | $\hat{L}(t+1)$ = particle with minimum fitness value in the swarm
- 23 | **if** $fitness(\hat{L}_j(t+1)) < fitness(\hat{L}_j(t))$ **then**
- 24 | $\hat{L}_j(t) = \hat{L}_j(t+1)$.
- 25 | Flag = True;
- 26 | count ++;
- 27 | **if** $Flag == True$ and $count > r$ **then**
- 28 | Return $\hat{L}_j(t)$.
- 29 |
- 30 | **else**
- 31 | **if** $fitness(\hat{L}_j(t+1)) == fitness(\hat{L}_j(t))$ **then**
- 32 | Flag = false;
- 33 | count1 ++;
- 34 | **if** $count1 > s$ **then**
- 35 | Return $\hat{L}_j(t)$.
- 36 |
- 37 | **else**
- 38 | Flag = false.
- 39 | count = count1= 0.
- 40 **until** k -times ;

We propose a Memory, Bandwidth, and Size based Virtual Machine (MBS-VM) selection algorithm considering the memory utilization, bandwidth utilization, and size of the virtual machine. The migration task adds load to both the target and

3.4 VM Selection and Migration using MBS-VM

host systems, which effects the performance of both the hosts as well as the data center. Hence, we carefully select the migratable virtual machines in the over-utilized hosts. In order to achieve this, we designed a cost function that effectively balances memory, bandwidth, and size factors. This function is further used to select a VM among the migratable VMs. The definitions for the utilization of various elements that we considered are discussed as follows:

Memory Utilization: Memory utilization of a VM reflects the usage of the physical machine by a VM. Generally, memory usage changes slowly when the host utilization is below the lower-threshold, and increases rapidly when approaching the upper-threshold. Considering these facts, we define the memory utilization of a VM (v_i) in a host (p_j) as shown in Equation 3.19.

$$VM_i(mm) = \frac{\text{Current allocated RAM for } v_i}{\text{Total RAM requested by } v_i}. \quad (3.19)$$

Bandwidth Utilization: Virtual machine migration involves transferring large amounts of data between hosts. VMs comprising a multi-tier application make complex load interactions among the underlying physical servers. For such applications, we need to consider bandwidth allocation and usage for a VM as shown in Equation 3.20

$$VM_i(bw) = \frac{\text{Current allocated bandwidth for } v_i}{\text{Total bandwidth requested by } v_i}. \quad (3.20)$$

VM Size: The amount of data transfer is directly related to the migration cost. The amount of data transfer is the only factor that has a linear relationship with the power consumption and also with migration time when physical machines have stable loads. In fact, the amount of data is controlled by the memory size of migratable VMs [225]. VM Size is defined as follows:

$$VM_i(size) = \text{Size of the virtual machine } (v_i). \quad (3.21)$$

Based on the concept of memory utilization, bandwidth utilization and VM size, the cost function for a VM is defined as the weighted sum, as follows:

$$\text{Cost}(VM_i) = a.VM_i(mm) + b.VM_i(bw) + c.VM_i(size). \quad (3.22)$$

where $a, b,$ and c are the random weights such that $a + b + c = 1$.

3.4 VM Selection and Migration using MBS-VM

Algorithm 3.6: VM Selection Using MBS-VM

Input: List of physical machines: *hostList*

```
1 Initialize lowerThreshold, upperThreshold.
2 foreach host  $\in$  hostList do
3   host-utilization = getHostUtilization (host);
4   VMlist = host.getVmlist ().
5   if host-utilization  $\geq$  upperThreshold then
6     repeat
7       foreach VM  $\in$  VMlist do
8         Initialize a, b, and c with random value between 0.0 and 1.0.
9         total = a + b + c.
10        a = a/total.
11        b = b/total.
12        c = c/total.
13        Cost(Vm) = a.VMi (mm) + b.VMi (bw) + c.VMi (size).
14        Sort the VMlist in decreasing order of their Cost.
15        probedVM = VM with highest cost.
16        if probedVM is not Aggressive and probedVM is not in migration
17          then
18            update host-utilization without probedVM.
19            migrationList.add (probedVM)
20            VMlist = VMlist - probedVM .
21          else
22            continue;
23        until host-utilization < upperThreshold ;
24    if host-utilization  $\leq$  lowerThreshold then
25      migrationList.addAll (VMlist).
26      Remove all the VMs form VMlist of the host.
27 return migrationList.
```

Algorithm 3.6 describes the process of virtual machine selection and migration using memory utilization, bandwidth utilization, and VM size. The MBS-VM algorithm calculates the cost of each migratable VM in a physical machine. It selects a virtual machine with the highest cost among the migratable VMs and then checks whether the selection of this VM leads to an aggressive consolidation which increases the possibility of over-utilization. If not, the algorithm adds this VM to the migration list. This process continues until the CPU utilization drops below the upper utilization threshold when the upper threshold is violated. If the physical machine is underutilized i.e utilization falls below the lower threshold, it is better to put the physical machine to sleep mode. To achieve this, we migrate all VMs to another physical machine to eliminate the idle power consumption. If a VM is running in an over utilized physical machine, the proposed algorithm automatically migrates the VM to a target host which gives the better performance. Once the process of VM selection using the MBS-VM algorithm is completed, the selected virtual machines are admitted for VM provisioning.

3.5 Performance Evaluation

We compare our proposed approaches with few other approaches concerning energy consumption, SLA violations, number of migrations and active hosts. We also conduct the convergence analysis of the proposed approaches and present the parallel efficiency for the proposed IPSOGA approach. We performed our experiments simulating a cloud data center using CloudSim [226]. CloudSim follows a layered implementation structure and provides support for simulating cloud data center environments.

3.5.1 Experimental Environment

In our experiments, we simulated a cloud data center comprising 100 heterogeneous physical machines. These physical machines were characterized according to their configurations as shown in Table 3.1. Virtual Machine Provisioner is responsible for allocating VMs to physical nodes. Resource requests are of heterogeneous nature and VM parameters include RAM, MIPS, and bandwidth. The characteristics for each type of VM are presented in Table 3.2. The user submits

3.5 Performance Evaluation

the requests for provisioning of VMs. Each virtual machine processes time varying workloads. This workload is designed in such a way that CPU is loaded according to a uniformly distributed random variable. This load runs for 4,32,000 million instructions (MI) which is equal to running a machine with 300 MIPS processing power for 24 minutes with 100% utilization. We assume the idle conditions of the

Host Type	MIPS	RAM	Bandwidth	Storage
G4 Xeon 3040	3000	8.5 GB	30 GBPS	125 GB
G4 Xeon 3075	3000	8.5 GB	30 GBPS	125 GB
G3 Pentium D 930	3000	8.5 GB	30 GBPS	125 GB

Table 3.1: Physical Machines Configurations

host upon starting. All the programs are written in Java 8 and simulations have been executed in a 64-bit Ubuntu Operating System running on an Intel Core i7-5500U, 2.40 GHz Dell Latitude with four cores, 540GB hard disk and 8 GB of RAM.

VM Type	MIPS	RAM	Bandwidth	Storage
Type 1	100	613 MB	100 MBPS	2.5 GB
Type 2	200	870 MB	100 MBPS	2.5 GB
Type 3	300	1740 MB	100 MBPS	2.5 GB

Table 3.2: Virtual Machines Requirements

Tests were performed to set the initial parameters for all the approaches. Initial populations size of 10, 20 and 30 are falling to a local optimum regardless of the swarm behavior where as the population size of 50 needs more iterations. Thus, the population size of swarm is set to 40 which has a linear convergence rate with increasing number of VMs and provides the optimal solution. From our experiments, we observed that the maximum number of iterations above 150 ensures a global optimum solution without falling to local optimum. So, the maximum number of iterations is set to 150. Learning parameters (k_1, k_2) are set to 3, 2 respectively, to prevent divergence of the particles. These values are chosen after several experiments and the way these weights coordinate during the searching process after each iteration is illustrated in Appendix A.4. For IPSOGA approach,

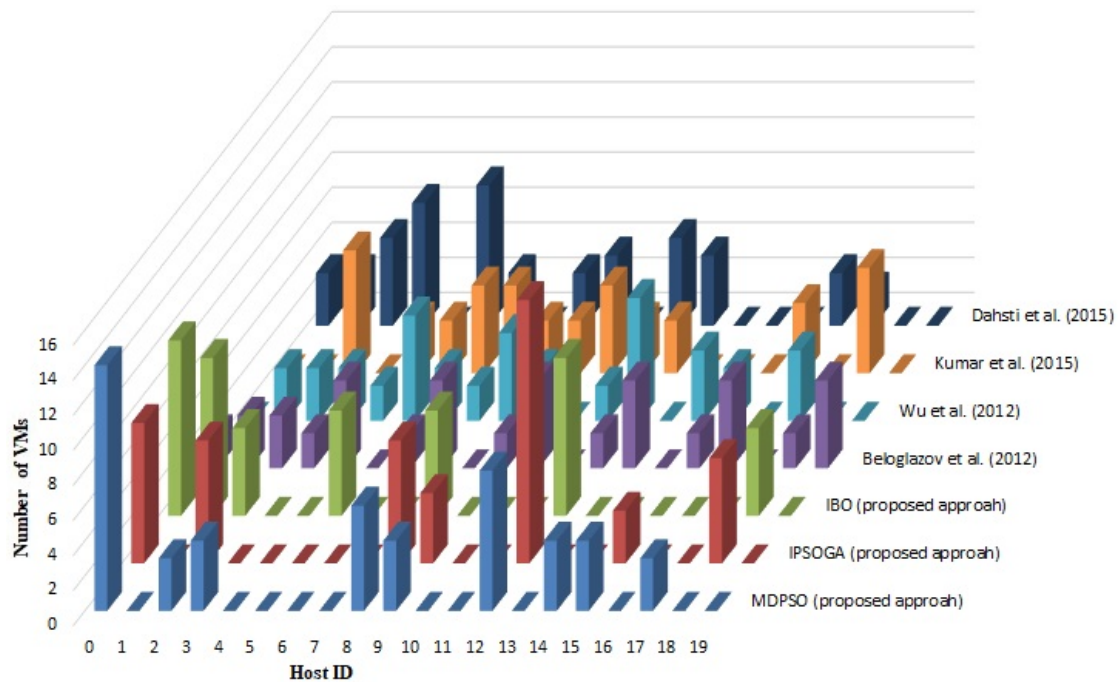


Figure 3.3: Comparison of active hosts in different algorithms

the number of threads used for each evaluation is set to 2 and all the approaches are allowed to run until there are ‘3’ continuous improvements in the optimal solution or there is no improvement in the optimal solution for 10 iterations or 150 function evaluations are completed.

3.5.2 Analysis of VM Allocation Approaches

We evaluated the performance of our proposed approaches by comparing the solutions of Beloglazov et al. [77], Wu et al. [122], Kumar et al. [110], and Dashti et al. [109]. The ultimate goal of any VM placement algorithm is to reduce the number of active hosts thereby reducing the energy consumption of the data center. We evaluated all the approaches with 20 physical machines and 50 virtual machines to find the number of active hosts in the data center. Figure 3.3 illustrates the number of virtual machines placed in each of the physical machines by the said algorithms. The proposed approaches reduce the number of active machines by placing more workloads in an efficient physical machine. From Figure 3.3, we see that there are 9, 7, and 7 active hosts with the proposed MDPSO, IPSOGA, and IBO approaches

3.5 Performance Evaluation

where as the number of active physical machines by [77], [122], [110], and [109] are 15, 14, 12, and 12 respectively. Thus, the proposed approaches use the minimum number of physical machines while satisfying the resource requirements of VMs.

VMs	Energy Consumption (kWh)						
	MDPSO (proposed)	IPSOGA (proposed)	IBO (proposed)	Beloglazov et al. (2012)	Wu et al. (2012)	Kumar et al. (2015)	Dashti et al. (2015)
50	0.91	0.91	0.91	1.10	1.09	0.92	0.92
100	1.01	1.01	0.95	1.52	1.31	1.03	1.20
150	1.30	1.10	1.3	2.02	1.82	1.40	1.60
200	1.45	1.50	1.35	2.35	1.89	1.65	1.75
250	1.70	1.65	1.52	2.64	2.50	1.92	1.97
300	2.25	2.12	2.08	3.15	2.93	2.47	2.50

Table 3.3: Performance comparison in terms of energy consumption

In an environment of 300 virtual machines, IBO, IPSOGA, and MDPSO lead to the minimum energy consumption of 2.08 kWh, 2.12 kWh and 2.25 kWh respectively where as Beloglazov et al. [77], Wu et al. [122], Kumar et al. [110], and Dashti et al. [109] consume 3.15 kWh, 2.93 kWh, 2.47 kWh, and 2.50 kWh respectively. Thus, with the proposed IBO, IPSOGA, and MDPSO approaches, the energy consumption is decreased upto 33%, 31% and 28% respectively compared to Beloglazov et al. [77]. We tested with the same setup varying the number of VMs, starting form 50 to 300. The results are presented in Table 3.3.

3.5.3 Performance Analysis in terms of Migrations

Excessive virtual machine migrations in a data center could affect the desired quality of service (QoS) level under peak loads and could increase the energy consumption of a data center. In Figure 3.4, we present a comparison of our proposed approaches with different VM allocation policies in a data center with 100 physical machines. The comparison is in terms of the number of migrations with increasing virtual machines. From the experiments we observed that the number of migrations for all the algorithms are almost equal in case of 50 virtual machines. But with the increase in number of virtual machines, the proposed approaches outperformed the other approaches. Figure 3.4 shows that there is no

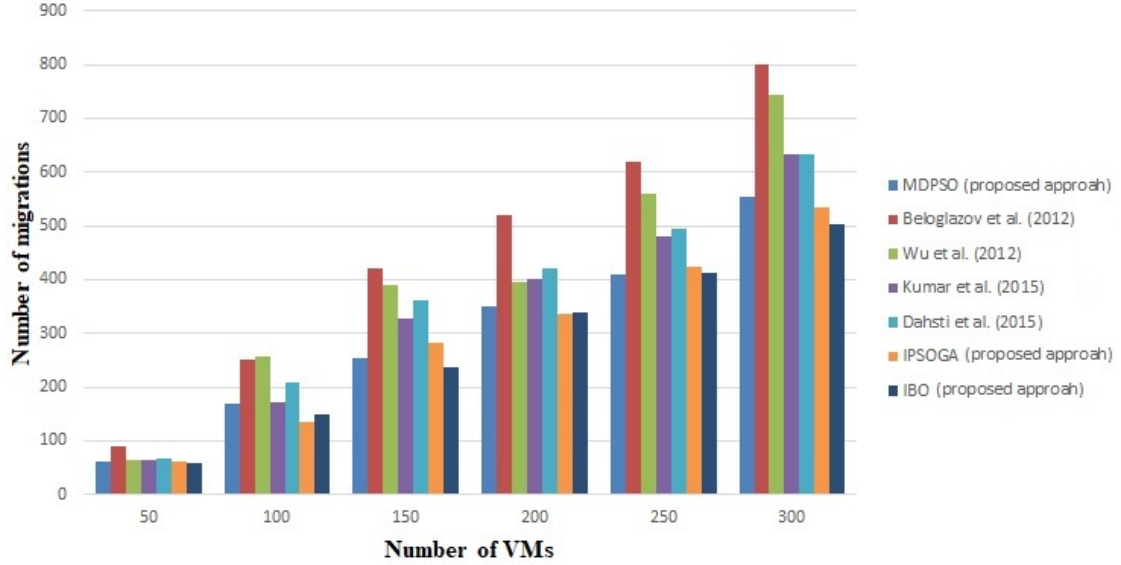


Figure 3.4: Number of Migrations vs. Number of Virtual Machines

significant difference in the number of migrations for IPSOGA and IBO. Further, IBO entails the minimum number of migrations in all cases as compared to MDPSO and other approaches.

3.5.4 Performance Analysis in terms of SLA Violations

Another important aspect to ensure the smooth operations of a data center environment is to guarantee the service level agreements (SLA) in terms of satisfying QoS requirements such as availability, reliability, and throughput. SLA violations are determined as follows :

$$SLA\ Violation = \frac{(Requested\ resource - Allocated\ resource)}{Requested\ resource} \quad (3.23)$$

The SLA is violated, if a physical machine fails to allocate the requested processor share to a VM. In our model, we calculate the SLA violation for each VM request using Equation 3.23. Overall SLA violations as well as resource requests that were not allocated are determined by CloudSim and the results are presented in Table 3.4. The overall SLA violations for 50 and 100 virtual machines are mostly same for all the approaches. The results in Table 3.4 shows that all our proposed

3.5 Performance Evaluation

	Overall SLA Violations (%)						
VMs	MDPSO (proposed)	IPSOGA (proposed)	IBO (proposed)	Beloglazov et al. (2012)	Wu et al. (2012)	Kumar et al. (2015)	Dashti et al. (2015)
50	0.72	0.72	0.41	0.75	0.75	0.94	1.02
100	1.25	1.21	1.14	1.27	1.34	1.38	1.75
150	1.32	1.26	1.12	1.4	1.54	1.52	1.69
200	1.25	1.19	1.2	1.5	1.66	1.7	1.84
250	2.03	1.72	1.67	1.95	2.21	2.5	2.4
300	2.26	2.17	2.09	2.2	2.41	2.78	3.13

Table 3.4: Overall SLA Violations

approaches have less SLA violations in comparison to other approaches. Further, IBO and IPSOGA demonstrate lesser overall SLA violations than MDPSO for the increase in the number of virtual machines.

3.5.5 Performance Analysis in terms of Convergence

The scalability of our approaches is specified in terms of the time to converge and the space complexity of all the proposed algorithms. These algorithms are influenced by three factors: number of VMs (M), number of physical machines (N) and number of migratable VMs. Starting with a swarm of size S and repeating the procedure for K times, to reach the global optimum solution, the computational complexity of the proposed algorithms is in the order of $O(SK)$ [227]. The time complexity of the MBS-VM algorithm is proportional to the product of the number of over-utilized physical machines and the number of migratable VMs in each of these physical machines. In our proposed approaches, we modeled each particle as a M -dimensional vector, so that the search space is limited to I^M where I is the population size.

For meta-heuristic approaches, the complexity is specified in terms of the hardness of the problem and the size of input than the algorithm because the complexity of the algorithm is affected by the representation, initialization of parameters, termination condition and the diversity of the initial solutions, etc. For each algorithm, the convergence rate is calculated as the average number of iterations it takes to reach the global optimum solution. We plot the convergence of IBO,

IPSOGA and MDPSO for varying design variables, ranging from 50 to 300. From

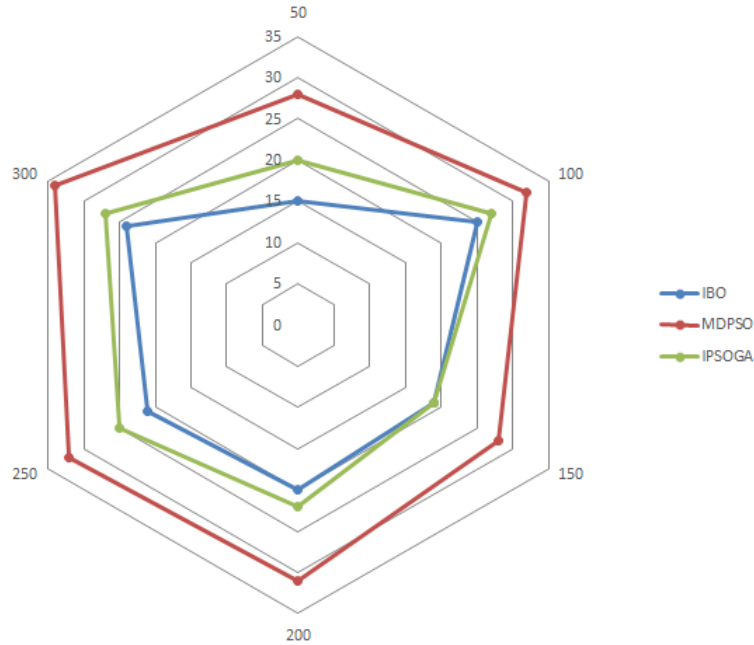


Figure 3.5: Convergence Analysis

the experiments, we observe that IBO approach performs better for higher number of virtual machines. IBO decreases the convergence time up to 35% and 13% compared to MDPSO and IPSOGA respectively. Further, it shows higher consistency and lower convergence rate with the increasing of design variables. The convergence time of IPSOGA, MDPSO and IBO for varying virtual machines are given in Figure 3.5.

3.5.6 Speedup and Parallel Efficiency of IPSOGA

The fundamental reason for parallelizing PSO and GA programs is to make them intractable and run quicker. To evaluate the advantage of parallel program, we generally use the speedup parameter. The speedup of a parallel code is how much faster it runs in parallel. Speedup expresses the relative speed of the parallel program compared to the sequential one [228]. The speedup of a parallel program is defined as the ratio of the rate at which work is done when a job is run on N processors to the rate at which it is done by just one. The speedup is calculated as

follows: where T_s is used to express the execution time for the given optimization problem when PSO and GA run in sequence. We took the fastest execution time among 20 runs. T_p is the execution time of the parallelized version of the program (IPSOGA) for solving the same problem.

$$S = \frac{\text{Time for sequential run } (T_s)}{\text{Time for parallel run } (T_p)}$$

If multiple threads are executed individually on a computer with p processors, the speedup is equal to 'p'. The maximum speedup is achievable if the load is balanced among the available processors and there is no communication cost. In such an environment, each processor requires T_s/p time units to complete the thread execution and the speedup will be:

$$\text{Speedup} = \frac{T_s}{T_s/p} = p$$

But, according to the Amdahl law [229], it is difficult to obtain a speedup value 'p' because every program has a fraction (δ) of code that cannot be parallelized. In our case the code for merging the two populations after each iteration cannot be parallelized. The remaining code can be executed in parallel. The execution time and the speedup are given as follows:

$$T_p = T_s * \delta + T_s * (1 - \delta)/p$$

$$\text{Speedup} = \frac{T_s}{T_p} = \frac{p}{\delta * (p - 1) + 1}$$

Parallel performance is calculated using the speedup as given in the said equations. We use the sequential execution time of PSO, GA and the execution time of the IPSOGA for calculating speedup. In general, parallel efficiency is the ratio of speedup to the number of threads used. We reported the speedup along with parallel efficiency for 50, 100, 150, 200, 250, and 300 design variables in Table 3.5.

3.5.7 Analysis of VM Selection Algorithms

We analyzed our proposed MBS-VM selection algorithm in terms of energy consumption for the given set up. We compared our proposed approach with the Minimization of Migrations (MM) policy [77] and the Minimum Memory Size

	IPSOGA	
Number of VMs	Parallel Efficiency (%)	Speedup
50	97.5	1.95
100	100	2
150	87.5	1.75
200	82.5	1.65
250	85	1.70
300	80	1.60

Table 3.5: Speedup and Parallel Efficiency of IPSOGA

(RAM) policy [43]. These selection approaches are evaluated in combination with the proposed MDPSO placement method.

The MM policy selects the minimum number of virtual machines to avoid a host being over-utilized using thresholds. The RAM policy selects a VM with the minimum memory size in a over utilized host which in turn gives less migration time under the same spare network bandwidth. The comparison of these algorithms in a data center with 300 virtual machines and 100 hosts is presented in Figure 3.6. The energy consumption by the data center using MBS-VM, RAM, and MM algorithms are 2.25 kWh, 3.30 kWh, and 3.87 kWh respectively. In our algorithm a , b , and c values are chosen randomly. We have simulated this experiment for 20 runs and the mean values of a , b , c are reported as follows: $a=0.37464$, $b=0.28705$, and $c=0.33817$. Figure 3.6 shows that our proposed approach performs better than the MM policy and the RAM policy with the increase in the number of virtual machines.

3.6 Summary

A data center is equipped with several thousands of physical machines, each of them handling several requests for virtual machines. An optimal resource allocation strategy for a data center helps achieving improved utilization of the data center. In this chapter, we presented three optimization approaches for VM placement and they achieved high utilization of physical machines by decreasing the number of active physical machines. We have also described a VM selection method for ef-

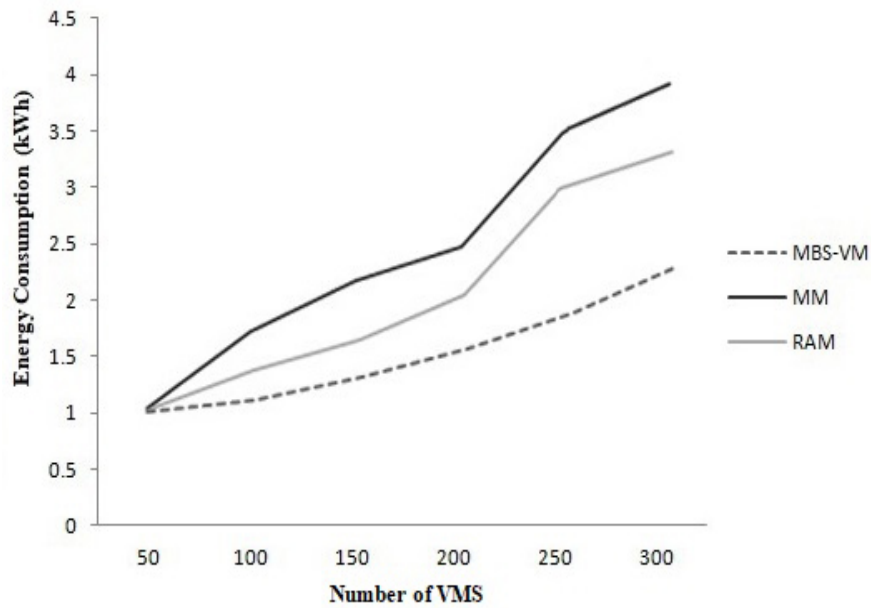


Figure 3.6: Comparison of VM selection algorithms

efficiently choosing virtual machines when a physical machine is either over-utilized or under-utilized. The experimental results illustrated that the combination of proposed placement and selection algorithms leads to significant reduction in energy consumption in a cloud data center up to 32%. Among our approaches, IBO has shown significant consistency in finding optimal solutions and achieved lowest convergence.

Chapter 4

Machine Learning Approaches for Forecasting Data Center Energy Demand

Energy costs are the fastest rising expense for today's data centers. The energy efficiency of a data center is influenced by many factors, such as data center layout design and characteristics, ambient weather conditions, rack density, the operation of HVAC systems and their behavior. This complex connection makes it hard to predict data center energy consumption. With sensor data and information about devices operations, forecasting energy consumption for data centers helps in planning and operations. Well-planned power consumption makes good Return on Investment (ROI) and elasticity of computing infrastructures. In this chapter, we forecast chiller energy consumption of a data center using Multi-Layer Feed Forward Neural Networks and Deep learning with Parallel Stochastic Gradient Descent. We evaluate the efficiency of proposed approaches using a power consumption data set collected over a period of 5 years from a data center.

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

Flexibility of model in multiple mechanisms, capability to work with high-dimensional data, and level of accuracy in prediction are needed for an applicable model. Neu-

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

ral networks are considered as more suitable tools for prediction of properties compared to usual correlations. Neural networks do not require a pre-assumption on the relationship between the input and output parameters and they are able to correlate parameters with any possible complexity [230, 231].

Deep learning changes its internal parameters to compute the intricate structure in each layer from the representation of the previous layer. Deep learning is an emerging area of machine learning, that learns from multiple levels of representations [232, 233]. Deep learning methods are representation-learning methods with multiple levels of abstraction, obtained by a deeper stack of coupling layers. Deep learning is making major advances discovering intricate structures in high-dimensional data. It has produced extremely promising results for solving complex problems in many domains such as signal processing, natural language understanding, etc. We see Deep learning as another possibly attainable strategy to forecast data center energy consumption. We describe two machine learning approaches for forecasting data center energy consumption: Multi-layer Feed Forward Neural Networks and Deep learning with Parallel Stochastic Gradient Descent.

4.1.1 Approach 1: Multi-layer Feed Forward Neural Networks (MFNN)

Multi-layer neural networks are composed of multiple levels of linear and non-linear operations which create accurate models of hierarchical feature extraction to learn from large-scale unlabeled representations of data. These models learn useful information of raw data and exhibits high performance. The basic unit of this architecture is the neuron. The weighted combination of inputs (α) is accumulated and passed through an activation function which generates an output signal ($f(\alpha)$), that is transmitted to connected neurons multiplied with the specific weight. Fig. 4.1 illustrates how information is processed through a single node.

$$\alpha = \sum_{i=1}^n w_i x_i + b \quad (4.1)$$

In Fig. 4.1, f represents a linear (OR) non linear activation function and bias b is the activation threshold.

We propose a Multi-layer Feed Forward Neural Networks (MFNN) for perform-

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

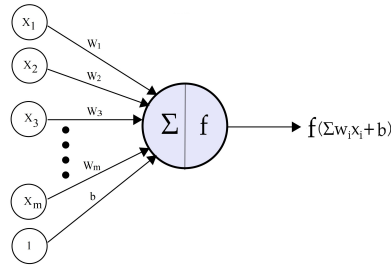


Figure 4.1: Structure of an Artificial Neuron

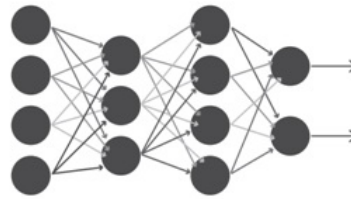


Figure 4.2: Basic structure of a four-layer feed forward network.

ing a forecasting task. The proposed MFNN comprises an input layer, two hidden layers, and an output layer as shown in Fig. 4.2. The width of each hidden layer is varied from 1 to 15 for finding the optimal configuration. A layer declaration for a trainable layer (the hidden or output layers) includes an activation function. Tanh, rectified linear, sigmoid and linear are the four different activation functions that can be used by the models. Definitions of the said activation functions are described as follows.

The hyperbolic tangent (Tanh) activation function is represented by:

$$f(\alpha) = \tanh(\alpha) = \frac{e^\alpha - e^{-\alpha}}{e^\alpha + e^{-\alpha}} \quad (4.2)$$

Rectified linear activation function is represented by:

$$f(\alpha) = \max(0, \alpha) \quad (4.3)$$

Sigmoid activation function refers to a special case of the logistic function and is defined by :

$$f(\alpha) = \frac{1}{1 + e^{-\alpha}} = \frac{e^\alpha}{e^\alpha + 1} \quad (4.4)$$

Linear activation function is a line of positive slope used to reflect the increase in firing rate that occurs as the input current increases. This gives a range of

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

activations. Thus, the linear activation function is defined as:

$$f(\alpha) = c\alpha \quad (4.5)$$

f becomes a straight line function where activation is proportional to input which is the weighted sum from neuron.

Our proposed model consists of a stack of interconnected neuron units. All intermediate layers of the network include biases. These biases and weights of linking neurons determine the output of the model. In this process, weights are fine-tuned to minimize the error on the labeled training data that in turn minimizes the loss function $Loss(W_j, B_j)$. Here, W is the weight matrix connecting two adjacent layers and B is the column vector of biases.

In our proposed model, we use a purely supervised training protocol for regression with MFNN for estimating the chiller consumption of a data center. We consider the Gaussian distribution for the response variable and the Mean Squared Error (MSE) for the loss functions. Equation 4.6 shows the calculation of the loss functions.

$$Loss(W_j, B_j) = \frac{1}{N} \sum_{j=1}^N (p^{(j)} - a^{(j)})^2 \quad (4.6)$$

where $a^{(j)}$ and $p^{(j)}$ are actual and predicted outputs respectively.

To minimize the loss function $Loss(W_j, B_j)$, we use the back-propagation (BP) training algorithm [234] where the weights connecting the layers are updated in an iterative manner. We initialize the weights of the network randomly. Then these weights are adjusted using the back-propagation algorithm, which follows the following four steps:

- Feed-forward computation : The input vector is presented to the network. We calculate and store the outputs of each node in the hidden layers. Then, we calculate the derivatives of the activation functions and store them at each node.
- Back-propagation to the output layer : By inspection, we collect all the multiplicative terms in the back propagation path from the output layer of the network to an intermediate output unit. Then we calculate the back propagated error. For calculating gradients for weight, we consider the following

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

3 values: node value feeding into weight, slope of the loss function with respect to node it fed into, and slope of the activation function with respect to output node [235].

- Back-propagation to the hidden layer : Back-propagation takes error from the output layer and propagates it back to the input layer. It calculates the necessary slope sequentially from the weights closest to the prediction through the hidden layers eventually back to the weights coming from the inputs.
- Weight updates : After calculating the partial derivative for each weight, we minimize the error function by performing a simple gradient descent as follows:

$$w_{ij}(t+1) = w_{ij}(t) - \epsilon \frac{\delta Loss}{\delta w_{ij}}(t)$$

where w_{ij} is the weight from the neuron j to the neuron i , and ϵ is the learning rate used to scale the derivative. We adjust the network weights iteratively to find a minimum of the error function, where the gradient of error function is sufficiently small. The training procedure of our proposed MFNN is described in Algorithm 4.1.

4.1.2 Approach 2: Deep Learning with Parallel Stochastic Gradient Descent (DPSGD)

Multi-layer feed forward neural networks are capable of performing just about any linear or nonlinear computation, and can approximate any reasonable function arbitrarily well. Such networks overcome the problems associated with perceptrons and linear networks. However, while the network being trained may be theoretically capable of performing correctly, backpropagation and its variations may not always find a optimal solution. Settling in a local minimum may be good or bad depending on how close the local minimum is to the global minimum and how low an error is required. In any case, we need to be cautioned that although a multilayer backpropagation network with enough neurons can implement just about any function, backpropagation will not always find the correct weights for the op-

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

Algorithm 4.1: Training a Multi-layer Feed Forward Neural Network

- 1 Clean the missing values in the dataset.
 - 2 Normalize the input data.
 - 3 Divide the dataset into training and test sets.
 - 4 Set up the network with three input units and fully connected to non-linear hidden units in hidden layer 1 and hidden layer 2 via connections with weights $w_{ij}^{(n)}$ and $w_{jk}^{(n)}$ respectively. Hidden units in layer 2 are in turn fully connected to output unit via connections with weights $w_{kl}^{(1)}$.
 - 5 Generate random initial weights.
 - 6 Select an appropriate error function and learning rate.
 - 7 **repeat**
 - 8 | Apply the weight update equation $\Delta w_{jk}^{(n)} = -\epsilon \frac{\delta Loss^{(n)}(t)}{\delta w_{ij}^{(n)}}$ to each weight $w_{jk}^{(n)}$ for each training sample.
 - 9 **until** *the network loss function is small enough* ;
-

timum solution. We need to reinitialize the network and retrain several times to guarantee that we have the best solution.

To overcome the said problems and to improve the accuracy of the predictions, we propose Deep Learning with Parallel Stochastic Gradient Descent (DPSGD), a purely supervised training protocol for regression, based on columnar compression and Map/Reduce with memory efficient java implementations. Multi-threaded and parallel computation is used on a single node where each core handles separate subsets or all of the training data.

The standard Stochastic Gradient Descent (SGD) computes the gradient $\nabla Loss(W_j, B_j)$ using backpropagation [236]. It is fast and memory-efficient but it is not parallelizable. The HOGWILD! approach proposed by Niu et al [237] is the lock-free parallelized version of SGD. In our proposed approach, to minimize the loss function $Loss(W_j, B_j)$, we used the lock free parallelized version of stochastic Gradient Descent. This algorithm follows a shared memory with multiple cores, each handling separate set of training data. Each core asynchronously contributes to the gradient update ($\nabla Loss(W_j, B_j)$) and is described in Algorithm 4.2.

For each iteration, we select an active subset (T_{na}) of training data. This subset is further divided into n_c (number of cores) partitions. For each core, the following

4.1 Machine Learning Approaches for Data Center Chiller Energy Consumption Prediction

Algorithm 4.2: Multi-threaded training with SGD

```

1 Initialize global model parameters  $W, B$ 
2 Distribute training data  $T$  across nodes (can be disjoint or replicated)
3 repeat
4   For nodes  $n$  with training subset  $T_n$ , do in parallel:
5     a. Obtain copy of the global model parameters  $W_n, B_n$ 
6     b. Select active subset  $T_{na} \subset T_n$  (user-given number of samples per
       iteration)
7     c. Partition  $T_{na}$  into  $T_{nac}$  by cores  $n_c$ 
8     d. For cores  $n_c$  on node  $n$ , do in parallel:
9       i. Get training example  $i \in T_{nac}$ 
10      ii. Update all weights  $w_{jk} \in W_n$ , biases  $b_{jk} \in B_n$ 
11           $w_{jk} := w_{jk} \alpha \frac{\nabla \text{Loss}(W_j, B_j)}{\nabla b_{jk}}$ 
12           $b_{jk} := b_{jk} \alpha$ 
13      Set  $W, B := \text{Avg}_n W_n, \text{Avg}_n B_n$ 
14      Optionally score the model on (potentially sampled) train/validation
       scoring sets
15 until convergence criterion reached ;

```

operations are performed in parallel in the inner for loop: (i) Read a training sample, (ii) Evaluate the equations in line 11 and 12 in Algorithm 4.2 and (iii) Update the weight and the bias.

The second operation does not change shared variables because it involves a series of arithmetic operations. However, for the first and the last operations, we use atomic instructions that are executed without considering the situation of other cores. Here, the weights and bias updates follow the asynchronous procedure to incrementally adjust the parameters W_n, B_n of each nodes after seeing each training sample. The final average of these local parameters across all cores is calculated to obtain the global model parameters. All available threads continuously execute the said procedure until the maximum number of iterations or the convergence criterion is reached.

The proposed deep learning framework supports regularization to prevent overfitting by modifying the loss function to minimize the loss according to Equation

4.7:

$$Loss'(W_j, B_j) = Loss(W_j, B_j) + \lambda_1 R_1(W_j, B_j) + \lambda_2 R_2(W_j, B_j) \quad (4.7)$$

where $R_1(W_j, B_j)$ is the sum of all l_1 norms for the weights and biases in the network and $R_2(W_j, B_j)$ is used to perform l_2 regularization which is the sum of squares of all the weights and biases in the network.

4.2 Description of Data Set and Preprocessing

For training the network, we used adequate data on measurements of a data center power consumption with a one-day sampling rate over a period of almost 5 years from January 2013 to October 2017. This data set has power consumptions of chillers, lighting, Air Handling Units (AHU), UPS, rack and loss incurred during the transformations. The dataset contains some missing values in the measurements. To handle the missing data, we used the mean substitution method where the mean is calculated for each variable over all examples in the data set.

Some inputs to neural network might not have a ‘naturally defined’ range of values. So, feeding the raw value to the neural network will not work very well. Hence, it is useful to normalize all the input and output data before this stage. Normalizing the input data avoids sticking in local optima and training becomes faster. For normalization of the input data, Min-Max normalizer, Binning normalizer, and Gaussian normalizer are generally considered [238].

Min-Max normalizer linearly transforms real data values such that the minimum and the maximum of the transformed data take certain values frequently 0 and 1 or -1 and 1. It is calculated as follows:

$$Y_i = (X_i - X_{min}) / (X_{max} - X_{min}). \quad (4.8)$$

where X_{max} and X_{min} are the maximum and minimum values of variable X_i in the training and testing data sets.

Gaussian normalization rescales the values of each feature to have mean 0 and variance 1. This normalizer calculates the mean and the standard deviation of a data set and normalizes by taking each data value, subtracting the mean, and then dividing by the standard deviation as follows :

$$Y_i = (X_i - mean) / Standard\ Deviation. \quad (4.9)$$

4.3 Experimental Analysis

Normalizer	MAE	RMSE	RAE	RSE	COD
Binning	10.811613	13.148259	0.788843	0.553419	0.446581
Min-Max	7.381487	9.83771	0.538572	0.309817	0.690183
Gaussian	8.217084	11.418274	0.599539	0.417367	0.582633
Without Normalization	15.313279	18.069382	1.117296	1.045212	-0.045212

Table 4.1: Comparison of different normalization functions

Binning normalizer creates the groups of equal size for the given input data and discovers a set of patterns in continuous variables. Once the bins are created, the information gets compressed into groups and then normalizes every value in each group to be divided by the total number of groups.

We performed experiments to find the best normalization method using a network with two hidden layers (with 2 and 4 nodes) varying the normalizers and the results are illustrated in Table 4.1 together with Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), Relative Absolute Error (RAE), Relative Squared Error (RSE), and Coefficient Of Determination (COD). From the results, we found that the developed network with Min-Max normalization has better performance compared to other normalization methods. Hence, we selected Min-Max normalization for designing our proposed methods.

4.3 Experimental Analysis

We train our models using 4 years data (2013-2016) and we forecast the energy consumption of chiller for the year 2017. We feed the network with last three days values of the chillers energy consumption to predict the next day value. Let a time series be $Y = \{y_1, y_2, \dots, y_n\}$. Then, the prediction can be obtained as $y_k = f(y_{k-1}, y_{k-2}, y_{k-3})$ where $f(\cdot)$ is a nonlinear function that maps previous three days values to current day value.

We perform experiments based on trial and error to determine the optimal arrangement of the feed forward network. We created, trained, and tested on a number of networks by varying number of neurons in each hidden layer, the activation function(s), and even the number of hidden layers. We designed the

4.3 Experimental Analysis

MFNN and DPSGD network proposed in Section 4.1 using Python running on an Intel Core i7, $2 \times 16\text{GB}$ DDR4, 1 TB SATA, and 2.60GHz Tyrone system with two cores. We have developed the MFNN and DPSGD models using the Tensorflow software (www.tensorflow.org).

Regression analysis is used to evaluate the network capability for chiller energy consumption prediction in a data center. The Root Mean Square Error (RMSE) is used as a measure to evaluate how the trained network estimation is correlated to the experimental data. We conducted several experiments to choose the best activation function for our data by changing the activation functions for a network containing 2 and 4 neurons in the first and second hidden layers respectively for both the proposed approaches. The results with various activation functions for the two proposed approaches are presented in Table 4.2 and Table 4.3. We observe that the softmax and tanh activation functions give better performance over other activation functions for the proposed MFNN and DPSGD approaches respectively.

Activation function	RMSE	COD
Tanh	9.3783	0.7184
sigmoid	9.03220	0.7388
linear	9.03220	0.7388
Rlinear	17.8983	0.02551
Softmax	8.14034	0.78787

Table 4.2: Comparison of different activation functions for MFNN

Activation function	RMSE	COD
Tanh	2.7923	0.9797
sigmoid	3.4652	0.9654
linear	4.2145	0.9543
Rlinear	3.0542	0.9878
Softmax	4.1939	0.9692

Table 4.3: Comparison of different activation functions for DPSGD

We choose the optimum performance for the network by changing the number

of neurons in the hidden layers iteratively for both approaches. The number of neurons are changed from 1 to 15 in the first and second hidden layers. The results of these networks are presented in Table 4.4 and Table 4.5, together with RMSE and COD where H1 and H2 indicate the number of neurons in the first and second hidden layer respectively. During this procedure, we selected the network with the least RMSE and the highest COD.

4.3.1 Performance Evaluation

Based on the experimental results, we observe that a neural network with two hidden layers gives better results compared to a network with one hidden-layer with the same number of neurons. We worked on deep networks with different structures and found better performance with a four-layer network that has two hidden layers. For the MFNN approach, the results of developed network with two hidden layers, containing 15 and 5 neurons, are better than the results of other neural networks, as presented in Table 4.4. For the DPSGD approach, the network with 7 and 6 neurons in hidden layers has shown better performance compared to other networks.

We compared the performance of our proposed MFNN and DPSGD models with boosted decision tree (BDT) regression, neural network, and linear regression approaches. DPSGD and MFNN are relatively advanced models and show better performance over other methods. The performance of the said models are given in Table 4.6. From Table 4.6, we can see that our proposed models have low Root Mean Square Error (RMSE) and high COD. This prediction accuracy is promising and comparable with the reported results. Fig. 4.3 and Fig. 4.4 present the prediction results of the said models for the test data and the training data respectively. These figures present the predicted versus the observed chiller energy consumption with varying loads. In Fig. 4.3 and Fig. 4.4, we see that the predicted values of the DPSGD approach has a similar pattern with the observed chiller energy consumption. Further, the proposed models have shown high accuracy in high and low power consumption conditions.

Statistical significance testing is an essential procedure to evaluate which model is best supported by the sample data [239]. Tests for statistical significance are used to check whether the prediction of the model is due only to random chance.

Table 4.4: Performance comparison of MFNN with two hidden layers

H1	H2	RMSE	COD	H1	H2	RMSE	COD	H1	H2	RMSE	COD	H1	H2	RMSE	COD
1	1	8.4503	0.7714	2	1	6.9985	0.8432	3	1	6.8373	0.8503	4	1	6.9056	0.8473
	2	8.4638	0.7707	2	2	6.9310	0.8462	3	2	6.6000	0.8606	4	2	6.3733	0.8700
	3	7.9632	0.7970	3	1	6.7110	0.8558	3	3	7.0271	0.8835	3	3	6.0333	0.8835
	4	7.6573	0.8123	4	1	6.7581	0.8538	4	4	6.7825	0.8527	4	4	6.4598	0.8664
	5	8.2840	0.7803	5	1	6.8978	0.8477	5	5	6.6879	0.8568	5	5	7.7434	0.8081
	6	7.9595	0.7972	6	1	7.4213	0.8237	6	6	6.9812	0.8440	6	6	5.5368	0.9019
	7	8.1016	0.7899	7	1	6.9070	0.8473	7	7	7.0223	0.8421	7	7	5.6110	0.8992
	8	7.9019	0.8001	8	1	6.7888	0.8525	8	8	6.9904	0.8436	8	8	6.2936	0.8732
	9	8.0846	0.7908	9	1	8.2426	0.7825	9	9	6.8228	0.8510	9	9	8.1052	0.7897
	10	7.7522	0.8076	10	1	7.0588	0.8405	10	10	7.1537	0.8362	10	10	6.3324	0.8716
	11	8.0781	0.7911	11	1	6.9086	0.8472	11	11	6.3725	0.8700	11	11	5.6576	0.8975
	12	7.9479	0.7978	12	1	7.3197	0.8285	12	12	6.1916	0.8773	12	12	5.5316	0.9020
	13	8.3807	0.7607	13	1	7.2357	0.8324	13	13	7.1537	0.8362	13	13	7.9985	0.7952
	14	7.9019	0.8001	14	1	7.5128	0.8193	14	14	7.5705	0.8165	14	14	6.0414	0.8832
	15	7.9696	0.7970	15	1	6.7789	0.8529	15	15	6.9094	0.8472	15	15	5.4900	0.9035
6	1	7.1741	0.8352	7	1	8.9324	0.7446	8	1	6.5897	0.8610	9	1	7.3887	0.8252
	2	5.5987	0.8997	2	2	6.4558	0.8666	2	2	7.7042	0.8100	2	2	7.1136	0.8380
	3	5.8624	0.8900	3	3	6.7330	0.8549	3	3	6.0716	0.8820	3	3	7.7218	0.8091
	4	7.4221	0.8237	4	4	5.9002	0.8886	4	4	5.8542	0.8903	4	4	6.6059	0.8603
	5	6.5571	0.8624	5	5	6.3541	0.8708	5	5	6.6693	0.8576	5	5	5.8313	0.8911
	6	7.7737	0.8065	6	6	5.2236	0.9127	6	6	5.4477	0.9050	6	6	5.7547	0.8940
	7	6.2077	0.8766	7	7	6.8535	0.8496	7	7	6.7184	0.8555	7	7	5.4770	0.9040
	8	6.4002	0.8689	8	8	4.9273	0.9223	8	8	6.5114	0.8643	8	8	6.3655	0.8703
	9	6.0215	0.8839	9	9	5.8685	0.8898	9	9	7.9015	0.8001	9	9	7.2874	0.8300
	10	6.9125	0.8470	10	10	5.6299	0.8985	10	10	5.9458	0.8968	10	10	5.6269	0.8986
	11	6.0425	0.8831	11	11	7.2747	0.8306	11	11	6.1603	0.8785	11	11	6.8083	0.8516
	12	6.8152	0.8513	12	12	5.7786	0.8931	12	12	7.3399	0.8275	12	12	7.0813	0.8395
	13	7.1982	0.8341	13	13	7.9748	0.7964	13	13	6.4072	0.8686	13	13	6.5575	0.8623
	14	5.5667	0.9008	14	14	5.4135	0.9062	14	14	7.5132	0.8193	14	14	6.4273	0.8678
	15	5.9984	0.8848	15	15	6.0939	0.8811	15	15	5.2676	0.9112	15	15	6.3961	0.8690
11	1	8.4799	0.7698	12	1	9.8214	0.6912	13	1	11.1433	0.6025	14	1	7.8668	0.8019
	2	5.8869	0.8891	2	2	7.5470	0.8177	2	2	7.4182	0.8238	2	2	11.3197	0.5898
	3	6.0553	0.8826	3	3	6.0952	0.8811	3	3	7.2144	0.8334	3	3	6.2621	0.8745
	4	6.0926	0.8812	4	4	6.4879	0.8653	4	4	5.2015	0.9134	4	4	9.2197	0.7279
	5	5.1336	0.9156	5	5	6.9179	0.8468	5	5	5.7681	0.8935	5	5	5.9145	0.8880
	6	6.4772	0.8657	6	6	7.6130	0.8145	6	6	6.6435	0.8587	6	6	6.4877	0.8653
	7	7.1940	0.8343	7	7	6.8820	0.8484	7	7	5.7090	0.8957	7	7	6.1354	0.8795
	8	5.3139	0.9096	8	8	7.9550	0.7974	8	8	7.0761	0.8397	8	8	6.5849	0.8612
	9	5.7176	0.8953	9	9	6.8851	0.8482	9	9	6.8021	0.8519	9	9	8.8610	0.7486
	10	8.0165	0.7943	10	10	5.4198	0.9060	10	10	7.2234	0.8330	10	10	5.5533	0.9013
	11	7.0616	0.8404	11	11	8.7865	0.7529	11	11	5.5020	0.9031	11	11	6.6630	0.8579
	12	5.5348	0.9019	12	12	5.8636	0.8899	12	12	6.0337	0.8835	12	12	5.8961	0.8887
	13	5.7276	0.8950	13	13	5.1193	0.9161	13	13	6.3575	0.8706	13	13	7.2662	0.8310
	14	5.7878	0.8928	14	14	6.1530	0.8788	14	14	5.8662	0.8898	14	14	5.6481	0.8979
	15	8.5106	0.7681	15	15	7.5376	0.8181	15	15	6.6269	0.8594	15	15	5.3274	0.9091

4.3 Experimental Analysis

Table 4.5: Performance comparison of DPSGD with two hidden layers

H1	H2	RMSE	COD	H1	H2	RMSE	COD	H1	H2	RMSE	COD	H1	H2	RMSE	COD
1	1	4.2041	0.9557	2	1	3.4002	0.9777	3	1	4.1099	0.9566	4	1	4.0643	0.9624
	2	4.2104	0.9571	2	2	2.9065	0.9888	3	2	2.8205	0.9917	4	2	3.4420	0.9675
	3	4.2209	0.9595	3	3	3.5395	0.9685	4	3	3.3259	0.9755	5	3	3.3259	0.9695
	4	4.2294	0.9614	4	4	2.8531	0.9872	5	4	3.1924	0.9836	6	4	3.2624	0.9703
	5	4.2349	0.9627	5	5	2.8145	0.9796	6	5	3.4348	0.9758	7	5	3.2771	0.9723
	6	3.2925	0.9750	6	6	2.5446	0.9898	7	6	2.5016	0.9876	8	6	2.3838	0.9945
	7	4.2401	0.9638	7	7	2.8463	0.9790	8	7	3.0699	0.9884	9	7	3.3610	0.9713
	8	4.2412	0.9640	8	8	3.1372	0.9738	9	8	3.1772	0.9795	10	8	3.2204	0.9795
	9	4.2419	0.9642	9	9	3.2745	0.9695	10	9	3.1390	0.9792	11	9	3.2733	0.9736
	10	4.2423	0.9643	10	10	3.3477	0.9643	11	10	3.1121	0.9795	12	10	3.3244	0.9753
	11	4.2427	0.9644	11	11	3.1590	0.9721	12	11	3.1604	0.9792	13	11	3.3146	0.9799
	12	4.2429	0.9644	12	12	3.2334	0.9715	13	12	3.2852	0.9756	14	12	3.3089	0.9825
	13	4.2431	0.9645	13	13	2.8076	0.9791	14	13	3.2072	0.9758	15	13	3.3240	0.9753
	14	4.2430	0.9644	14	14	2.7923	0.9797	15	14	3.3327	0.9747	16	14	3.2703	0.9726
	15	4.2429	0.9644	15	15	3.5111	0.9633	16	15	3.0770	0.9872	17	15	3.3013	0.9654
6	1	3.0365	0.9884	7	1	3.1351	0.9887	8	1	3.1491	0.9879	9	1	4.1983	0.9688
	2	3.2251	0.9853	2	2	3.1369	0.9877	3	2	3.2928	0.9762	4	2	3.6522	0.9625
	3	3.2385	0.9894	3	3	3.3475	0.9754	4	3	3.1484	0.9878	5	3	3.2813	0.9763
	4	3.3566	0.9719	4	4	3.3083	0.9772	5	4	3.2373	0.9785	6	4	3.2204	0.9850
	5	3.3033	0.9761	5	5	3.2900	0.9785	6	5	3.2469	0.9777	7	5	3.2659	0.9787
	6	2.3572	0.9925	6	6	2.3228	0.9937	7	6	2.3984	0.9943	8	6	2.3685	0.9925
	7	3.4075	0.9798	7	7	3.2554	0.9768	8	7	3.2617	0.9766	9	7	3.2387	0.9787
	8	3.2993	0.9832	8	8	3.2536	0.9772	9	8	3.3238	0.9855	10	8	3.2133	0.9796
	9	3.2939	0.9816	9	9	3.2880	0.9759	10	9	2.8759	0.9885	11	9	3.2516	0.9775
	10	3.2837	0.9785	10	10	3.1958	0.9783	11	10	3.2791	0.9753	12	10	3.2890	0.9793
	11	3.2466	0.9675	11	11	3.2620	0.9786	12	11	3.2452	0.9775	13	11	3.2891	0.9773
	12	3.3187	0.9890	12	12	3.2670	0.9785	13	12	3.2652	0.9763	14	12	3.2450	0.9777
	13	3.2481	0.9679	13	13	3.2648	0.9786	14	13	3.0689	0.9886	15	13	3.2504	0.9768
	14	3.1832	0.9868	14	14	3.2808	0.9743	15	14	3.0025	0.9892	16	14	3.2405	0.9785
	15	3.2303	0.9894	15	15	3.2240	0.9789	16	15	3.0542	0.9876	17	15	3.2232	0.9789
11	1	4.2039	0.9695	12	1	4.2040	0.9634	13	1	3.1543	0.9782	14	1	3.2626	0.9770
	2	3.4476	0.9741	2	2	3.1892	0.9777	3	2	3.3741	0.9643	4	2	3.0633	0.9857
	3	3.2169	0.9788	3	3	3.0299	0.9857	4	3	3.3567	0.9766	5	3	3.4124	0.9785
	4	3.3206	0.9700	4	4	3.4930	0.9686	5	4	3.2599	0.9728	6	4	3.1263	0.9898
	5	3.2072	0.9758	5	5	3.1759	0.9822	6	5	3.1903	0.9783	7	5	2.9461	0.9879
	6	2.3722	0.9920	6	6	2.4614	0.9915	7	6	2.4356	0.9917	8	6	2.4325	0.9917
	7	3.2360	0.9768	7	7	3.2658	0.9842	8	7	3.1799	0.9788	9	7	3.2550	0.9776
	8	3.2183	0.9798	8	8	3.1525	0.9781	9	8	3.2138	0.9786	10	8	3.2114	0.9798
	9	3.1933	0.9812	9	9	3.1656	0.9781	10	9	3.2153	0.9857	11	9	3.2432	0.9783
	10	3.2031	0.9831	10	10	3.0571	0.9728	11	10	3.1699	0.9781	12	10	3.0059	0.9884
	11	3.1940	0.9814	11	11	3.3361	0.9792	12	11	3.2317	0.9783	13	11	3.2564	0.9800
	12	3.2966	0.9840	12	12	3.1578	0.9814	13	12	3.1426	0.9794	14	12	3.2151	0.9853
	13	3.2090	0.9860	13	13	3.3141	0.9763	14	13	3.2074	0.9785	15	13	3.3062	0.9727
	14	3.2691	0.9848	14	14	3.2387	0.9785	15	14	3.2061	0.9788	16	14	3.2850	0.9733
	15	3.2598	0.9852	15	15	3.2626	0.9787	16	15	3.1543	0.9782	17	15	3.2293	0.9768

4.3 Experimental Analysis

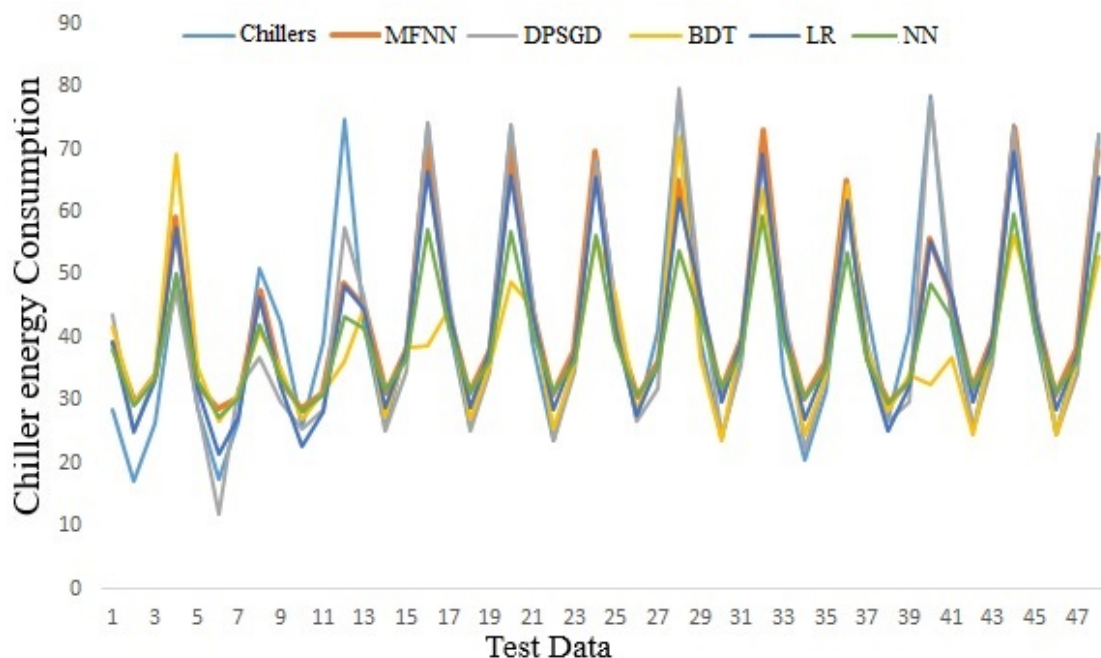


Figure 4.3: Experimental data vs predictions for test data

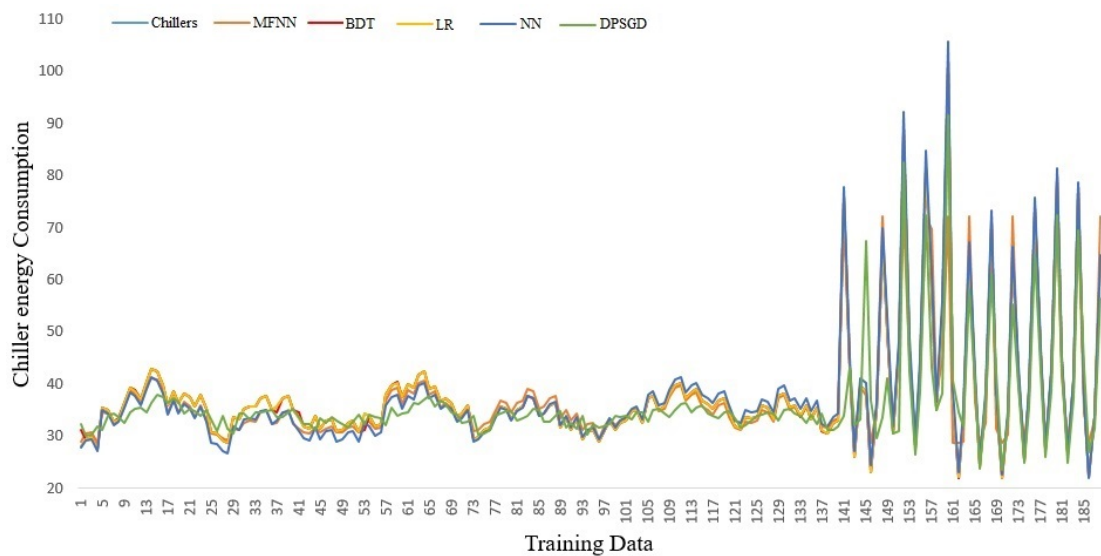


Figure 4.4: Experimental data vs predictions for training data

These tests tell us the probability of making an error if we use the said model. As our data is interval based and we work with more than one sample with a repeated-measures (within-subjects) design, we performed T-test [239] and Wilcoxon signed

	Root Mean Squared Error	Coefficient of Determination
NN	7.641626	0.813066
LR	10.329192	0.658453
BDT	7.458507	0.821917
MFNN	4.874361	0.923941
DPSGD	2.322854	0.993756

Table 4.6: Performance comparison of various prediction models

	MFNN	DPSGD	BDT	NN	LR
MFNN	*	*	*	*	*
DPSGD	$\frac{1.9127}{0.030364}$	*	*	*	*
BDT	$\frac{2.124239}{0.0423}$	$\frac{2.827}{0.003219}$	*	*	*
NN	$\frac{5.453751}{<0.00001}$	$\frac{5.268}{<0.00001}$	$\frac{-0.866085}{0.393552}$	*	*
LR	$\frac{3.925519}{0.00049}$	$\frac{4.587}{0.000012}$	$\frac{-0.047383}{0.962532}$	$\frac{1.470144}{0.15229}$	*

Table 4.7: Values for T and P for significance level of 0.05

rank test [240] for our data. We compared our methods with each of other methods by performing these statistical tests. The results obtained by T-test justify whether the predicted mean values of all approaches have a distinguished difference with 58 degrees of freedom at 0.05% level of significance. The statistical test values for all approaches are presented in Table 4.7. Each cell in the table represents $\frac{T-value}{P-value}$. Based on Table 4.7, the prediction values are statistically significant (all P-values are less than 0.00001). Similarly, we performed Wilcoxon Signed-Rank test and the results are presented in Table 4.8. Based on Table 4.8, the obtained values are statistically significant (all Z-values are obtained by positive ranks and P values are mostly 0.000).

4.4 Summary

Minimizing the energy consumptions is inevitable with the increasing of CO₂ emissions and rising energy prices in recent years. It became very important to predict the energy demands of the data centers. But the prediction task for data centers

4.4 Summary

	LR	NN	BDT	MFNN	DPSGD
LR	*	*	*	*	*
NN	Z = 4.7821 P = 0	*	*	*	*
BDT	Z = 2.3139 P = 0.02088	Z = 4.6587 P = 0	*	*	*
MFNN	Z = 4.7821 P = 0.	Z = 4.7821 P = 0	Z = 2.189 P = 0.02859	*	*
DPSGD	Z = 5.2621 P = 0.	Z = 4.9621 P = 0	Z = 4.2134 P = 0	Z = 3.564 P = 0	*

Table 4.8: Statistical analysis using Wilcoxon Signed Rank test (Z and P values)

involves complex numerical applications. In this chapter, two machine learning approaches are described for predicting the chiller energy consumption of a data center. The proposed methods have shown a real significance because of their reliability and accuracy in prediction of chiller energy. Among the proposed approaches, DPSGD has shown consistency in forecasting the chillers energy consumption and also achieved high accuracy.

Chapter 5

Metrics for Sustainable Data Centers

Most of us are familiar with the quote that if you cannot measure you cannot manage. In all fields, spanning technology and management, a set of metrics are established for measuring the stated objectives. There is a multitude of metrics available to analyze individual key performance indicators of data centers. In order to predict growth or set effective goals it is important to choose the correct metric and be aware of its expressivity and limitations. Continuous monitoring and measuring helps data center operators pro-actively identify and resolve potential issues, serve the growing business demands and thus improve the financial growth of an organization. This chapter describes a taxonomy of metrics for sustainable data centers and analyzes interrelationships among the metrics for data centers.

5.1 A Taxonomy of Data Center Metrics

Measuring how resources are used in a data center is crucial to understand the efficiency, reduce the costs of operations and achieve sustainability goals. Organizations are continuously searching for information and insights that offer control over their data centers. To remain competitive with their peers in the industry, they must ensure optimal utilization of resources in order to increase efficiency while minimizing environmental impact. This is only possible if there is information available that is meaningful and actionable. Organizations have to collate

5.1 A Taxonomy of Data Center Metrics

and analyze the information in order to evaluate themselves according standards and metrics. Properly defined and organized metrics will increase the organization productivity and assist in making the correct management decisions.

Data center design is at an evolutionary crossroad. Massive use of cloud services and challenging the physical limitations of power, cooling, and space are exerting pressure on the data center operators. Some of the challenges faced by the operators are as follows: enhancing resource usage, maintaining information assurance and providing high availability of service. Other challenges include reducing thermal inefficiencies and cooling utilization, decreasing the deployment time for new and existing services, enabling innovation through new cooling and consumption models, adopting on-site power generation and the proliferation of recycling in the architecture of data centers.

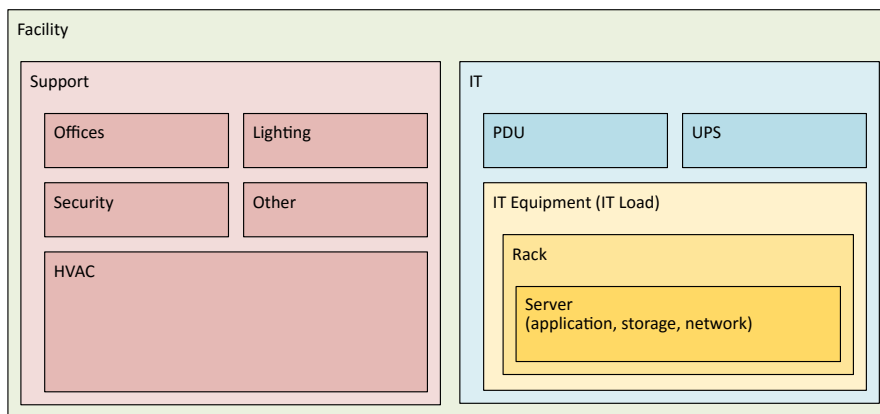


Figure 5.1: Categories of Components of Data Centers

In recent years data centers are experiencing a steady growth [241]. These data center require high availability and reliability for their daily operations, which in turn has an impact on the required resources. The next generation of data centers demand solutions that can lower the total cost of ownership and decrease the complexity of management. There are several techniques and tools to better utilize *always-on* data center infrastructure and reduce the recurring cost of Information Technology (IT) and facility management [242]. The operational costs of data centers are vastly different from those of other enterprises as less than 5% of the costs are personnel related. Servers are responsible for 45% of the amortized costs, followed by infrastructure (25%), power draw (15%) and networking (15%) [243].

5.1 A Taxonomy of Data Center Metrics

Understanding and analyzing data center metrics allows the operators to have a better view on how efficiency can be achieved in a data center by choosing the appropriate metrics for a given situation. Metrics also allow architects and operators to measure the performance and effect of changes made to subsystems in order to improve the efficiency of data center. A metric is generally defined as *the empirical, objective assignment of numbers, according to a rule derived from a model or theory, to attributes of objects or events with the intent of describing them* [244]. Poorly defined metrics could impede business innovation while still failing to meet environmental sustainability goals.

During last few years, significant research efforts and technological developments have been devoted to data centers targeting energy efficiency and eco-friendliness. The primary step in developing a model to capture the effects of data center management is to decide which dimensions are relevant, define the metrics, and populate the the metrics [245]. The Green Grid consortium proposed the Power Usage Effectiveness (PUE) [7], which currently is the most prevailing metric. The Green Grid consortium also proposed the Partial Power Usage Effectiveness (pPUE) which is based on PUE, and the Data Center Infrastructure Efficiency (DCiE) [246] which measures the efficiency of data centers by relating power consumption of data center and IT equipment. PUE and DCiE metrics help data center operators to know the efficiency of the data center where pPUE measures the energy efficiency of a zone in a data center. The consortium also proposed the metrics such as Carbon Usage Effectiveness (CUE) [247], Water Usage Effectiveness (WUE) [248], and Electronics Disposal Efficiency (EDE) [249] to measure the CO₂ footprint, the water consumption per year, and the disposal efficiency of the data centers respectively. Air flow performance in a data center plays an important role in improving cooling efficiency and space utilization and the metrics to monitor and control the air flow in a data center are discussed in [250, 251]. Munteanu et al. [252] proposed two different metrics based on energy consumption and Central Processing Unit (CPU) usage for calculating useful work done by Internet Data Centers (IDC). They proposed EnergyTIC Usage Effectiveness (EUE) considering total IDC power, IT power and load levels. They also proposed EUE(CPU), EUE(kWh) and EUE(kWh)-IT. Schaeppi et al.[253] explored energy related metrics for IT equipment, data storage and network equipment. Fiandrino et al.[254] proposed new metrics for computing energy efficiency

5.1 A Taxonomy of Data Center Metrics

of the data center communication systems, processes and protocols which includes communication network energy efficiency, network power usage effectiveness and network performance related metrics.

The European Union created an 8-project cluster of over 50 different partners in order to develop new environmental efficiency metrics and methodologies. The projects included are All4Green, CoolEmAll, GreenDataNet, RenewIT, GENiC, GEYSER, Dolfin and DC4Cities. The cluster has published multiple works [255, 256] in which they analyze existing metrics and also propose novel metrics to assess the performance of data centers. Capozzoli et al. [257] reviewed thermal, power and energy consumption metrics. Aravanis et al. [256] introduced new metrics for the assessment of flexibility and sustainability of data centers. Siso et al. [258] propose and evaluate several metrics for the CoolEmAll project.

Daim et al. [245] explored measurable components of a data center and proposed a new metric that fills the gap in measuring the data center equipment power and uses a credit-based system for data centers that do not meet the standard. Chen et al. [259] identified and presented usage-centric green performance indicators at various levels such as server and storage. Wang et al. [260] presented a set of performance metrics for a green data center. They focused on available benchmarks and how performance metrics can be used to measure the greenness of a data center.

Wiboonrat [261] discussed the effect of a data center outage and provided the solutions to minimize the data center downtime. He proposed improvements on the data center topologies to reduce the failure rate. American Society of Heating Refrigeration and Air Conditioning Engineers (ASHRAE) [262, 263] provides a common set of environmental guidelines for data processing environments, equipment and guidance on server metrics which enables data center operators to develop their own envelop that matches their business values.

To the best of our knowledge there is no such classification available which presents the dimensions of a data center from efficiency to security from the metrics perspective. Furthermore, we derive relationships between metrics, and discuss the advantages and disadvantages of each, in order to expose the research gaps and illustrate the latest research trends in computing the efficiency of a data center. We present a taxonomy of state of the art metrics used in the data center industry which is useful for the researchers and practitioners working on monitoring and

5.1 A Taxonomy of Data Center Metrics

Dimension	Metrics
Energy Efficiency	Power Usage Effectiveness (PUE), Data Center Infrastructure Efficiency (DCIE), partial PUE (pPUE), Corporate Average-Data Center Efficiency (CADE), Facility Efficiency (FE), Asset Efficiency (AE), IT-Power Usage Effectiveness (ITUE), Total-Power Usage Effectiveness (TUE), Operating System Workload Efficiency (OSWE), Deployed Hardware Utilization-Ratio (DH-UR), Deployed Hardware Utilization Efficiency (DH-UE), Datacentre Performance Per Energy (DPPE), Server compute efficiency (ScE), Data Center Compute Efficiency (DCcE), PAR ⁴ , Data Center Power Density (DCPD), Data Center Lighting Density (DCLD), Site Infrastructure Power Overhead Multiplier (SI-POM), IT Hardware Power-Overhead Multiplier (H-POM), Data Center Efficiency (DCE), Data Center Energy Productivity (DceP), Performance per-Watt (PpW), Compute Power Efficiency (CPE), Data Center Performance Efficiency (DCPE), Data center Workload Power-Efficiency (DWPE), Power Density Efficiency (PDE).
Cooling	HVAC System Effectiveness (HSE), Coefficient of Performance Ensemble (CoP), Data Center Cooling System Efficiency (DCSE), Data center Cooling System Sizing Factor (DCSSF), Air Economizer Utilization Factor (AEUF), Water Economizer Utilization Factor (WEUF), Recirculation Index (RI).
Greenness	Electronics Disposal Efficiency (EDE), The Green Index (TGI), Water Usage Energy (ω), Water Usage Effectiveness (WUE), Energy Reuse Effectiveness (ERE), Energy Reuse Factor (ERF), Material Recycling Ratio (MRR), Carbon Usage Effectiveness (CUE), Green Energy Coefficient (GEC), Technology Carbon Efficiency (TCE).
Performance	Data Center Energy Efficiency and Productivity Index (DEEPI), IT Productivity Per Embedded Watt (IT-PEW), Site-Infrastructure Energy Efficiency Ratio (SI-EER), Data Center Productivity (DCP), CPU usage, UPS losses, UPS-Power-Factor (PF), Crest Factor (CF) and Surge Factor (SF), Data Center Utilization (U_{Dc}), Server Utilization (U_{server}), SWaP: Space, Watts and Performance, ACE Score, UPS energy efficiency, FLOPS per Watt.
Thermal and Air Management	CRAC Flow (M_c), Negative Pressure Flow (M_n), Capture Index (CI), ByPass Air Flow (M_{bp}), Recirculation air flow (M_r), Return Heat Index (RHI), Supply Heat Index (SHI), β -index, Rack Cooling Index (RCI _{HI} , RCI _{LO}), Airflow Efficiency, Relative Humidity, Return Temperature Index (RTI), Data Center Temperature (T), Dew Point (DP), Heat Flux (HF), Heating Degree Days (HDD), Cooling Degree Days (CDD), Mahalanobis Generalized Distance (D^2).
Network	Communication Network Energy Efficiency (CNEE), Network Power Usage Effectiveness (NPUE), Maximum Relative-Size (RS_{max}), Diameter Stretch (DS), Path Stretch (PS), Network traffic/Kwh, Bits per Joule Capacity (BJC), Energy Consumption Rating Variable Load (ECR-VL), Telecommunications Energy Efficiency Ratio (TEER), Data center Fixed to Variable Energy Ratio metric (DC-FVER), Network Utilization ($U_{network}$).
Storage	Low-Cost Storage Percentage (LSP), Storage Usage, Slot Utilization, Overall Storage Efficiency (OSE), Throughput, Memory Usage, Response Time, Capacity, $U_{Storage}$.
Security	Connection Establishment Rate (CER), Latency, Concurrent Connections (CC), Illegal Traffic Handling (ITH), Application-Transaction Rate (ATR), Connection Tear down Rate (CTR), IP throughput, IP Fragmentation Handling - (IPHE), HTTP transfer rate.
Financial Impact	Total Cost of Ownership (TCO), Business value of Converged Infrastructure (BVCI), Mean Time To Repair (MTTR), Mean Time Between Failures (MTBF), Mean Time To Failure (MTTF), Reliability (λ), Availability (A), Carbon Credit, Return of Green Investment (RoGI).

Table 5.1: Taxonomy of Data Center Metrics

5.1 A Taxonomy of Data Center Metrics

improving the energy efficiency of data centers.

For efficient and eco-friendly operation of data centers, we need to monitor and measure all the components of a data center. These components are visualized in Figure 5.1. At the top level, we define the entire facility, which encompasses energy and other resources going into IT related components and into support components such as lighting, HVAC, and offices. The IT power flows to the Power Distribution System (PDS) and Uninterruptible Power Supply (UPS), which further distributes the power among IT equipments. The IT equipment consists of servers which are organized into racks. Servers can include application servers, networking equipment such as switches, routers, and storage servers. This division allows us to assign categories to each metric and group them based on these categories.

We propose a categorization of metrics based on the following dimensions (see Table 5.1): Energy Efficiency, Cooling, Greenness, Performance, Thermal and Air management, Network, Storage, Security and Financial Impact. There exists different metrics, each with their own approaches to measuring the efficiency of the data center, each with their own advantages, drawbacks and limitations. We describe the unit of each metric, objective, optimal value as well as the scale at which the metric operates. Objective specifies the optimization that should be done for the considered metrics (ex: minimize, maximize). Optimal value is a ideal or target value for the metric. Furthermore, there are interdependencies between different metrics as some are based on, or have a strong relationship with, other metrics. We provide an overview of the metrics collection for data centers. We also describe the challenges which are associated with certain metrics. We take a look at the relationships that exist between different metrics. The relationships between metrics can be defined as ‘uses’-relationship and ‘based on’-relationship. The ‘uses’-relationship exists when a metric uses another metric directly as input for the calculation. The ‘based on’-relationship indicates that a metric is based on the principles of another metric.

5.1.1 Energy Efficiency Metrics

Efficiency is defined as the ratio of useful work done by a system to the total energy delivered to the system. For data centers the energy efficiency translates to the useful work performed by different subsystems. This can be measured using energy

5.1 A Taxonomy of Data Center Metrics

Table 5.2: Energy Efficiency Metrics

Acronym	Full Name	Unit	Objective	Optimal	Category	Reference
APC	Adaptability Power Curve	Ratio	Maximize	1.0	Facility	[204]
CADE	Corporate Average Data Center Efficiency	Percentage	Maximize	1.0	Facility	[212]
CPE	Compute Power Efficiency	Percentage	Maximize	1.0	Facility	[213]
DCA	DCAdapt	Ratio	Minimize	$-\infty$	Facility	[204]
DCcE	Data Center Compute Efficiency	Percentage	Maximize	1.0	Server	[214]
DCeP	Data Center Energy Productivity	UW / kWh	Maximize	∞	Facility	[215]
DCiE	Data Center Infrastructure Efficiency	Percentage	Maximize	1.0	Facility	[194]
DCLD	Data Center Lighting Density	kW / ft ²	Minimize	0.0	Facility	[216]
DCPD	Data Center Power Density	kW / Rack	Maximize	∞	Rack	[216]
DCPE	Data Center Performance Efficiency	UW / Power	Maximize	∞	Facility	[217]
DC-FVER	Data Center Fixed to Variable Energy Ratio	Ratio	Minimize	1.0	Facility	[218]
DH-UE	Deployed Hardware Utilization Efficiency	Percentage	Maximize	1.0	Server	[219]
DH-UR	Deployed Hardware Utilization Ratio	Percentage	Maximize	1.0	Server	[219]
DPPE	Data Center Performance Per Energy	Ratio	Maximize	1.0	Facility	[220]
DWPE	Data center Workload Power Efficiency	Perf / Watt	Maximize	∞	Server	[221]
EES	Energy ExpenseS	Ratio	Maximize	1.0	Facility	[204]
EWB	Energy Wasted Ratio	Ratio	Minimize	0.0	Facility	[220]
GEC	Green Energy Coefficient	Percentage	Maximize	1.0	Facility	[220]
H-POM	IT Hardware Power Overhead Multiplier	Ratio	Minimize	1.0	IT Equipment	[219]
ITEE	IT Equipment Energy	Cap / kW	Maximize	∞	IT Equipment	[220]
ITEU	IT Equipment Utilization	Percentage	Maximize	1.0	IT Equipment	[220]
OSWE	Operating System Workload Efficiency	OS / kW	Maximize	∞	Facility	[222]
PDE	Power Density Efficiency	Percentage	Maximize	1.0	Rack	[223]
PEsavings	Primary Energy Savings	Percentage	Maximize	1.0	Rack	[204]
PUE ₁₋₄	Power Usage Effectiveness Level 1-4	Ratio	Minimize	1.0	Facility	[194, 203]
PUE _{scalability}	Power Usage Effectiveness Scalability	Percentage	Maximize	1.0	Facility	[224]
pPUE	Partial Power Usage Effectiveness	Ratio	Minimize	1.0	Facility	[194]
PpW	Performance per Watt	Perf / Watt	Maximize	∞	Server	[225]
ScE	Server Compute Efficiency	Percentage	Maximize	1.0	Server	[214]
SI-POM	Site Infrastructure Power Overhead Multiplier	Ratio	Minimize	1.0	Facility	[219]
SPUE	Server Power Usage Efficiency	Ratio	Minimize	1.0	Facility	[221]
SWaP	Space, Watts and Performance	Ratio	Maximize	∞	Rack	[226]
TUE	Total-Power Usage Effectiveness	Ratio	Minimize	1.0	Facility	[227]

5.1 A Taxonomy of Data Center Metrics

efficiency metrics. An overview of available energy efficiency metrics is presented in Table 5.2. The unit of each metric is listed, including the objective, optimal value and the category to which it belongs. We analyze the relationship between these metrics and present them in Figure 5.2. We organize the metrics horizontally based on their category in Figure 5.2 and visualize the relationships that exist among them. Figure 5.2 shows that the most popular energy efficiency metric, PUE, is used by a large number of other metrics either directly or as a derivation, e.g., Server Power Usage Efficiency (SPUE) and pPUE metrics are based on same principles as the PUE metric. The Data Center Performance Per Energy (DPPE) metric is also interesting as the metric is based on combining four other metrics namely DCiE, Green Energy Coefficient (GEC), IT Equipment Energy (ITEE), and IT Equipment Utilization (ITEU). Details and definitions of these metrics are given in Appendix B.1. To calculate the ITEU, one needs to know the exact power used by fans, voltage regulators and other components inside IT equipment. It is not clear how to measure the total energy that goes into IT equipment accurately. Defining coefficients for different types of IT equipment is also challenging especially in the heterogeneous environments of co-location data centers. To accurately calculate the Operating System Workload Efficiency (OSWE) metric the numbers of operating systems needs to be known, including operating systems in virtual machines. Thus, some of these metrics require accurate and very hardware specific data in order to be useful.

5.1.2 Cooling Metrics

The heat generated by the IT equipment in a data center must be controlled to maintain high levels of operational performance. Therefore, cooling plays a vital role in any data center. The complex interconnection of HVAC systems ensures optimal conditions for the computing environment in a data center, guaranteeing the life span, scalability and flexibility of the servers [264]. An overview of the available cooling metrics that can be applied in the context of data centers can be found in Table 5.3. Details and definitions of these metrics are given in Appendix B.2.

Table 5.3: Cooling Metrics

Acronym	Full Name	Unit	Objective	Optimal	Category	Reference
AEUF	Air Economizer Utilization Factor	Percentage	Maximize	1.0	HVAC	[228]
CoP	Coefficient of Performance Ensemble	Ratio	Maximize	∞	Facility	[229]
DCCSE	Data Center Cooling System Efficiency	kW/ton	Minimize	0.0	HVAC	[230]
DCSSF	Data center Cooling System Sizing Factor	Ratio	Minimize	1.0	HVAC	[230]
EER	Energy Efficiency Ratio	Ratio	Maximize	∞	Facility	[203]
HSE	HVAC System Effectiveness	Ratio	Maximize	3.5	HVAC	[231]
RI	Recirculation Index	Ratio	N/A	N/A	HVAC	[232]
WEUF	Water Economizer Utilization Factor	Percentage	Maximize	1.0	HVAC	[233]

5.1 A Taxonomy of Data Center Metrics

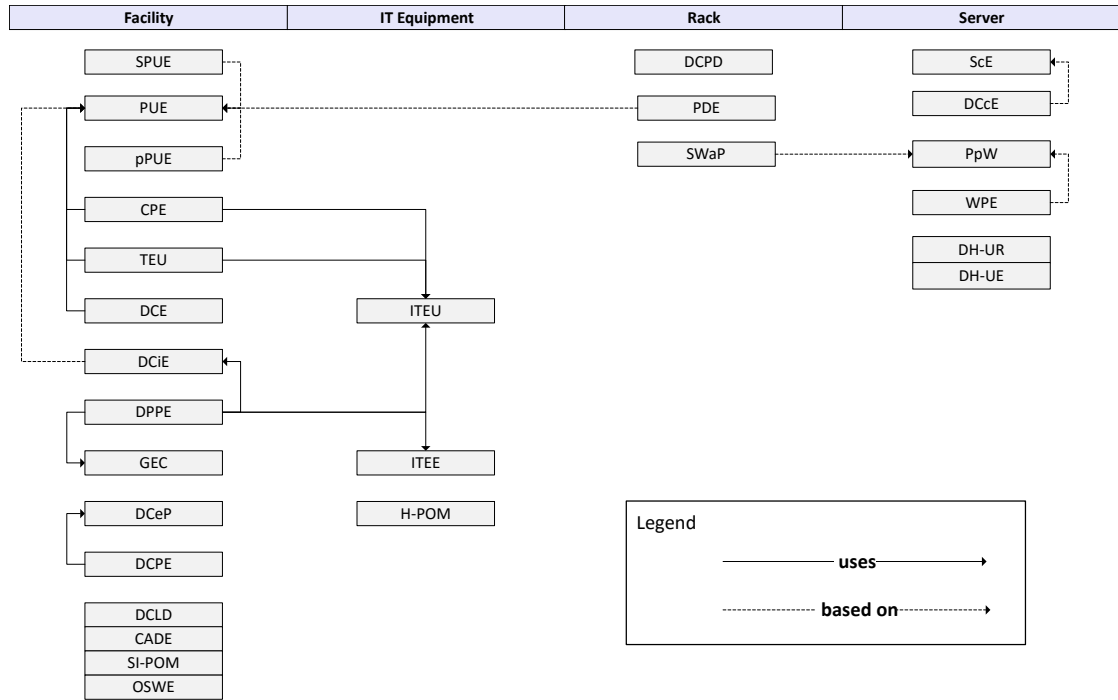


Figure 5.2: Relationships between Energy Efficiency Metrics

5.1.3 Greenness Metrics

Green Data centers are emerging, based on the design and operations of data centers in a more efficient and eco-friendly manner. Green initiatives and practices not only reduces GHG emission but also achieves measurable benefits. The carbon footprint and greenhouse gases are becoming subject to governmental regulations and taxes. As a result, the greenness of a data center is becoming increasingly important. A green data center is “a system in which the mechanical, lighting, electrical and IT equipment are designed for maximum energy efficiency and minimum environmental impact” [265, 266, 267]. Table 5.4 presents various green metrics which reflect the greenness of the data center in terms of reducing carbon footprint, reusing heat, the efficiency of water consumption and use of renewable resources. Figure 5.3 illustrates the relationships between these metrics considering four concepts: Reduce (reducing resources), Reuse (reusing resources), Recycle (recycling resources) and Renewable (use of renewable resources). In Figure 5.3, we organize the green metrics horizontally according to the these four concepts and vertically based on the category in which they operate. Details and definitions of

5.1 A Taxonomy of Data Center Metrics

Table 5.4: Green Metrics

Acronym	Full Name	Unit	Objective	Optimal	Category	Reference
-	CO ₂ Savings	Ratio	Maximize	1.0	Facility	204
CUE	Carbon Usage Effectiveness	KgCO ₂ /kWh	Minimize	0.0	Facility	195
EDE	Electronics Disposal Efficiency	Percentage	Maximize	1.0	Facility	197
ERE	Energy Reuse Effectiveness	Percentage	Minimize	0.0	Facility	234
ERF	Energy Reuse Factor	Percentage	Maximize	1.0	Facility	234
GEC	Green Energy Coefficient	Percentage	Maximize	1.0	Facility	220
GUF	Grid Utilization Factor	Percentage	Minimize	0.0	Facility	204
MRR	Material Recycling Ratio	Percentage	Maximize	1.0	Facility	235
Omega	Water Usage Energy / ω	Ratio	Minimize	0.0	Facility	236
TCE	Technology Carbon Efficiency	Pounds of CO ₂ /kWh	Minimize	0.0	Facility	[237]
TGI	The Green Index	Ratio			Facility	238
WUE	Water Usage Effectiveness	Liters/kWh	Minimize	0.0	Facility	196

5.1 A Taxonomy of Data Center Metrics

these metrics are given in Appendix B.3.

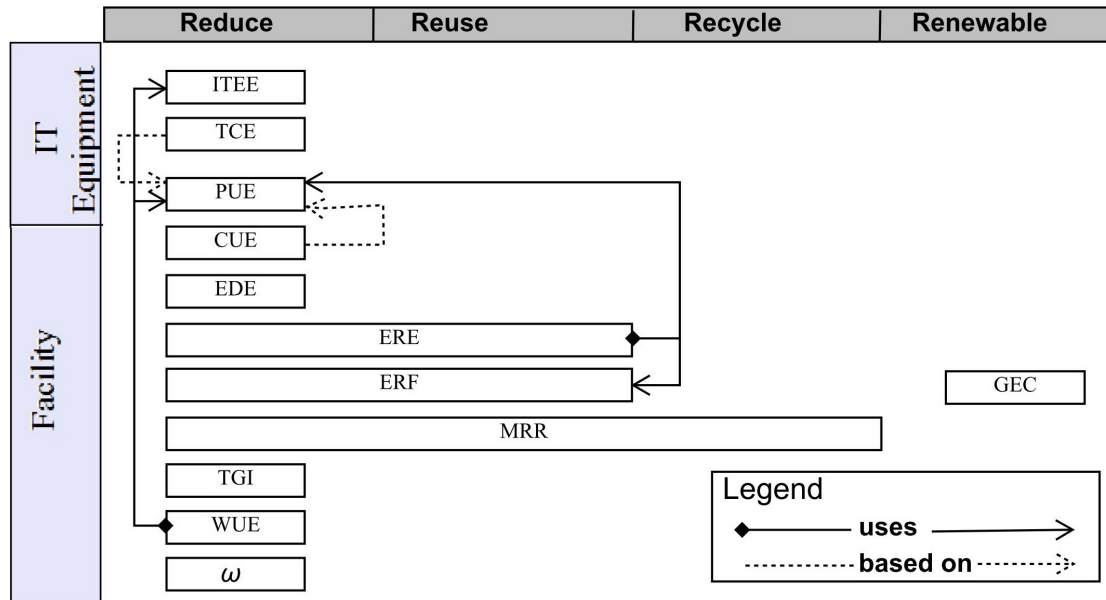


Figure 5.3: Relationship between Green Metrics

5.1.4 Performance / Productivity Metrics

The performance of a data center is the total effectiveness of the system, including throughput, response time, and availability [260]. Measuring performance and productivity is crucial as sub-optimal performance has operational and financial implications for a data center. When determining the performance of a data center one can encounter several difficulties, such as distinguishing significant workloads, overhead of performance measurements, energy distribution losses, and measuring the energy consumption at various levels of the data center. However, with the ability to measure performance and productivity, data center operators can determine how to improve the performance and plan for future work loads. An overview of the metrics which measure the performance of various components in data centers is presented in Table 5.5. Details and definitions of these metrics are given in Appendix B.4.

5.1 A Taxonomy of Data Center Metrics

Table 5.5: Performance Metrics

	Full Name	Unit	Objective	Opt.	Cat.	Reference
ACE	Availability, Capacity, and Efficiency Performance Score	Ratio	Maximize	1.0	HVAC	[243]
CPU	Central Processing Unit Usage	Percentage	Maximize	1.0	Server	[244]
DCP	Data Center Productivity	$\frac{Usefulwork}{Watt}$	Maximize	∞	Facility	[245]
DEEPI	Data Center Energy Efficiency and Productivity Index	Prod. / Watt	Maximize	∞	Facility	[246]
DR	Dynamic Range	Ratio	Maximize	1.0	Server	[247]
EP	Energy Proportionality	Ratio	Maximize	1.0	Server	[248]
FpW	Flops per Watt	Float. ops/Joule	Maximize	∞	Server	[244]
IPR	Idle-to-peak Power Ratio	Ratio	Minimize	0.0	Server	[249]
LD	Linear Deviation	Ratio	Minimize	0.0	Server	[247]
LDR	Linear Deviation Ratio	Ratio	Minimize	0.0	Server	[249]
PG	Proportionality Gap	Ratio	Minimize	0.0	Server	[247]
SWaP	Space, Watts and Performance	Ratio	Maximize	∞	Facility	[250]
U_{DC}	Data Center Utilization	Percentage	Maximize	1.0	Facility	[251]
U_{server}	Server Utilization	Percentage	Maximize	1.0	Server	[251]
UCF	Uninterruptible Power Supply Crest Factor	Ratio	Optimize	1.4	UPS	[252]
UPF	Uninterruptible Power Supply Power Factor	Ratio	Maximize	1.0	UPS	[252]
UPFC	Uninterruptible Power Supply Power Factor Corrected	Ratio	Maximize	1.0	UPS	[253]
UPS	Uninterruptible Power Supply Energy Efficiency	Percentage	Maximize	1.0	UPS	[254]
USF	Uninterruptible Power Supply Surge Factor	Ratio	Optimize	1.5	UPS	[252]

5.1.5 Thermal and Air Management Metrics

Cooling has been the major issue consuming nearly one third of the data centers energy consumption. High performance computing servers are bringing in new thermal and power challenges for data center operators. Data center operators must ensure minimum amount of energy for cooling which can be done through efficient air movement and using free cooling.

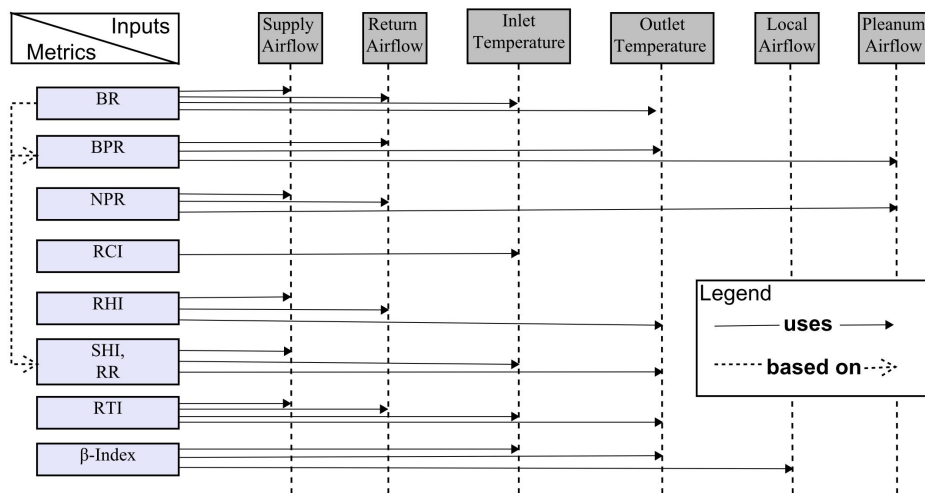


Figure 5.4: Relationship between Thermal and Air Management Metrics

Thermal and air management metrics are metrics which measure environmental conditions of the data center and also determine how air flows within a data center, from cooling units to the vents. These metrics assist with the diagnostic analysis to determine, for example, the amount of recirculation by-pass air. In general these metrics are based on the relationship between air flow rate and ambient temperature. They can be influenced by internal parameters and location [2]. Metrics like temperature, humidity, dew point and heat flux are used to prevent the over-heating, maintain the humidity levels, capture the current condition of the cooling system and to assist with making the correct decisions. The dimension, objective, optimal value of the outcomes, and the scale at which these metrics operate are presented in Table 5.6. Details and definitions of these metrics are given in Appendix B.5.

Air management metrics address air flow efficiency and separation of hot and

Table 5.6: Thermal and Air Management Metrics

	Full Name	Unit	Objective	Optimal	Category	Reference
-	Airflow Efficiency	W/cfm	Minimize	0.0	Facility	[230]
BPR	Bypass Ratio	ratio	N/A	N/A	Facility	[199]
BR	Balance Ratio	ratio	N/A	N/A	Facility	[199]
CI	Capture Index	Percentage	Maximize	1.0	HVAC	[255]
DC	Data Center Temperature	$^{\circ}C$ or $^{\circ}F$	Optimize	18 - 27 $^{\circ}C$	Facility	[210]
DP	Dew Point	$^{\circ}C$ or $^{\circ}F$	Optimize	17 $^{\circ}C$	Facility	[210]
HF	Heat Flux	W/m ²	Minimize	0.0	Facility	[210]
IoT	Imbalance of Temperature	Percentage	Minimize	0.0	Rack - Server	[203]
-	Mahalanobis Generalized Distance (D^2)	Unit	Minimize	0.0	Facility	[256]
M	Mass Flow $M_c, M_n, M_{bp}, M_r, M_s$	cfm	N/A	N/A	Facility	[257]
RCI	Rack Cooling Index	Percentage	Maximize	1.0	Rack	[258]
-	Relative Humidity	Percentage	Optimize	60%	Facility	[259]
RHI	Return Heat Index	Ratio	Maximize	1.0	Facility	[260]
RR	Recirculation Ratio	ratio	N/A	N/A	Facility	[199]
RTI	Return Temperature Index	Percentage	Optimize	1.0	Rack	[261]
SHI	Supply Heat Index	Ratio	Maximize	1.0	Facility	[260]
-	β -index	Ratio			Rack	[262]

5.1 A Taxonomy of Data Center Metrics

cold air streams. We observed that most of the air management metrics depend on common inputs. We analyzed these metrics considering the inputs, airflow path, and purpose of the metric and presented in Figure 5.4. Return Heat Index (RHI) and Supply Heat Index (SHI) differs in the way they have chosen the airflow path. Recirculation Ratio (RR) and SHI coincides with each other. Balance Ratio (BR) can be developed as a function of RR and Bypass Ratio (BPR).

5.1.6 Network Metrics

The data center network acts as a core component for providing numerous services. Networking equipment is responsible for up to 15% of a data center's amortized cost [243]. To increase the efficiency of data centers, operators should improve the energy efficiency of the network of data centers. The performance variability in network harms the application performance and causes the revenue loss. Data center network performance can typically be characterized using well known metrics such as bandwidth, NPUE, CNEE, reliability and throughput [268]. Some of these metrics are interrelated. An overview of the network metrics with the unit of each metric, objective, optimal value and the scale at which these metrics operate are presented in Table 5.7. Details and definitions of these metrics are given in Appendix B.6.

5.1.7 Storage Metrics

It is important to monitor and notify the measurements that boost efficiency to meet storage requirements of a data center. Conveying productive and enhanced storage execution for cloud data centers can be troublesome as it requires interaction with many components in the infrastructure such as application servers, storage devices, and network equipment. With a set of metrics for storage operations in the data centers, storage performance can be increased with continuous monitoring of these metrics [269]. Metrics such as OSE and slot utilization provide for better visibility into how proficiently our capacity is being utilized to store client information. Traditional metrics are unable to capture the improved efficiency achieved using new tools and methods such as trim storage and just-in-time allocations. We perceive the requirement for a single set of metrics that reflects storage utilization across changing technology base. We analyze and present the

Table 5.7: Network Metrics

Acronym	Metric	Unit	Objective	Optimal	Category	Reference
BJC	Bits per Joule Capacity	bits/joule	Maximize	∞	IT Equipment	[263]
CNEE	Communication Network Energy Efficiency	Joule/bit	Minimize	0.0	IT Equipment	[202]
DS	Diameter Stretch	Ratio	Optimize	1.0	IT Equipment	[264]
ECR-VL	Energy Consumption Rating Variable Load	Watts/Gbps	Minimize	0.0	IT Equipment	[265]
NPUE	Network Power Usage Effectiveness	Ratio	Minimize	1.0	IT Equipment	[202]
-	Network Traffic per Kilowatt-Hour	Bits / kWh	Maximize	∞	Facility	[266]
PS	Path Stretch	Ratio	Optimize	1.0	IT Equipment	[264]
RS _{max}	Maximum Relative Size	Ratio	Maximize	1.0	IT Equipment	[264]
TEER	Telecommunications Energy Efficiency Ratio	Ratio	Maximize	∞	IT Equipment	[267]
U _{network}	Network Utilization	Percentage	Maximize	1.0	IT Equipment	[251]

Table 5.8: Storage Metrics

Acronym	Metric	Unit	Objective	Optimal	Category	Reference
-	Capacity	GB / Watt	Maximize	∞	Storage	[268]
LSP	Low-cost Storage Percentage	Percentage	Maximize	1.0	Storage	[269]
-	Memory Usage	Ratio	Maximize	1.0	Storage	[270]
OSE	Overall Storage Efficiency	Ratio	Maximize	1.0	Storage	[269]
RT	Response Time	Milliseconds	Minimize	0.0	Storage	[268]
SU	Slot Utilization	Percentage	Maximize	1.0	Storage	[268]
-	Throughput	Bytes/second	Maximize	∞	Storage	[270]
$U_{storage}$	Storage Usage	Percentage	Maximize	1.0	Storage	[251]

5.1 A Taxonomy of Data Center Metrics

storage metrics along with their units as well as the objective, optimal value of the outcomes and the scale at which these metrics operate in Table 5.8. Details and definitions of these metrics are given in Appendix B.7.

5.1.8 Security Metrics

Security is one of the major concerns for business operations. As the data center houses the core assets of owners and data of clients, they must be safeguarded against physical as well as software threats. A firewall must have the capacity to handle the quickly advancing, network intensive service environment of the data center. Quality of the firewall policy is the key to any cyber defense perimeters. For security metrics, number of blocked attacks or intrusions detected will become the logical starting points. Table 5.9 lists the metrics for complexity and performance of firewalls, intrusion detection and prevention systems. Details and definitions of these metrics are given in Appendix B.8.

5.1.9 Financial Impact Metrics

Financial impact metrics describe the costs associated with designing and operation of a data center, financial impact of data center outage and return on investments on management tools and technologies for sustainable data center. Employing financial metrics into the data center provides many improvements in the IT part of the business. Employing financial metrics in the balanced score card will help the operators put other key metrics such as outage reports and service quality metrics. An overview of the financial impact metrics are presented in Table 5.10 in which the unit of each metric is listed including the objective, optimal value, and the category to which it belongs. Definition and detailed description of these metrics are given in Appendix B.9.

We analyze the relationship between these metrics and present them in Figure 5.5. Figure 5.5 shows that the Total Cost of Ownership (TCO), is calculated as a sum of Capital Expenditure (CapEx) and Operational Expenditure (OpEx) of the data center either directly or indirectly. Component failure rate (λ) and component repair rate (μ) are calculated using Mean Time Between Failures (MTBF) and Mean Time To Repair (MTTR) respectively.

5.1 A Taxonomy of Data Center Metrics

Table 5.9: Security Metrics

Acronym	Metric	Unit	Objective	Optimal	Category	Reference
ACPR	Average Comparisons Per Rule	Count	Minimize	0.0	IT Equipment	[272]
AS	Accessibility Surface	Count	Optimize	-	IT Equipment	[272]
ATR	Application Transaction Rate	Bits / sec	Maximize	∞	IT Equipment	[273]
CC	Concurrent Connections	Count	Maximize	∞	IT Equipment	[272]
CER	Connection Establishment Rate	Connections / sec	Maximize	∞	IT Equipment	[274]
CTR	Connection Tear down Rate	Connections / sec	Optimize	-	IT Equipment	[274]
DeD	Defense Depth	Count	Maximize	∞	Facility	[273]
DeP	Detection Performance	-	Maximize	1.0	IT Equipment	[275]
DTE	Data Transmission Exposure	Count	Minimize	0.0	IT Equipment	[275]
FC	Firewall Complexity	Ratio	Optimize	-	IT Equipment	[272]
-	HTTP Transfer Rate	Bits / sec	Maximize	0.0	IT Equipment	[274]
IAS	Interface Accessibility Surface	Count	Optimize	-	IT Equipment	[272]
IPFH	IP Fragmentation Handling	-	Maximize	∞	IT Equipment	[274]
-	IP throughput	Bits / sec	Maximize	∞	IT Equipment	[276]
ITH	Illegal Traffic Handling	Percentage	Maximize	∞	IT Equipment	[274]
-	Latency	Milli-seconds	Minimize	0.0	IT Equipment	[273]
RA	Rule Area	Count	Optimize	-	IT Equipment	[272]
RC	Reachability Count	Count	Minimize	0.0	Facility	[275]
RCD	Rogue Change Days	Days	Minimize	0.0	IT Equipment	[275]
T	Vulnerability Exposure	days	Minimize	0.0	IT Equipment	[275]

5.1 A Taxonomy of Data Center Metrics

Table 5.10: Financial Impact Metrics

Acronym	Metric	Unit	Objective	Optimal	Category	Reference
A	Availability	Ratio	Maximize	1.0	Facility	277
BVCI	Business Value of Converged Infrastructure	Dollars	Maximize	∞	Facility	[278]
CapEx	Capital Expenditure	Dollars	NA	NA	Facility	279
CCr	Carbon Credit	Tons of Carbon	Maximize	∞	Facility	268
MTBF	Mean Time Between Failures	Hours	Maximize	∞	Facility	280
MTTF	Mean Time To Failure	Hours	Maximize	∞	Storage	277
MTTR	Mean Time To Repair	Hours	Minimize	0.0	Facility	277
OpEx	Operational Expenditure	Dollars	Minimize	0.0	Facility	279
ROI	Return On Investment	Ratio	Maximize	∞	Facility	281
TCO	Total Cost of Ownership	Dollars	NA	NA	Facility	282
λ	Reliability	Faults / Hour	Minimize	0.0	Facility	277

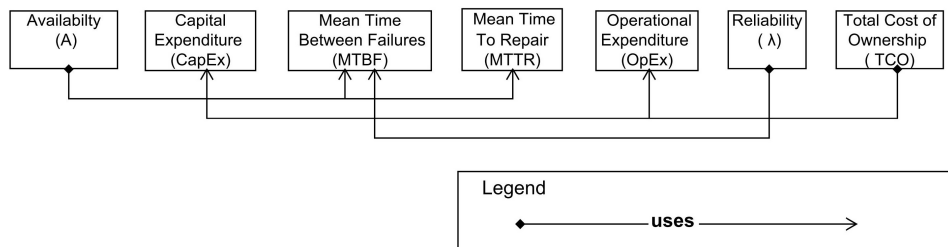


Figure 5.5: Relationship between Financial Metrics

5.2 Analysis of Metrics

There are a multitude of metrics to measure and monitor different aspects of data centers. When looking at the relationship between the metrics and challenges associated with using them, it becomes apparent that there is no single metric which covers all dimensions of the data center's performance. Even per dimension, there are several metrics promising to provide insight into the same area, through similar or different methods. However, none of the metrics are designed for comparing data centers amongst each other. Although the PUE metric is currently used for this purpose, it was never intended to be used as a comparison metric [7]. Instead the metric was envisioned to be an internal measurement to steer the data center towards higher levels of efficiency, by knowing which areas have a low efficiency in terms of energy consumption. For example, the IT load of a data center influences the PUE significantly. Furthermore, the PUE is also influenced by the weather and the location of the data center. Therefore, comparisons between data centers using PUE are most often not representative of the actual situation.

It is not possible for a single metric to represent the energy efficiency for all of the possible combinations of an IT environment. The Corporate Average Data Center Efficiency (CADE) metric can be extended by considering how efficiently servers, storage, and network equipment are utilized. Data Center Infrastructure Efficiency (DCiE) metric is effective at discovering the initial problem and helps justify the need to implement energy saving changes. However, the DCiE metric varies for each data center as it depends on the IT electrical load, which is a variable and site specific function of the IT software, architecture, hardware, load and efficiency. Due to this variability, we can not predict the impact of changes to

the data center using DCiE. The Green Index (TGI) metric allows for flexibility in green benchmarking as it can be used and viewed in different ways by its end-users. Even though we have specified the performance-per-watt metric for computing TGI, it can be computed with any other energy-efficient metric. TGI does not consider the power consumed outside of the IT equipment context. Therefore it can be extended by including components such as the cooling infrastructure.

Overhead metrics such as IT Hardware Power Overhead Multiplier (H-POM), Site Infrastructure Power Overhead Multiplier (SI-POM) give an understanding of a data center's energy use considering variations in IT equipment energy and the committed power to a facility. These metrics provide useful insights to the operators where modular provisioning is used. Because of the complexity and the unpredictable nature of data centers, a credible quantitative measure of security risk is not currently feasible. Security managers should chose a set of metrics which allows for better decision making and actual security improvements. The specific metrics discussed in this work can be refined and expanded to reduce the risk of a successful cyber-attack. The study of the given metrics has identified the need for improved measurement tools.

The number of inter-dependencies between different metrics on the facility level is large. It is important to be aware of these relationships, as metrics can have certain limitations that affect other metrics associated with them. When combining existing metrics into new ones, or basing new metrics on existing ones, the flaws of the existing metrics are usually not overcome, and sometimes even increased. Therefore, it is useful to understand these shortcomings and know what a metric can and cannot measure. Applying metrics is even more difficult for co-location data centers as the equipment, space and bandwidth are available for rental in these types of data centers.

Furthermore, there is a need for a metric which is designed with comparison in mind from its inception. Ideally, this metric should attempt to normalize the data in such a way that a fair comparison between data centers can be performed. The metric should take into account the utilization of data centers, as well as the location and weather. A metric which is highly dependent on those factors is PUE: a change in the utilization efficiency of the data center is immediately reflected in the value variations of the PUE. The location of the data center also has an influence on the outcome of this metric, as does the weather.

By applying a well-defined set of metrics which measure energy consumption and environmental impact during data center operation, and while making choices at various levels, it is possible for data centers to be planned, designed, implemented and operated in an energy-aware and more eco-friendly manner. Energy efficiency metrics measure the computing and non-computing energy used in a data center. These metrics measure the efficiency at various levels of granularity starting from operating system to data center. But it is difficult to measure the energy consumption at operating system level. Also, it is challenging to measure the energy consumption at sub-component level of a data center, as these low level measurements are often not available. Cooling metrics are used to specify the performance of the Computer Room Air Conditioning (CRAC) / Heating, Ventilation and Air Conditioning (HVAC) units and proper sizing of the cooling units. These metrics also measure the efficiency of the cooling systems. Estimating power and cooling capacity requirements using the ratings found on nameplates of IT equipment may not be accurate. Another issue is that heat densities change within racks, and also differ from one rack to the next.

Green metrics measure the environmental impact of a data center and its components. They highlight the importance of green energy and measure the efficiency of recycling and reuse in a data center. Efficient measurement of these metrics require capturing regional and seasonal changes to enable comparison of different data centers. A data center can increase its productivity by clearly defining performance metrics. These metrics help to measure IT performance and productivity of the data center and also identify problem areas. Metrics can range from low level UPS performance to high level data center utilization. Across all the components, a single fault may affect many other systems and ultimately decrease the overall performance of the data center. Operators rely on nameplate capacities and modelled load which do not accurately represent the actual capacity requirements. It is challenging to understand the impact of changes that are made in real-time.

Thermal and Air Management metrics monitor environmental conditions inside the data center. These metrics give an overview of how efficiently air flows within a data center and also quantify the extent of cold and hot air mixing. Continuous monitoring of these metrics allows the operators to reduce fan speed and increase cooling set points in real-time, which increases cooling efficiency and energy savings. It is difficult to determine the correct values for temperature and humidity

in the data center, as the environment is dynamic and constantly changing. Network metrics cover the network energy efficiency, network utilization and traffic demands of a data center. Networking equipment is responsible for a large portion of a data center's energy consumption, therefore it is important to optimize the efficiency of the equipment.

Storage metrics capture the performance of storage operations. These metrics assist the operators in reducing storage cost, improving storage utilization and increasing the overall storage performance. The distributed nature of cloud computing makes it critical to learn what workloads customers are accessing and the level of importance of the accessed data. Security metrics cover aspects such as the firewall performance. These metrics are highly dependent on internal governance, compliance standards and service level agreements of the data center in question. Another issue is authorization: the visibility of and control over resources in a data center. Financial Impact metrics help achieve a data center's financial and strategic objectives. These metrics range from total cost of ownership to return on investments. Measuring business value may vary from one organization to another due to different definitions, and Carbon Credit may vary based on a country's policies.

5.3 Summary

Metrics are important for planning, designing, building and operating a data center in an efficient manner. Our classification of metrics provides deep insights into the state-of-the-art of measuring different data center components and allows for quick access to the right subset of metrics from a huge collection that fits the desired context.

We observed that existing metrics are mainly focused on measuring the energy efficiency of IT equipment or facilities. Older facilities may not be able to capture the raw data that feeds today's more sophisticated metrics. There are very few metrics defined which can integrate different components of a data center that have a single numerical value to report the efficiency of the data center in all perspectives. Also, there is no metric which reflects the changes made to a data center and its sub-components. Furthermore, there is a need for new metrics that consider different factors such as the location and age of the data center, in order

to allow comparison across different data centers.

Chapter 6

Best Practices for Sustainable Data Centers

Data centers are an essential utility of modern societal infrastructures. Data centers are responsible for an important share of global energy consumptions, which could be up to 2% according to some studies. It is thus important to consider any opportunity to reduce the energy consumption of the data centers, both in design and operations. In this chapter, we analyzed seven data centers in India and the Netherlands and, based on our findings and industry standards, we propose a set of best practices to improve the energy efficiency of the data centers which spans the categories of Energy Efficiency, Cooling, Thermal, and Air management, Greenness, Storage and Networks. Following some of these best practices, data centers surveyed in our study have achieved 10-20% improvements on their energy consumption. The present study targets IT professionals, operators, and managers, providing efficient alternatives in daily operations of the data centers and identifying cost saving opportunities.

6.1 Research Methodology

Data centers constitute the core of an organization's information system by centralizing IT operations and equipment. Data centers experiencing steady growth in last decade. It demands high availability and reliability, straining the resources which may lead to the poor performance. The growth in hyper-scale cloud data

6.2 Data Center Management Best Practices

centers is one of the major contributors for the increase in electricity consumption across the globe. Next generation data centers demand solutions with lower Total Cost of Ownership (TCO) and take down complexity of management. There are several techniques and tools to make better utilization of "always-on" data center infrastructure and reduce the recurring cost of IT and facility management.

This chapter is based on the multiple case research approach [270, 271]. The features and infrastructure details of the data centers used for the study are presented in Table 6.1. Among the seven data centers, two are colocation data centers and provide various services including public cloud services. Three data centers are privately owned by companies and used in the financial services sector running banking, financial and insurance applications. The remaining two data centers are data centers of academic institutions. A uniform and standard data collection methodology was adopted in each case which included a standard questionnaire, review of procedures, benchmarks (to confirm gaps) and key staff interviews (5-8 key personnel having different designations in each data center). The standard questionnaire was based on the assessment of the following key dimensions: Energy Efficiency, Cooling, Thermal and Air management, Greenness, and Network and Storage. For each dimension, issues and practices were monitored directly in the data center for the purpose of the present study. Based on interview transcripts, we developed an ad-hoc case study report. The study report was then distributed and discussed with the interviewees and other staff to gain insights and tailored feedback for the correct understanding of the status of the center.

6.2 Data Center Management Best Practices

We consider the whole chain of operations of a data center and study the best practices for sustainability followed in the data center across the following dimensions: Energy efficiency, Cooling, Air and Thermal management, Greenness, Storage and Network. Data center users are interested in the performance of their applications at different scales of utilization, whereas data center operators are interested in efficiency of the resources.

6.2 Data Center Management Best Practices

Table 6.1: Data centers configurations

	Rack Space	Data Center Space (in SFT)	Power Capacity	Security level (Zones)	TIER	PUE
Data Center 1	5,000	230,000	30 MW	8	4	1.6
Data Center 2	1,400	40,000	10 MW	6	4	1.7
Data Center 3	800	20,000	6 MW	6	3	1.6
Data Center 4	3,000	125,000	20 MW	6	3	1.6
Data Center 5	1,000	30,000	10 MW	6	3	1.7
Data Center 6	160	3,500	450 kW	Redacted	2	1.25
Data Center 7	100	3,000	300 kW	Redacted	2	1.25

6.2.1 Energy Efficiency Practices

Efficiency is defined as the ratio of the useful work done by a system to the total energy delivered to it. For data centers, energy efficiency translates to the useful work performed by different subsystems. Following are some of the key steps to achieve energy efficiency in a data center.

6.2.1.1 Employ Automation Tools

Data center automation tools help automating tasks such as provisioning, configuration, patching, release management and compliance. Most of the data centers that we have studied rely on automation tools that enable real-time optimization, reduce error rates and improve the performance of the applications.

6.2.1.2 Use Virtualization and Consolidation

Virtualization enables abstracting physical servers in a data center facility along with storage, networking and other infrastructure devices and equipment. Consolidation combines workloads from different machines into a smaller number of systems when servers are under-utilized and consuming more energy [272, 273]. All the data centers in our study are virtualized and use different virtual machine consolidation and placement techniques to reduce power consumption and improve server utilization.

6.2.1.3 Dynamic Voltage and Frequency Scaling (DVFS)

DVFS reduces the power consumption of a processor on the fly by adjusting clock frequency according to current load indirectly showing the reduction in the supply voltage [274]. Power-saving can be achieved either by scheduling schemes with the capability of dynamic voltage and frequency or by consolidation techniques. Decommission servers without any computing Comatose / Zombie servers are those that run applications no longer required or unused, yet plugged in and operating continuously. Data centers operators have to audit and root out comatose servers and duplicate applications which account for up to 30% of the entire servers deployed. Based on our study, we observed that decommissioning of servers results in 50% energy savings.

6.2.1.4 Controlled Lighting

Installing a lighting control system in conjunction with more efficient fixtures and occupancy sensors can help reduce energy usage. Only three of the data centers in our study are using resource-friendly timers that dim or shut off lighting when people are not present.

6.2.1.5 On-site Power Generation

The critical need for clean and economical sources of energy is transforming data centers that are primarily energy consumers to also energy producers. On-site renewable power generation is an economical and eco-friendly solution for regions with high electricity and low natural gas prices and for campus-like facilities that

6.2 Data Center Management Best Practices

can re-utilize excess heating and cooling [275, 276]. In these cases, one can utilize the grid power as a backup in combination with on-site generation systems such as solar panels and fuel cells as the primary source. However, none of the data centers in our study use on-site power generation.

6.2.1.6 Integrated Monitoring

Energy monitoring allows greater visibility into overall data center energy usage while providing solutions to maximize server and infrastructure equipment operating efficiency. Many data centers use Data Center Infrastructure Management (DCIM) tools to monitor energy and cooling efficiency and claim achieving 20% savings in operational expenses [277].

Table 6.2 summarizes the practices for energy efficiency followed by the seven data centers (DC1, DC2, ..., DC7) in our study.

6.2.2 Cooling, Thermal, and Air Management Practices

The data center cooling system can be considered a giant air pump where the cooled air flows finally reaching the server inlets [278]. Therefore, energy efficiency of a data center can be improved by using better air-management practices in data centers with raised floor designs as well as non-raised floor air-conditioning designs as follows.

6.2.2.1 Using Central Air Handler

Efficient airflow can be achieved only by eliminating bypass and recirculation air flows as this is where the air flow is wasted in data center. For efficient air-management, some of the data centers under study make use of custom designed central air handler systems using variable speed drives. Further, many of the data centers in our study use loop design using median temperature of 10-15°C and load monitoring sensors for chiller plants [279].

6.2.2.2 Use Hot Aisle / Cold Aisle Containment

In the aisle containment approach (see Figure 6.1), all the hardware in a row of cabinets faces the same way, so hot air is expelled on one side while cold air blows

6.2 Data Center Management Best Practices

Table 6.2: Energy efficiency practices for data centers

	DC1	DC2	DC3	DC4	DC5	DC6	DC7
Automation tools	✓	✓	✓	✓	✓	X	X
Virtualization & consolidation	✓	✓	✓	✓	✓	✓	✓
Dynamic Voltage and Frequency Scaling	✓	✓	✓	✓	✓	✓	✓
Handling Comatose or Zombie servers	✓	✓	✓	✓	✓	✓	✓
Raised thermostat set point	✓	X	X	✓	X	✓	✓
Controlled lighting with sensors or other technologies	✓	✓	X	✓	X	X	X
On-Site power plant	X	X	X	X	X	X	X
DCIM tools	✓	✓	X	✓	X	✓	✓

from the other side. All of the data centers in our study follow this containment approach that allows proper flow of cold air to the destination, in turn reducing energy consumption.

6.2.2.3 Implementing Liquid Cooling Solutions

Liquid cooling solutions provide effective cooling and isolate equipment from the existing cooling system using a liquid in room-level and row-level systems. Only one data center in our study is using liquid cooling as it is more expensive than air based ones.

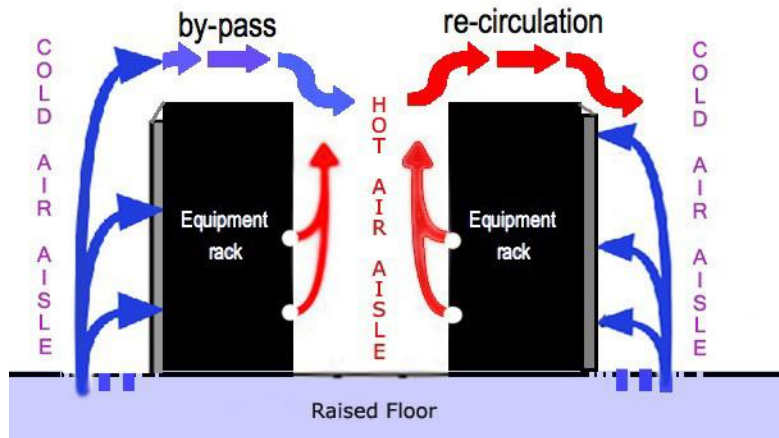


Figure 6.1: Hot Aisle / Cold Aisle Containment in Data Centers

6.2.2.4 Server Inlet Temperature and Humidity Control

The control of humidity in data centers is essential to achieve high availability and reduce maintenance costs. The level recommended is around 50% or higher. However, data centers without large high-speed fans can safely operate at 40% humidity levels, thus decreasing water and energy consumption [280]. Humidity can be best controlled knowing both inside and outside environmental conditions. Adiabatic humidification technology provides high efficiency than infrared or isothermal technologies [281, 282].

6.2.2.5 Air-flow Management

All the data centers in our study use horizontal, vertical, under rack panels and PVC curtains for isolation with the final goal of minimizing the recirculation of hot air. All the data centers in our study also use high raised-floors, overhead cabling, cable grouping, the placing of cable trays below the hot aisle, and cabling within the cabinets and racks to avoid air blockages. The best practice is to have dedicated horizontal airflows rather than a mixture of vertical and horizontal airflows because dedicated horizontal airflows provide much more uniform airflow distribution compared with the mixture. During the inspection of the return air ducts for HVAC, we observed inadequate ceiling height or undersized hot air return plenum in few data centers. Increasing the size of the return duct to match the air handler avoids this problem. Use of high overhead plenum and several feet of

6.2 Data Center Management Best Practices

Table 6.3: Cooling, Thermal, and Air management practices for data centers

	DC1	DC2	DC3	DC4	DC5	DC6	DC7
Central air handler	✓	X	✓	✓	X	✓	✓
Liquid cooling	✓	X	X	X	X	X	X
Sensors for chiller plant	✓	✓	✓	✓	✓	✓	✓
Hot aisle / Cold aisle containment	✓	✓	✓	✓	✓	✓	✓
Loop design for chillers	✓	X	✓	✓	✓	✓	✓
Adjustable speed drive chillers	✓	✓	✓	X	X	X	X

clearance under the raised-floor give better maintenance. Variable frequency fans in the Computer Room Air Handler (CRAH) units would allow for self adjusting, thus resulting in energy savings. Most of the data centers under study are operating at baseline temperature but raising the baseline temperature would save 4% in energy costs with each degree of increase in the set point [142]. Table 6.3 presents the practices for cooling, thermal, and air management followed in the data centers under our study.

6.2.3 Green Practices

A Green data center is one for which design requirements include energy efficiency non-functional ones and with precise and measurable power efficiency and sustainability targets for its operation. This entails the incorporation of energy-efficient design together with high efficiency power delivery, high efficiency cooling, and increased utilization of renewable energy sources [267, 283]. Some of the approaches and practices for greening data centers and IT are as follows.

6.2.3.1 Make Use of Economizers (Free Cooling)

Data centers can achieve significant energy savings either through water-side or air-side economizers. Economizers have impact when wet bulb temperatures outside the data center are less than 13°C for more than 3,000 hours/year. Some of the data centers in our study claim up to a 20% and 7% decrease in energy costs and

maintenance costs respectively since deploying economizers. However, the use of economizers depends on the geography, site conditions and economizer design.

6.2.3.2 Using Reclaimed Water for Data Center Cooling

The use of reclaimed or “gray” water is neither harmful to the environment nor to human health. Using gray water for cooling is considered eco-friendly because it reduces demands for ground water and does not require energy for the recycling process at waste water treatment sites. None of the data centers under this study are using reclaimed water for data center cooling.

6.2.3.3 Cooling Water Recirculation

Using the same water for several cycles of cooling operations reduces the water consumption. This improves data center energy efficiency and lowers environmental impact.

6.2.3.4 Using Renewable Resources

Renewable energy comes from solar panels, wind turbines or hydroelectric installations. As renewable energy production is intermittent in nature and depends on location, it is often combined with energy storage facilities. Nevertheless, these are still expensive installations. Two data centers under this study are using green energy.

Table 6.4 presents the green practices followed in the data centers under our study.

6.2.4 Storage and Network Practices

Generally, a data center is seen as a facility used to house computer systems and associated components, such as telecommunications and storage systems. Some of the best practices followed for storage and network communications of a data center are the following ones. Most of the data centers surveyed have centralized control over the servers, storage, and databases for storage optimization. Some of the data centers in our study are using pooling storage, hybrid storage and flash cache. Most of the data centers are using geo-replication for storage backup,

6.2 Data Center Management Best Practices

Table 6.4: Green practices for data centers

	DC1	DC2	DC3	DC4	DC5	DC6	DC7
Economizers (Free Cooling)	✓	X	X	✓	X	✓	✓
Reclaimed water for cooling	X	X	X	X	X	X	X
Cooling water recirculation	X	X	X	X	X	✓	✓
Using renewable resources	X	X	X	X	X	✓	✓

e-discovery and data mapping for archiving whereas flash storage is used only for specific applications in very few data centers of our study.

The data centers in our study use automation tools capable of predicting network loads to avoid outages and make continuous adjustments according to their requirements. All the data centers in our study have backbone connectivity to the global and region providers. Some of them are using network virtualization that allows each customer to have their own with different controller applications and balance the performance, port utilization and traffic demands. As Software Defined Networking (SDN) is the virtualization of networking and storage infrastructure, it offers resource flexibility, optimal resources usage, and scalability [284]. Some of our data centers are either adopted SDN or planning to adopt SDN in the near future. It is suggested that a data center should have a network strategy with the long-term view on network, servers, and storage to improve efficiency. Table 6.5 summarizes the storage and network practices followed in the data centers under our study.

6.2.5 Security Practices

Next generation data centers need to come up with built in security to meet high availability requirements of the business. Security must be planned to improve the risk management capabilities and to achieve compliance objectives. It is suggested to involve security team from design phase itself and they must develop controlled

6.2 Data Center Management Best Practices

Table 6.5: Storage and Network practices for data centers

	DC1	DC2	DC3	DC4	DC5	DC6	DC7
Storage tiering	✓	✓	✓	✓	✓	✓	✓
Automation to waste storage management	✓	✓	✓	✓	✓	✓	✓
Centralized control & Storage optimization	✓	✓	✓	✓	✓	✓	✓
Software-defined Flash storage	✓	✓	X	✓	X	✓	✓
Storage pooling and geo-replication	✓	✓	✓	✓	✓	✓	✓

access mechanism for each and every component of the data center which must agree on a standard policy environment. Data centers ought to be built utilizing a strong configuration that gives compelling protection against storms and natural disasters. Some data centers in our study established fire compartments and monitoring with the help of aspirating smoke detectors, to handle fire break out.

Most of the data centers in our study are using firewalls of different manufacturers to secure data. In case of component failure, stability is maintained in most of the data centers by re-directing the load to the redundant components. Most of the data centers in our study follow the security policies which are scalable and flexible with different kinds of applications and environments considering future as well as current needs. These data centers implement secure trust zones [285] in place of physical trust zones to make applications accessible from any device and at any time. All the studied data centers have single-person access and mantrap systems that avoid tailgating. Outside individuals are authenticated by means of biometric scans, configured ID cards, etc. They have continuous monitoring of physical and virtual assets either active or passive. Data centers should monitor for the changes which may introduce vulnerabilities that can be exploited using logs and configurations. Centralized management will provide the data center detailed reports across all controls necessary for risk management. In the case of multi-tenant environments, policies must be able to segment across different tenants.

6.3 Recommendations for Data Center Operators and IT Professionals

Table 6.6: Security practices for data centers

	DC1	DC2	DC3	DC4	DC5	DC6	DC7
Security zones	eight	six	six	six	six	Reda- -cted	Reda- -cted
Seismic zone	2	3	2	2	2	0	0
Redundancy	N+N	N+N	N+1	N+1	N+1	N	N
Single-person access and mantrap systems	✓	✓	✓	✓	✓	✓	✓
Monitoring inside and outside individuals	✓	✓	✓	✓	✓	✓	✓
Smoke detectors and alarming systems	✓	✓	✓	✓	✓	✓	✓

Table 6.6 summarizes the security practices followed in the data centers under our study.

6.3 Recommendations for Data Center Operators and IT Professionals

Data center life cycle management help owners to understand the key management tasks, connection between different phases and the pitfalls that exist in each phase. Generally, the data center life cycle comprises five phases: Plan & Analyze, Design, Build, Operate, and Continuous Evaluation. For initial phases, it is better to use reference designs to validate the early project choices and develop system concepts. Considering the whole chain of operations, data center operations become the base layer that has the goal of optimizing not only energy and cost, but also help the long-term planning and provisioning of equipment and resources. Table 6.7 summarizes various implementation issues in the said categories that lead to inefficiencies in data centers we studied.

6.3 Recommendations for Data Center Operators and IT Professionals

Table 6.7: Implementation issues and challenges in data centers

Area	Implementation Issues
Energy Efficiency	<ul style="list-style-type: none"> ● Some virtual machines are allocated more resources than requested to achieve high performance, leaving fewer resources available for other virtual machines. This may lead to inefficient use of resources. ● There are several power conversion steps while delivering power to the IT equipment. This leads to losses in power distribution. ● Most data centers do not monitor the detailed energy use. But detailed energy use will help the operators to know actual point of losses and inefficiencies. ● Power savings and performance requirements may lead to service level agreement violations.
Cooling, Thermal, and Air management	<ul style="list-style-type: none"> ● High density racks may result in one or more areas of excess temperature known as Hot spots which result in equipment damage. ● Poor air management in data centers may lead to bypass and recirculation air flow which reduces the cooling efficiency. ● Air blockages lead to poor flow of cool air to the server inlets. ● High heat loads from racks restrict the space utilization. ● Improper configurations lead to inefficient use of chillers. ● Oversized cooling infrastructure limit the operating capacity. ● Inadequate duct size results in poor air flow.
Greenness	<ul style="list-style-type: none"> ● Data centers consume a lot of water but they fail to recycle and use it in an efficient manner. ● Recycling of electronic equipment and other materials are not in place.
Network and Storage	<ul style="list-style-type: none"> ● Storage over-provisioning and massive volumes of redundant data lead to inefficient use of storage and consume more energy. ● Increasing virtual machines makes more snapshots that consume more storage. ● High throughput and performance requirements limit the efficient use of storage and network. ● Connectivity issues may lead to outages.
Security	<ul style="list-style-type: none"> ● Network attacks are able to compromise conventional security devices. ● Most traffic vulnerabilities are not visible in virtual data center which are known as blind spots that lead to security issues and chances of inter-virtual machine attacks. ● It is difficult to handle the mixed trust level virtual machines deployed in a server. ● Security and isolation should be maintained in shared multi-tenant environments of the public cloud.

6.3 Recommendations for Data Center Operators and IT Professionals

Nowadays, cloud computing is of strategic importance benefiting both providers and their customers. If a new data center is under-utilized, with the acceptance and popularity of cloud computing, it can act as a cloud provider for other data centers and customers. To accommodate the growing demands of users and other background processes using the same physical resources, data centers are required to make optimal use of all the resources by increasing utilization and visibility. Proper selection of virtual machines for migration minimizes the number of power-on nodes. Designing and implementing fast energy-efficient virtual machine allocation and selection algorithms considering multiple resources can build energy efficient data centers as discussed in Chapter 3. Maintaining a separate direct current feed to power the telecommunications and storage systems will reduce the energy consumption, real estate cost, and conversion losses. DCIM or automation tools are able to achieve considerable energy savings ranging from 5-20% [277]. Centralized cooling system in a large scale data center can be optimized by maintaining median temperature of 10°C to 15°C, using adjustable-speed drive chillers, storing excess thermal energy and installing energy and load monitoring sensors. Following these recommendations for chiller plants will have quick return of investments in the order of 2-3 years in terms of energy savings and diagnosis. Hot and cold aisles containment, increasing the data center supply air temperatures, using waterside economizer, and increasing the room temperature reduces the need for cooling provided by CRAH units. Further, using higher temperature-chilled water supply can provide sufficient cooling for a data center that reduces the number of hours of compressor-based cooling.

To improve the Power Usage Effectiveness (PUE) [286], it is important to understand normal power consumption and scrubble out unusual conditions, across the data centers. Further, it is important to develop compute efficiencies by server type, adjust operations according to peak power utilization and shift resource usage based on known use profiles. Operators need to correlate infrastructure investments more closely to actual application requirements. In short, optimizing cooling plant, operational parameters of data center, UPS load, and zombie servers along with controlled lighting, continuous monitoring, and proper air-flow management improve the energy-efficiency of a data center. Placing archived data on slower and larger drives that use less power saves some of the associated energy costs.

In our study, DC3 claims a reduction of 20% in their PUE score by following the best practices in cooling, thermal and air management and storage and network practices. DC1 claims a PUE reduction of 20% by strictly adhering to energy efficiency practices such as virtualization, consolidation, and automation tools. DC2 and DC4 claim a PUE reduction of 10% by effectively following the best practices of energy efficiency and cooling, thermal and air management.

Based on the metrics and best practices, we develop a checklist of recommended actions to increase energy efficiency of data centers, described in Appendix C. Designed for data center owners and operators, this checklist provides actionable guidance to prioritize and implement energy saving measures in data centers.

6.4 Summary

This chapter describes the best practices to improve operational efficiency and the ways to create a sustainable data center. We studied the issues in the current environment of the seven data centers in India and the Netherlands and described the best practices to deal with these issues. Implementing these practices in new and existing facilities will improve the efficiency of the data centers. These practices should be suitably adapted and fine-tuned for implementation in other data centers. Further, we provided a set of recommendations for data center operators and IT professionals to improve the efficiency of the data centers.

Chapter 7

Conclusions and Future Directions

The number of data centers has been increasing over the last decade to meet the rapidly exploding computation demand. This high demand, in turn, translates into elevated electricity bills, which have become one of the largest components of data center operational expenses. As a result, it has become extremely important to manage the data center resources in an energy efficient manner to reduce the operating costs and CO₂ emissions.

An energy efficient data center needs to optimize all the resources in the data center in terms of its energy consumption. This thesis addresses the problem of energy efficient resource management in cloud data centers. It focuses on the techniques for energy aware VM placement and selection, investigates prediction of energy demand for better capacity planning and improving resource utilization, and analyzes metrics and best practices for sustainable data centers. In this context, this thesis presents the following contributions:

In Chapter 3, we presented three VM placement algorithms and a selection algorithm to minimize energy consumption, number of VM migrations, and SLA violations. The first approach is Modified Discrete Particle Swarm Optimization (MDPSO) based on DPSO with new operators for updating velocity and particle positions. The second approach is based on Particle Swarm Optimization and Genetic Algorithm. We proposed a parallelized optimization algorithm called Interactive PSO-GA (IPSOGA), using multi-threading and shared memory for

information exchange to enhance convergence time and global exploration. The third approach is based on the imitating behavior of humans. We developed an Imitation Based Optimization (IBO) algorithm for virtual machine placement. This technique generates an optimal solution that tend to satisfy the structural information and provides consistency. Further, we presented a novel virtual machine selection method considering the factors such as memory, bandwidth and size of the VM (MBS-VM). This method optimally selects the virtual machines from a under/over utilized server and performs migration to further improve the energy efficiency in a data center. The experimental results of our proposed approaches revealed significantly improved performance as compared to the state-of-art methods.

Considering the estimated power requirements, keeping energy costs within budget and operating within the available capacities of power distribution are becoming important requirements for data centers. In Chapter 4, we developed two machine learning approaches capable of predicting the chiller energy consumption of a data center. We developed a Multi layer Feed Forward Neural Networks (MFNN) with the popular back-propagation training. We found that backpropagation may not always find the correct weights for the optimum solution. Thus, to improve the accuracy of the predictions, we developed a Deep learning approach with Parallel Stochastic Gradient Descent (DPSGD) training for forecasting the energy demand of chillers in a data center.

By clearly defining metrics, any business can increase its productivity. Whether the goal is to make a data center energy efficient, green, or more resilient, metrics and standards shape the action plan and provide insights into improvements. In Chapter 5, we presented an analysis of metrics that are commonly used in data centers, starting from the power grid and going all the way up to the service delivery. Furthermore, we derive relationships between metrics, and discuss the advantages and disadvantages of each metric in order to expose the research gaps. In Chapter 6, we described a set of best practices to improve the energy efficiency of the data centers which spans the categories of energy efficiency, cooling, thermal and air management, greenness, storage, and networks. We presented efficient alternatives in daily operations of the data centers and costs saving opportunities.

There are several open research challenges that need to be addressed to further advance the area of energy efficiency (sustainability) in data centers.

Virtual machines communicate with one another in a network of different topologies. If the allocation of resources is not done in an optimized way, then several migrations of processes will occur. To eliminate delays in data transfer and reduce power consumption, observing the communication pattern among CPUs is important. More efforts are needed to study the relationships among virtual machines and their communication patterns.

Another important future direction is to take advantage of modern sensing and communicating devices, referred to as Internet of Things (IoT), in order to acquire data from operations. All software aspects, such as scheduling, load balancing, and all the computations performed by the devices can be considered as cyber components. The supported infrastructure, such as servers and switches are can be modeled as physical components. In this context of cyberphysical systems, the collected data can be processed with pattern matching / machine learning approaches in order to deduce facts about the operations. These facts are useful to control the future operations with the precise goal of reducing resource consumption to the strictly necessary, thus identifying optimal processes of operation in terms of performance/resource consumption.

Fog computing is an extension to the cloud computing model that enables general purpose computing on traffic routing nodes so as to process data as it is transmitted between user devices and a data center. Researchers can identify the scenarios where running an application from nano servers are more energy-efficient than running the same applications from a data center. Another research direction is to reduce the delay in provisioning resources for resource limited fog nodes by designing scheduling algorithms considering priority and mobility model.

Researchers can model the problem of VM placement with various other soft computing approaches considering different objectives such as real-time scheduling, network congestion, access latency, turn around time etc. in addition to energy consumption. Further, it is required to develop VM placement and selection algorithms for private or hybrid cloud data centers considering the given objectives.

Data centers can benefit from machine/deep learning in order to have more optimized resource management. The potential of powerful Artificial Neural Network architectures such as extreme learning machine, spiking neural network, functional-link network, higher order neural networks, etc. is to be investigated for forecasting resource and energy demand of data centers.

References

- [1] C. Garnier, M. Aggar, M. Banks, J. Dietrich, B. Shatten, M. Stutz, and E. Tong-Viet, “White paper on data center life cycle assessment guidelines,” *The Green Grid*, 2012.
- [2] C. Wu and R. Buyya, *Cloud Data Centers and Cost Modeling: A Complete Guide To Planning, Designing and Building a Cloud Data Center*. Morgan Kaufmann, 2015.
- [3] R. Snevely, *Enterprise data center design and methodology*. Prentice Hall Press, 2002.
- [4] *Telecommunications Infrastructure Standard for Data Centers, Revision of TIA-942*. Technical Report, Telecommunication Industry Association, <http://www.tiaonline.org/standards/catalog/>, 2012.
- [5] “Data center site infrastructure tier standard: Topology,” *Technical Report, Uptime Institute Professional Services LLC*, 2012.
- [6] W. P. Turner IV, J. H. Seader, P. V. Renaud, and K. G. Brill, “Tier classification define site infrastructure performance,” *Uptime Institute*, vol. 17, 2006.
- [7] C. Belady, A. Rawson, J. Pflueger, and T. Cader, “The green grid data center power efficiency metrics: Pue and dcie,” *The Green Grid, White paper, 6*, 2007.
- [8] M. P. Mills, “The cloud begin with coal-an overview of the electricity used by the global digital ecosystem,” 2013.

REFERENCES

- [9] W. Van Heddeghem, S. Lambert, B. Lannoo, D. Colle, M. Pickavet, and P. De-meester, “Trends in worldwide ICT electricity consumption from 2007 to 2012,” *Computer Communications*, vol. 50, pp. 64–76, 2014.
- [10] J. Koomey, “Growth in data center electricity use 2005 to 2010,” *A report by Analytical Press, completed at the request of the New York Times*, pp. 9–10, 2011.
- [11] B. Pernici, M. Aiello, J. vom Brocke, B. Donnellan, E. Gelenbe, and M. Kretsis, “What is can do for environmental sustainability: A report from caise’11 panel on green and sustainable is,” *CAIS*, vol. 30, no. 18, 2012.
- [12] B. Pernici, C. Cappiello, M. G. Fugini, P. Plebani, M. Vitali, I. Salomie, T. Cioara, I. Anghel, E. Henis, and R. Kat, “Setting energy efficiency goals in data centers: the games approach,” in *Energy Efficient Data Centers*, pp. 1–12, Springer, 2012.
- [13] V. Dinesh Reddy, G. R. Gangadharan, and G. S. V. R. K. Rao, “Energy-aware virtual machine allocation and selection in cloud data centers,” *Soft Computing*, 2017.
- [14] S. K. Garg, A. N. Toosi, S. K. Gopalaiyengar, and R. Buyya, “Sla-based virtual machine management for heterogeneous workloads in a cloud datacenter,” *Journal of Network and Computer Applications*, vol. 45, pp. 108–120, 2014.
- [15] L. Ramachandran, N. C. Narendra, and K. Ponnalagu, “Dynamic provisioning in multi-tenant service clouds,” *Service Oriented Computing and Applications*, vol. 6, no. 4, pp. 283–302, 2012.
- [16] S. Ricciardi, D. Careglio, J. Sole-Pareta, U. Fiore, F. Palmieri, *et al.*, “Saving energy in data center infrastructures,” in *Proceedings of the First International Conference on Data Compression, Communications and Processing (CCP)*, pp. 265–270, IEEE, 2011.
- [17] R. Kruse, C. Borgelt, C. Braune, S. Mostaghim, and M. Steinbrecher, *Computational intelligence: a methodological introduction*. Springer, 2016.

REFERENCES

- [18] J. Fulcher, “Computational intelligence: an introduction,” in *Computational intelligence: a compendium*, pp. 3–78, Springer, 2008.
- [19] A. Konar, *Computational intelligence: principles, techniques and applications*. Springer Science & Business Media, 2006.
- [20] R. G. Michael and S. J. David, *Computers and intractability: a guide to the theory of NP-completeness*. New York: W. H. Freeman & Co, 1979.
- [21] F. Hao, M. Kodialam, T. Lakshman, and S. Mukherjee, “Online allocation of virtual machines in a distributed cloud,” *IEEE/ACM Transactions on Networking*, vol. 25, no. 1, pp. 238–249, 2017.
- [22] G. J. Woeginger, “Exact algorithms for np-hard problems: A survey,” in *Combinatorial Optimization-Eureka, You Shrink!*, pp. 185–207, Springer, 2003.
- [23] C. Jatoth, G. Gangadharan, and R. Buyya, “Computational intelligence based qos-aware web service composition: A systematic literature review,” *IEEE Transactions on Services Computing*, vol. 10, no. 3, pp. 475–492, 2017.
- [24] X.-S. Yang, “Harmony search as a metaheuristic algorithm,” in *Music-inspired harmony search algorithm*, pp. 1–14, Springer, 2009.
- [25] G. Portaluri, D. Adami, A. Gabbrielli, S. Giordano, and M. Pagano, “Power consumption-aware virtual machine placement in cloud data center,” *IEEE Transactions on Green Communications and Networking*, vol. 1, no. 4, pp. 541–550, 2017.
- [26] K. Zhang, T. Wu, S. Chen, L. Cai, and C. Peng, “A new energy efficient vm scheduling algorithm for cloud computing based on dynamic programming,” in *Proceedings of the 4th International Conference on Cyber Security and Cloud Computing*, pp. 249–254, 2017.
- [27] S. K. Garg, A. N. Toosi, S. K. Gopalaiyengar, and R. Buyya, “Sla-based virtual machine management for heterogeneous workloads in a cloud datacenter,” *Journal of Network and Computer Applications*, vol. 45, pp. 108–120, 2014.

REFERENCES

- [28] D. Shen, J. Luo, F. Dong, and J. Zhang, "Appbag: Application-aware bandwidth allocation for virtual machines in cloud environment," in *Proceedings of 45th International Conference on Parallel Processing (ICPP)*, pp. 21–30, 2016.
- [29] F.-H. Tseng, Y.-M. Jheng, L.-D. Chou, H.-C. Chao, and V. C. Leung, "Link-aware virtual machine placement for cloud services based on service-oriented architecture," *IEEE Transactions on Cloud Computing*, 2017.
- [30] H. Goudarzi and M. Pedram, "Energy-efficient virtual machine replication and placement in a cloud computing system," in *Proceedings of the 5th International Conference on Cloud Computing*, pp. 750–757, 2012.
- [31] S. K. Addya, A. K. Turuk, B. Sahoo, A. Satpathy, and M. Sarkar, "A game theoretic approach to estimate fair cost of vm placement in cloud data center," *IEEE Systems Journal*, 2017.
- [32] Z. Zhang, C.-C. Hsu, and M. Chang, "Cool cloud: A practical dynamic virtual machine placement framework for energy aware data centers," in *Proceedings of the IEEE 8th International Conference on Cloud Computing*, pp. 758–765, 2015.
- [33] N. Quang-Hung, N. Thoai, and N. T. Son, "Epobf: energy efficient allocation of virtual machines in high performance computing cloud," *Transactions on Large-Scale Data-and Knowledge-Centered Systems XVI*, vol. 8960, pp. 71–86, 2014.
- [34] X. Li, Z. Qian, S. Lu, and J. Wu, "Energy efficient virtual machine placement algorithm with balanced and improved resource utilization in a data center," *Mathematical and Computer Modelling*, vol. 58, no. 5, pp. 1222–1235, 2013.
- [35] X. Wang, X. Liu, L. Fan, and X. Jia, "A decentralized virtual machine migration approach of data centers for cloud computing," *Mathematical Problems in Engineering*, Hindawi Publishing Corporation, no. 878542, 2013.
- [36] Y. Mansouri, A. N. Toosi, and R. Buyya, "Cost optimization for dynamic replication and migration of data in cloud data centers," *IEEE Transactions on Cloud Computing*, 2017.

REFERENCES

- [37] Y. Xia, M. Tsugawa, J. A. Fortes, and S. Chen, “Large-scale vm placement with disk anti-colocation constraints using hierarchical decomposition and mixed integer programming,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 28, no. 5, pp. 1361–1374, 2017.
- [38] Y. Cui, Z. Yang, S. Xiao, X. Wang, and S. Yan, “Traffic-aware virtual machine migration in topology-adaptive dcn,” *IEEE/ACM Transactions on Networking*, vol. 25, no. 6, pp. 3427–3440, 2017.
- [39] B. P. Rimal and M. Maier, “Workflow scheduling in multi-tenant cloud computing environments,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 28, no. 1, pp. 290–304, 2017.
- [40] S. Imai, S. Patterson, and C. A. Varela, “Elastic virtual machine scheduling for continuous air traffic optimization,” in *Proceedings of 16th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, pp. 183–186, 2016.
- [41] M. Verma, G. R. Gangadharan, N. C. Narendra, R. Vadlamani, V. Inamdar, L. Ramachandran, R. N. Calheiros, and R. Buyya, “Dynamic resource demand prediction and allocation in multitenant service clouds,” *Concurrency and Computation: Practice and Experience*, vol. 28, no. 17, pp. 4429–4442, 2016.
- [42] Y. Wang and Y. Xia, “Energy optimal vm placement in the cloud,” in *Proceedings of the IEEE 9th International Conference on the Cloud Computing*, pp. 84–91, 2016.
- [43] Z. Zhou, Z. Hu, and K. Li, “Virtual machine placement algorithm for both energy-awareness and sla violation reduction in cloud data centers,” *Scientific Programming*, 2016.
- [44] P. Agrawal, D. Borgetto, C. Comito, G. Da Costa, J. M. Pierson, P. Prakash, S. Rao, D. Talia, C. Thiam, and P. Trunfio, “Scheduling and resource allocation,” in *Large-scale Distributed Systems and Energy Efficiency: A Holistic View* (P. JM, ed.), pp. 225–262, 2015.

REFERENCES

- [45] N. Quang-Hung, D. K. Le, N. Thoai, and N. T. Son, “Heuristics for energy-aware vm allocation in hpc clouds,” in *Proceedings of international conference on Future Data and Security Engineering*, pp. 248–261, 2014.
- [46] N. A. Singh and M. Hemalatha, “Energy efficient virtual machine placement technique using banker algorithm in cloud data centre,” in *Proceedings of the International Conference on Advanced Computing and Communication Systems (ICACCS)*, IEEE, 2013.
- [47] A. Beloglazov and R. Buyya, “Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in cloud data centers,” *Concurrency and Computation: Practice and Experience*, vol. 24, no. 13, pp. 1397–1420, 2012.
- [48] S. Bose, S. Brock, R. Skeoch, and S. Rao, “Cloudspider: Combining replication with scheduling for optimizing live migration of virtual machines across wide area networks,” in *Proceedings of the 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, pp. 13–22, 2011.
- [49] W. Li, J. Tordsson, and E. Elmroth, “Virtual machine placement for predictable and time constrained peak loads,” in *Proceedings of the International Workshop on Economics of grids, clouds, systems, and services*, pp. 120–34, 2011.
- [50] A. Beloglazov and R. Buyya, “Adaptive threshold-based approach for energy-efficient consolidation of virtual machines in cloud data centers.,” in *Proceedings of the 8th International Workshop on Middleware for Grids, Cloud and e-Science*, p. 4, 2010.
- [51] R. Buyya, A. Beloglazov, and J. Abawajy, “Energy-efficient management of data center resources for cloud computing: a vision, architectural elements, and open challenges,” in *Proceedings of the International Conference on Parallel and Distributed Processing Techniques and Applications*, pp. 6–17, CSREA Press, 2010.
- [52] X. Meng, V. Pappas, and L. Zhang, “Improving the scalability of data center networks with traffic-aware virtual machine placement,” in *Proceedings of the IEEE INFOCOM*, pp. 1–9, 2010.

REFERENCES

- [53] G. Dhiman, G. Marchetti, and T. Rosing, “vgreen:a system for energy efficient computing in virtualized environments,” in *Proceedings of the ACM/IEEE international symposium on Low power electronics and design*, pp. 243–248, 2009.
- [54] M. Cardoso, M. R. Korupolu, and S. A. S. and, “Shares and utilities based power consolidation in virtualized server environments,” in *Proceedings of the IFIP/IEEE International Symposium on Integrated Network Management*, pp. 327–334, 2009.
- [55] A. Verma, P. Ahuja, and A. Neogi, “pmapper: power and migration cost aware application placement in virtualized systems,” in *Proceedings of the 9th ACM/IFIP/USENIX International Conference on Middleware*, pp. 243–264, 2008.
- [56] R. Nathuji and K. Schwan, “Virtualpower: coordinated power management in virtualized enterprise systems,” *ACM SIGOPS Operating Systems Review*, vol. 41, no. 6, pp. 265–278, 2007.
- [57] S. Lee, R. Panigrahy, V. Prabhakaran, V. Ramasubramanian, K. Talwar, L. Uyeda, and U. Wieder, “Validating heuristics for virtual machines consolidation,” *Microsoft Research, MSR-TR-2011-9*, pp. 1–14, 2011.
- [58] Z. Xiao, W. Song, and Q. Chen, “Dynamic resource allocation using virtual machines for cloud computing environment,” *IEEE transactions on parallel and distributed systems*, vol. 24, no. 6, pp. 1107–1117, 2013.
- [59] Y. Wang and Y. Xia, “Energy optimal vm placement in the cloud,” in *Proceedings of the 9th International Conference on Cloud Computing*, pp. 84–91, IEEE, 2016.
- [60] L. Grange, G. Da Costa, and P. Stolf, “Green it scheduling for data center powered with renewable energy,” *Future Generation Computer Systems*, vol. 86, pp. 99–120, 2018.
- [61] A. K. Soomro, M. A. Shaikh, and H. Kazi, “Ffd variants for virtual machine placement in cloud computing data centers,” *International Journal of Advanced Computer Science and Applications*, vol. 8, no. 10, pp. 261–269, 2017.

REFERENCES

- [62] M. Ghobaei-Arani, M. Shamsi, and A. A. Rahmanian, “An efficient approach for improving virtual machine placement in cloud computing environment,” *Journal of Experimental & Theoretical Artificial Intelligence*, vol. 29, no. 6, pp. 1149–1171, 2017.
- [63] A. Zhou, S. Wang, B. Cheng, Z. Zheng, F. Yang, R. N. Chang, M. R. Lyu, and R. Buyya, “Cloud service reliability enhancement via virtual machine placement optimization,” *IEEE Transactions on Services Computing*, vol. 10, no. 6, pp. 902–913, 2017.
- [64] P. S. Pillai and S. Rao, “Resource allocation in cloud computing using the uncertainty principle of game theory,” *IEEE Systems Journal*, vol. 10, no. 2, pp. 637–648, 2016.
- [65] K. Su, L. Xu, C. Chen, W. Chen, and Z. Wang, “Affinity and conflict-aware placement of virtual machines in heterogeneous data centers,” in *Proceedings of the IEEE 12th International Symposium on Autonomous Decentralized Systems (ISADS)*, pp. 289–294, IEEE, 2015.
- [66] X. Zhang, K. Li, and Y. Zhang, “Minimum-cost virtual machine migration strategy in data center,” *Concurrency and Computation: Practice and Experience*, vol. 27, no. 17, pp. 5177–5187, 2015.
- [67] Q. Liang, J. Zhang, Y.-h. Zhang, and J.-m. Liang, “The placement method of resources and applications based on request prediction in cloud data center,” *Information Sciences*, vol. 279, pp. 735–45, 2014.
- [68] X. Dai, J. M. Wang, and B. Bensaou, “Energy-efficient virtual machine placement in data centers with heterogeneous requirements,” in *Proceedings of the IEEE 3rd International Conference on Cloud Networking (CloudNet)*, pp. 161–166, 2014.
- [69] W. Fang, X. Liang, S. Li, L. Chiaraviglio, and N. Xiong, “Vmplanner: Optimizing virtual machine placement and traffic flow routing to reduce network power costs in cloud data centers,” *Computer Networks*, vol. 57, no. 1, pp. 179–196, 2013.

REFERENCES

- [70] Z. Zhuang and G. C. Ocpa:, “An algorithm for fast and effective virtual machine placement and assignment in large scale cloud environments,” in *Proceedings of the International conference on Cloud Computing and Big Data (CloudCom-Asia)*, IEEE, 2013.
- [71] J. Dong, X. Jin, H. Wang, Y. Li, P. Zhang, and S. Cheng, “Energy-saving virtual machine placement in cloud data centers,” in *Proceedings of the IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, vol. 13, pp. 618–624, 2013.
- [72] S. Wang, P. P. Huang, C. Wen, and L.-C. Wang, “EQVMP: Energy-efficient and qos-aware virtual machine placement for software defined datacenter networks,” in *Proceedings of the International Conference on Information Networking (ICOIN)*, pp. 220–225, IEEE, 2014.
- [73] T. Wood, P. J. Shenoy, A. Venkataramani, and M. S. Yousif, “Black-box and gray-box strategies for virtual machine migration,” vol. 7, pp. 17–17, 2007.
- [74] D. Dong and J. Herbert, “Energy efficient vm placement supported by data analytic service,” in *Proceedings of the 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, pp. 648–655, 2013.
- [75] K. Y. Chen, Y. Xu, K. Xi, and H. J. Chao, “Intelligent virtual machine placement for cost efficiency in geo-distributed cloud systems,” in *Proceedings of the International Conference on Communications (ICC)*, pp. 3498–3503, IEEE, 2013.
- [76] G. Somani, P. Khandelwal, and K. Phatnani, “Vupic: Virtual machine usage based placement in iaas cloud,” *arXiv preprint, arXiv:1212.0085*, 2012.
- [77] A. Beloglazov, J. Abawajy, and R. Buyya, “Energy-aware resource allocation heuristics for efficient management of data centers for cloud computing,” *Future generation computer systems*, , 28(5), pp. 755–768, 2012.
- [78] H. Jin, D. Pan, J. Xu, and N. Pissinou, “Efficient vm placement with multiple deterministic and stochastic resources in data centers,” in *Proceedings of Global Communications Conference (GLOBECOM)*, IEEE, pp. 2505–2510, 2012.

REFERENCES

- [79] J. L. Lucas Simarro *et al.*, “Dynamic placement of virtual machines for cost optimization in multi-cloud environments,” in *Proceedings of IEEE International conference on High Performance Computing and Simulation (HPCS)*, 2011.
- [80] F. Machida, M. Kawato, and Y. Maeno, “Redundant virtual machine placement for fault-tolerant consolidated server clusters,” in *Proceedings of the Network Operations and Management Symposium*, pp. 32–39, 2010.
- [81] A. J. Younge, G. V. Laszewski, L. Wang, S. Lopez-Alarcon, and W. Carithers, “Efficient resource management for cloud computing environments,” in *Proceedings of International Green Computing Conference, IEEE*, pp. 357–364, 2010.
- [82] B. Li, J. Li, J. Huai, T. Wo, Q. Li, and Z. L. Enacloud:, “Enacloud: An energy-saving application live placement approach for cloud computing environments,” in *Proceedings of the IEEE International Conference on Cloud Computing*, pp. 17–24, 2009.
- [83] N. Bobroff, A. Kochut, and K. Beaty, “Dynamic placement of virtual machines for managing sla violations,” in *Proceedings of the 10th IFIP/IEEE International Symposium on Integrated Network Management*, pp. 119–128, 2007.
- [84] Z. Cao and S. Dong, “An energy-aware heuristic framework for virtual machine consolidation in cloud computing,” *The Journal of Supercomputing*, vol. 69, no. 1, pp. 429–451, 2014.
- [85] T. Wood, P. Shenoy, A. Venkataramani, and M. Yousif, “Sandpiper: Black-box and gray-box resource management for virtual machines,” *Computer Networks*, vol. 53, no. 17, pp. 2923–2938, 2009.
- [86] M. Mishra and A. Sahoo, “On theory of vm placement: Anomalies in existing methodologies and their mitigation using a novel vector based approach,” in *Proceedings of the IEEE International Conference on Cloud Computing*, pp. 275–282, 2011.
- [87] X. Li, Z. Qian, R. Chi, B. Zhang, and S. Lu, “Balancing resource utilization for continuous virtual machine requests in clouds,” in *Proceedings of the*

REFERENCES

- Sixth International Conference on Innovative Mobile and Internet Services in Ubiquitous Computing (IMIS)*, pp. 266–273, 2012.
- [88] X. Li, Z. Qian, S. Lu, and J. Wu, “Energy efficient virtual machine placement algorithm with balanced and improved resource utilization in a data center,” *Mathematical and Computer Modelling*, vol. 58, no. 5-6, pp. 1222–1235, 2013.
- [89] S. Fang, R. Kanagavelu, B.-S. Lee, C. H. Foh, and K. M. M. Aung, “Power-efficient virtual machine placement and migration in data centers,” in *Proceedings of the International Conference on Green Computing and Communications (GreenCom)*, pp. 1408–1413, IEEE, 2013.
- [90] K. Su, L. Xu, C. Chen, W. Chen, and Z. Wang, “Affinity and conflict-aware placement of virtual machines in heterogeneous data centers,” in *Proceedings of the 12th International Symposium on Autonomous Decentralized Systems (ISADS)*, pp. 289–294, IEEE, 2015.
- [91] H. Zhao, J. Wang, F. Liu, Q. Wang, W. Zhang, and Q. Zheng, “Power-aware and performance-guaranteed virtual machine placement in the cloud,” *IEEE Transactions on Parallel and Distributed Systems*, 2018.
- [92] H. Duan, C. Chen, G. Min, and Y. Wu, “Energy-aware scheduling of virtual machines in heterogeneous cloud computing systems,” *Future Generation Computer Systems*, vol. 74, pp. 142–150, 2017.
- [93] M. Tang and S. Pan, “A hybrid genetic algorithm for the energy-efficient virtual machine placement problem in data centers,” *Neural Processing Letters*, vol. 41, no. 2, pp. 211–221, 2014.
- [94] Q. Zheng, R. Li, X. Li, N. Shah, J. Zhang, F. Tian, K.-M. Chao, and J. Li, “Virtual machine consolidated placement based on multi-objective biogeography-based optimization,” *Future Generation Computer Systems*, vol. 54, pp. 95–122, 2016.
- [95] S. Wang, A. Zhou, C.-H. Hsu, X. Xiao, and F. Yang, “Provision of data-intensive services through energy-and qos-aware virtual machine placement in national cloud data centers,” *IEEE Transactions on Emerging Topics in Computing*, vol. 4, no. 2, pp. 290–300, 2016.

REFERENCES

- [96] X.-F. Liu, Z.-H. Zhan, J. D. Deng, Y. Li, T. Gu, and J. Zhang, “An energy efficient ant colony system for virtual machine placement in cloud computing,” *IEEE Transactions on Evolutionary Computation*, 2016.
- [97] S. Sawant, *A genetic algorithm scheduling approach for virtual machine resources in a cloud computing environment*. San Jose State University: Master’s Project, 2011.
- [98] Í. Goiri, J. L. Berral, J. O. Fitó, F. Julià, R. Nou, J. Guitart, R. Gavaldà, and J. Torres, “Energy-efficient and multifaceted resource management for profit-driven virtualized data centers,” *Future Generation Computer Systems*, vol. 28, no. 5, pp. 718–731, 2012.
- [99] W. Ding, C. Gu, F. Luo, Y. Chang, U. Rugwiro, X. Li, and G. Wen, “Dfa-vmp: An efficient and secure virtual machine placement strategy under cloud environment,” *Peer-to-Peer Networking and Applications*, vol. 11, no. 2, pp. 318–333, 2018.
- [100] A. Aryania, H. S. Aghdasi, and L. M. Khanli, “Energy-aware virtual machine consolidation algorithm based on ant colony system,” *Journal of Grid Computing*, pp. 1–15, 2018.
- [101] S. Kaur and A. Kaur, “The virtual machine migration in cloud computing using firefly and gravitational algorithm,” 2017.
- [102] Z. Chen, Y. Zhu, Y. Di, and S. Feng, “Optimization of virtual machine placement based on constrained immune memory and immunodominance clone in iaas cloud mode equipment training,” *International Journal of Modeling, Simulation, and Scientific Computing*, vol. 08, no. 01, p. 1750008, 2017.
- [103] G. Portaluri, D. Adami, A. Gabbrielli, S. Giordano, and M. Pagano, “Power consumption-aware virtual machine allocation in cloud data center,” in *Proceedings of the Globecom Workshops (GC Wkshps)*, pp. 1–6, IEEE, 2016.
- [104] S. Mu and Z. Li, “Deployment method of virtual machine based on pso algorithm modified by gauss strategy,” *Advanced Science and Technology Letters*, vol. 142, pp. 34–39, 2016.

REFERENCES

- [105] N. J. Kansal and I. Chana, “Energy-aware virtual machine migration for cloud computing—a firefly optimization approach,” *Journal of Grid Computing*, vol. 14, no. 2, pp. 327–345, 2016.
- [106] C. Gao, H. Wang, L. Zhai, Y. Gao, and S. Yi, “An energy-aware ant colony algorithm for network-aware virtual machine placement in cloud computing,” in *Proceedings of the IEEE 22nd International Conference on Parallel and Distributed Systems*, pp. 669–676, IEEE, 2016.
- [107] X. Wang, Y. Wang, and Y. Cui, “An energy-aware bi-level optimization model for multi-job scheduling problems under cloud computing,” *Soft Computing*, vol. 20, no. 1, pp. 303–317, 2016.
- [108] Q. Zheng *et al.*, “Virtual machine consolidated placement based on multi-objective biogeography-based optimization,” *Future Generation Computer Systems*, vol. 54, pp. 95–122, 2016.
- [109] S. E. Dashti and A. M. Rahmani, “Dynamic vms placement for energy efficiency by pso in cloud computing,” *Journal of Experimental & Theoretical Artificial Intelligence*, pp. 1–16, 2015.
- [110] D. Kumar and Z. Raza, “A pso based vm resource scheduling model for cloud computing,” in *Proceedings of IEEE International Conference on Computational Intelligence & Communication Technology (CICCT)*, pp. 213–219, IEEE, 2015.
- [111] S. Joshi and S. Kaur, “Cuckoo search approach for virtual machine consolidation in cloud data centre,” in *International Conference on Computing, Communication & Automation (ICCCA)*, pp. 683–686, IEEE, 2015.
- [112] M. Tang and S. Pan, “A hybrid genetic algorithm for the energy-efficient virtual machine placement problem in data centers,” *Neural Processing Letters*, vol. 41, no. 2, pp. 211–221, 2015.
- [113] J. p. Luo, X. Li, and M.-r. Chen, “Hybrid shuffled frog leaping algorithm for energy-efficient dynamic consolidation of virtual machines in cloud data centers,” *Expert Systems with Applications*, vol. 41, no. 13, pp. 5804–5816, 2014.

REFERENCES

- [114] B. Kruekaew and W. Kimpan, "Virtual machine scheduling management on cloud computing using artificial bee colony," in *Proceedings of the International MultiConference of engineers and computer scientists*, vol. 1, pp. 12–14, 2014.
- [115] D. J-k *et al.*, "Virtual machine placement optimizing to improve network performance in cloud data centers," *The Journal of China Universities of Posts and Telecommunications*, vol. 21, no. 3, pp. 62–70, 2014.
- [116] T. Yang, Y. C. Lee, and Z. Ay., "Energy-efficient data center networks planning with virtual machine placement and traffic configuration," in *Proceedings of the 6th international conference on Cloud Computing Technology and Science (CloudCom)*, pp. 284–291, IEEE, 2014.
- [117] J. Gao and G. Tang, "Virtual machine placement strategy research," in *Proceedings of the International Conference on Cyber-Enabled Distributed Computing and Knowledge Discovery (CyberC)*, pp. 294–297, IEEE, 2013.
- [118] S. Wang, Z. Liu, Z. Zheng, Q. Sun, and F. Yang, "Particle swarm optimization for energy-aware virtual machine placement optimization in virtualized data centers," in *Proceedings of the International Conference on Parallel and Distributed Systems (ICPADS)*, pp. 102–109, 2013.
- [119] F. Ma, F. Liu, and Z. Liu, "Multi-objective optimization for initial virtual machine placement in cloud data center," *Journal of Information & Computational Science*, vol. 9, no. 16, pp. 5029–5038, 2012.
- [120] Y. Wu, M. Tang, and W. A. Fraser, "simulated annealing algorithm for energy efficient virtual machine placement," in *Proceedings of the IEEE international conference on Systems Man and Cybernetics (SMC)*, pp. 1245–1250, IEEE, 2012.
- [121] H. Goudarzi and M. Pedram, "Energy-efficient virtual machine replication and placement in a cloud computing system," in *Proceedings of the 5th International conference on Cloud Computing (CLOUD)*, pp. 750–757, IEEE, 2012.

REFERENCES

- [122] G. Wu, M. Tang, Y. Tian, and W. Li, “Energy-efficient virtual machine placement in data centers by genetic algorithm,” in *Proceedings of International Conference on Neural Information Processing*, pp. 315–323, 2012.
- [123] R. Jeyarani, N. Nagaveni, and R. V. Ram, “Self adaptive particle swarm optimization for efficient virtual machine provisioning in cloud,” *International Journal of Intelligent Information Technologies*, vol. 7, no. 2, pp. 25–44, 2011.
- [124] M. Cct, D. Niyato, and T. Chen-Khong, “Evolutionary optimal virtual machine placement and demand forecaster for cloud computing,” in *Proceedings of the International conference on Advanced Information Networking and Applications*, pp. 348–355, IEEE, 2011.
- [125] J. Xu and F. Ja., “Multi-objective virtual machine placement in virtualized data center environments,” in *Proceedings of the IEEE/ACM international conference on Green Computing and Communications (GreenCom)*, pp. 179–188, IEEE, 2010.
- [126] J. T. Piao and J. A. Yan in *Proceedings of the 9th International Conference on Grid and Cooperative Computing (GCC)*, pp. 87–92, IEEE.
- [127] F. Palmieri and D. Castagna, “Swarm-based distributed job scheduling in next-generation grids,” in *Advances and Innovations in Systems* (C. Sciences and S. Engineering, eds.), pp. 137–143, Springer, 2007.
- [128] M. H. Ferdous, M. Murshed, R. N. Calheiros, and R. Buyya, “Virtual machine consolidation in cloud data centers using aco metaheuristic,” in *Proceedings of the European Conference on Parallel Processing*, pp. 306–317, 2014.
- [129] F. Palmieri, U. Fiore, S. Ricciardi, and A. Castiglione, “Grasp-based resource re-optimization for effective big data access in federated clouds,” *Future Generation Computer Systems*, vol. 54, pp. 168–179, 2016.
- [130] J. Levine and F. Ducatelle, “Ant colony optimization and local search for bin packing and cutting stock problems,” *Journal of the Operational Research Society*, vol. 55, no. 7, pp. 705–716, 2004.

REFERENCES

- [131] B. Brugger, K. F. Doerner, R. F. Hartl, and M. Reimann, “Antpacking—an ant colony optimization approach for the one-dimensional bin packing problem,” in *Proceedings of the European Conference on Evolutionary Computation in Combinatorial Optimization*, pp. 41–50, 2004.
- [132] E. Feller, L. Rilling, and C. Morin, “Energy-aware ant colony based workload placement in clouds,” in *Proceedings of the IEEE/ACM 12th International Conference on Grid Computing*, pp. 26–33, 2011.
- [133] Y. Gao, H. Guan, Z. Qi, Y. Hou, and L. Liu, “A multi-objective ant colony system algorithm for virtual machine placement in cloud computing,” *Journal of Computer and System Sciences*, vol. 79, no. 8, pp. 1230–1242, 2013.
- [134] S. Petersen and S. Svendsen, “Method and simulation program informed decisions in the early stages of building design,” *Energy and buildings*, vol. 42, no. 7, pp. 1113–1119, 2010.
- [135] R. H. Katz, D. E. Culler, S. Sanders, S. Alspaugh, Y. Chen, S. Dawson-Haggerty, P. Dutta, M. He, X. Jiang, and L. Keys, “An information-centric energy infrastructure: The berkeley view,” *Sustainable Computing: Informatics and Systems*, vol. 1, no. 1, pp. 7–22, 2011.
- [136] Y. Bengio *et al.*, “Learning deep architectures for ai,” *Foundations and trends® in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [137] A. Krioukov, P. Mohan, S. Alspaugh, L. Keys, D. Culler, and R. Katz, “Napsac: Design and implementation of a power-proportional web cluster,” *ACM SIGCOMM Computer Communication Review*, vol. 41, no. 1, pp. 102–108, 2011.
- [138] T. F. Abdelzaher, K. G. Shin, and N. Bhatti, “Performance guarantees for web server end-systems: A control-theoretical approach,” *IEEE transactions on parallel and distributed systems*, vol. 13, no. 1, pp. 80–96, 2002.
- [139] L. Wang, X. Wang, M. Tornatore, K. J. Kim, S. M. Kim, D.-U. Kim, K.-E. Han, and B. Mukherjee, “Scheduling with machine-learning-based flow detection for packet-switched optical data center networks,” *Journal of Optical Communications and Networking*, vol. 10, no. 4, pp. 365–375, 2018.

REFERENCES

- [140] A. Yu, H. Yang, W. Bai, L. He, H. Xiao, and J. Zhang, “Leveraging deep learning to achieve efficient resource allocation with traffic evaluation in datacenter optical networks,” in *Proceedings of the Optical Fiber Communication Conference*, 2018.
- [141] H. Shoukourian, T. Wilde, D. Labrenz, and A. Bode, “Using machine learning for data center cooling infrastructure efficiency prediction,” in *Proceedings of the Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pp. 954–963, 2017.
- [142] A. Nishi, S. Chuan, S. Yanbing, S. Xiaogang, D. Abishai, K. Rahul, Z. Tianyu, Z. Xiang, and Z. Lifei, “Power variation trend prediction in modern datacenter,” in *Proceedings of the 16th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, pp. 977–980, IEEE, 2017.
- [143] C. Liu, J. Han, Y. Shang, C. Liu, B. Cheng, and J. Chen, “Predicting of job failure in compute cloud based on online extreme learning machine: A comparative study,” *IEEE Access*, vol. 5, pp. 9359–9368, 2017.
- [144] I. K. Kim, S. Zeng, C. Young, J. Hwang, and M. Humphrey, “icsi: A cloud garbage vm collector for addressing inactive vms with machine learning,” in *Proceedings of IEEE International Conference on Cloud Engineering (IC2E)*, pp. 17–28, 2017.
- [145] J. Ahmed, A. Johnsson, F. Moradi, R. Pasquini, C. Flinta, and R. Stadler, “Online approach to performance fault localization for cloud and datacenter services,” in *Proceedings of IFIP/IEEE Symposium on Integrated Network and Service Management (IM)*, pp. 873–874, IEEE, 2017.
- [146] A. Bashar, “Bn-based approach for predictive admission control of cloud services,” in *Proceedings of the 7th International Advance Computing Conference (IACC)*, pp. 59–64, IEEE, 2017.
- [147] A.-y. Son and E.-N. Huh, “Study on a migration scheme by fuzzy-logic-based learning and decision approach for qos in cloud computing,” in *Proceedings of Ninth International Conference on Ubiquitous and Future Networks (ICUFN)*, pp. 507–512, 2017.

REFERENCES

- [148] A. Jobava, A. Yazidi, B. J. Oommen, and K. Begnum, “On achieving intelligent traffic-aware consolidation of virtual machines in a data center using learning automata,” *Journal of Computational Science*, vol. 24, pp. 290–312, 2017.
- [149] R. S. Sidhu, “Machine learning based datacenter monitoring framework,” *Masters Thesis, University of Texas*.
- [150] V. Moskalenko and S. Pimonenko, “Optimizing the parameters of functioning of the system of management of data center it infrastructure,” *Eastern-European Journal of Enterprise Technologies*, vol. 5, no. 2, pp. 21–29, 2016.
- [151] Y. Ukidave, X. Li, and D. Kaeli, “Mystic: Predictive scheduling for gpu based cloud servers using machine learning,” in *Proceedings of the IEEE International Parallel and Distributed Processing Symposium*, pp. 353–362, IEEE, 2016.
- [152] Z. Shen, C. C. Young, S. Zeng, K. Murthy, and K. Bai, “Identifying resources for cloud garbage collection,” in *Proceedings of the 12th International Conference on Network and Service Management*, pp. 248–252, 2016.
- [153] C. Sieber, A. Basta, A. Blenk, and W. Kellerer, “Online resource mapping for sdn network hypervisors using machine learning,” in *Proceedings of the NetSoft Conference and Workshops (NetSoft)*, pp. 78–82, IEEE, 2016.
- [154] Y. Tarutani, K. Hashimoto, G. Hasegawa, Y. Nakamura, T. Tamura, K. Matsudax, and M. Matsuoka, “Reducing power consumption in data center by predicting temperature distribution and air conditioner efficiency with machine learning,” in *Proceedings of the IEEE International Conference on Cloud Engineering*, pp. 226–227, 2016.
- [155] N. T. Hieu, M. Di Francesco, and A. Ylä-Jääski, “Virtual machine consolidation with usage prediction for energy-efficient cloud data centers,” in *Proceedings of the IEEE 8th International Conference on Cloud Computing (CLOUD)*, pp. 750–757, 2015.
- [156] S. S. Masoumzadeh and H. Hlavacs, “A cooperative multi agent learning approach to manage physical host nodes for dynamic consolidation of virtual

REFERENCES

- machines,” in *Proceedings of the IEEE Fourth Symposium on Network Cloud Computing and Applications (NCCA)*, pp. 43–50, 2015.
- [157] M. Dabbagh, B. Hamdaoui, M. Guizani, and A. Rayes, “Energy-efficient resource allocation and provisioning framework for cloud data centers,” *IEEE Transactions on Network and Service Management*, vol. 12, no. 3, pp. 377–391, 2015.
- [158] M. Fallah, M. G. Arani, and M. Maeen, “Nasla: Novel auto scaling approach based on learning automata for web application in cloud computing environment,” *International Journal of Computer Applications*, vol. 113, no. 2, 2015.
- [159] D. Versick, I. Waßmann, and D. Tavangarian, “Power consumption estimation of cpu and peripheral components in virtual machines,” *ACM SIGAPP Applied Computing Review*, vol. 13, no. 3, pp. 17–25, 2013.
- [160] J. L. Berral, R. Gavaldà, and J. Torres, “Power-aware multi-data center management using machine learning,” in *Proceedings of the 42nd International Conference on Parallel Processing (ICPP)*, pp. 858–867, IEEE, 2013.
- [161] K. Sato, M. Samejima, and N. Komoda, “Dynamic optimization of virtual machine placement by resource usage prediction,” in *Proceedings of the 11th IEEE International Conference on Industrial Informatics (INDIN)*, pp. 86–91, IEEE, 2013.
- [162] N. Vasić, D. Novaković, S. Miućin, D. Kostić, and R. Bianchini, “Dejavu: accelerating resource allocation in virtualized environments,” *ACM SIGARCH computer architecture news*, vol. 40, no. 1, pp. 423–436, 2012.
- [163] J. Yuan, X. Jiang, L. Zhong, and H. Yu, “Energy aware resource scheduling algorithm for data center using reinforcement learning,” in *Proceedings of the fifth International Conference on Intelligent Computation Technology and Automation (ICICTA)*, pp. 435–438, IEEE, 2012.
- [164] L. Li, C. J. M. Liang, J. Liu, S. Nath, A. Terzis, and C. Faloutsos, “Thermocast: a cyber-physical forecasting model for datacenters,” in *Proceedings of the 17th ACM international conference on Knowledge discovery and data mining*, pp. 1370–1378, ACM, 2011.

REFERENCES

- [165] Q. Zhang, L. Cherkasova, and E. Smirni, “A regression-based analytic model for dynamic resource provisioning of multi-tier applications,” in *Proceedings of the Fourth International Conference on Autonomic Computing, 2007. ICAC’07*, pp. 27–27, IEEE, 2007.
- [166] J. Moore, J. S. Chase, and P. Ranganathan, “Weatherman: Automated, on-line and predictive thermal mapping and management for data centers,” in *Proceedings of the International Conference on Autonomic Computing*, pp. 155–164, IEEE, 2006.
- [167] M. Arif, A. K. Kiani, and J. Qadir, “Machine learning based optimized live virtual machine migration over wan links,” *Telecommunication Systems*, vol. 64, no. 2, pp. 245–257, 2017.
- [168] S. Zhang, W. Meng, J. Bu, S. Yang, Y. Liu, D. Pei, J. Xu, Y. Chen, H. Dong, and X. Qu, “Syslog processing for switch failure diagnosis and prediction in datacenter networks,” in *Proceedings of the IEEE/ACM 25th International Symposium on Quality of Service (IWQoS)*, pp. 1–10, IEEE, 2017.
- [169] S. A. Yousif and A. Al-Dulaimy, “Clustering cloud workload traces to improve the performance of cloud data centers,” in *Proceedings of the World Congress on Engineering*, vol. 1, 2017.
- [170] Y. Li, Y. Wen, K. Guan, and D. Tao, “Transforming cooling optimization for green data center via deep reinforcement learning,” *arXiv preprint, arXiv:1709.05077*, 2017.
- [171] M. Nakamura, “Learning and optimization models for energy efficient cooling control in data center,” in *Proceedings of the 55th Annual Conference of the Society of Instrument and Control Engineers of Japan (SICE)*, pp. 395–400, IEEE, 2016.
- [172] R. Shaw, “An artificial intelligence model for autonomous resource allocation in cloud computing environments,” *Masters Thesis, National University of Ireland, Galway*, 2016.

REFERENCES

- [173] M. Dabbagh, B. Hamdaoui, M. Guizani, and A. Rayes, "Toward energy-efficient cloud computing: Prediction, consolidation, and overcommitment," *IEEE network*, vol. 29, no. 2, pp. 56–61, 2015.
- [174] F.-H. Tseng, X. Chen, L.-D. Chou, H.-C. Chao, and S. Chen, "Support vector machine approach for virtual machine migration in cloud data center," *Multimedia Tools and Applications*, vol. 74, no. 10, pp. 3419–3440, 2015.
- [175] Z. Tang, Y. Mo, K. Li, and K. Li, "Dynamic forecast scheduling algorithm for virtual machine placement in cloud computing environment," *The Journal of Supercomputing*, vol. 70, no. 3, pp. 1279–1296, 2014.
- [176] H. Shoukourian, T. Wilde, A. Auweter, and A. Bode, "Predicting the energy and power consumption of strong and weak scaling hpc applications," *Supercomputing frontiers and innovations*, vol. 1, no. 2, pp. 20–41, 2014.
- [177] F. Farahnakian, P. Liljeberg, and J. Plosila, "Energy-efficient virtual machines consolidation in cloud data centers using reinforcement learning," in *Proceedings of the 22nd Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP)*, pp. 500–507, IEEE, 2014.
- [178] R. Mijumbi, J.-L. Gorricho, J. Serrat, M. Claeys, F. De Turck, and S. Latré, "Design and evaluation of learning algorithms for dynamic resource management in virtual networks," in *Proceedings of the Network Operations and Management Symposium (NOMS)*, pp. 1–9, IEEE, 2014.
- [179] Y. Liu, B. Gong, C. Xing, and Y. Jian, "A virtual machine migration strategy based on time series workload prediction using cloud model," *Mathematical Problems in Engineering*, vol. 2014, 2014.
- [180] D. Dong and J. Herbert, "Energy efficient vm placement supported by data analytic service," in *Proceedings of the 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, pp. 648–655, IEEE, 2013.
- [181] F. Ramezani, J. Lu, and F. Hussain, "An online fuzzy decision support system for resource management in cloud environments," in *Proceedings of the*

REFERENCES

- Joint IFSA World Congress and NAFIPS Annual Meeting (IFSA/NAFIPS)*, pp. 754–759, IEEE, 2013.
- [182] C. Canali and R. Lancellotti, “Automatic virtual machine clustering based on bhattacharyya distance for multi-cloud systems,” in *Proceedings of the international workshop on Multi-cloud applications and federated clouds*, pp. 45–52, ACM, 2013.
- [183] X. Zhao, J. Yin, Z. Chen, and S. He, “Workload classification model for specializing virtual machine operating system,” in *Proceedings of the IEEE Sixth International Conference on Cloud Computing (CLOUD)*, pp. 343–350, IEEE, 2013.
- [184] F. Farahnakian, P. Liljeberg, and J. Plosila, “Lircup: Linear regression based cpu usage prediction algorithm for live migration of virtual machines in data centers,” in *Proceedings of the 39th EUROMICRO Conference on Software Engineering and Advanced Applications (SEAA)*, pp. 357–364, 2013.
- [185] C. Z. Xu, J. Rao, and X. Bu, “Url: A unified reinforcement learning approach for autonomic cloud management,” *Journal of Parallel and Distributed Computing*, vol. 72, no. 2, pp. 95–105, 2012.
- [186] W. Fang, Z. Lu, J. Wu, and Z. Cao, “Rpps: A novel resource prediction and provisioning scheme in cloud data center,” in *Proceedings fo the Ninth International Conference on Services Computing (SCC)*, pp. 609–616, IEEE, 2012.
- [187] I. S. Moreno and J. Xu, “Neural network-based overallocation for improved energy-efficiency in real-time cloud environments,” in *Proceedings of the 15th International Symposium on Object/Component/Service-Oriented Real-Time Distributed Computing (ISORC)*, pp. 119–126, IEEE, 2012.
- [188] S. Kundu, R. Rangaswami, A. Gulati, M. Zhao, and K. Dutta, “Modeling virtualized applications using machine learning techniques,” *ACM Sigplan Notices*, vol. 47, no. 7, pp. 3–14, 2012.
- [189] L. Wang, G. von Laszewski, F. Huang, J. Dayal, T. Frulani, and G. Fox, “Task scheduling with ann-based temperature prediction in a data center: a

REFERENCES

- simulation-based study,” *Engineering with Computers*, vol. 27, no. 4, pp. 381–391, 2011.
- [190] N. Roy, A. Dubey, and A. Gokhale, “Efficient autoscaling in the cloud using predictive models for workload forecasting,” in *Proceedings of the IEEE International Conference on Cloud Computing (CLOUD)*, pp. 500–507, IEEE, 2011.
- [191] G. Kousiouris, T. Cucinotta, and T. Varvarigou, “The effects of scheduling, workload type and consolidation scenarios on virtual machine performance and their prediction through optimized artificial neural networks,” *Journal of Systems and Software*, vol. 84, no. 8, pp. 1270–1291, 2011.
- [192] J. L. Berral, R. Gavaldà, and J. Torres, “Adaptive scheduling on power-aware managed data-centers using machine learning,” in *Proceedings of the 12th IEEE/ACM International Conference on Grid Computing (GRID)*, pp. 66–73, IEEE, 2011.
- [193] B. Hu, Z. Lei, Y. Lei, D. Xu, and J. Li, “A time-series based precopy approach for live migration of virtual machines,” in *Proceedings of the 17th International Conference on Parallel and Distributed Systems (ICPADS)*, pp. 947–952, IEEE, 2011.
- [194] S. Kraft, G. Casale, D. Krishnamurthy, D. Greer, and P. Kilpatrick, “Io performance prediction in consolidated virtualized environments,” in *Proceedings of the 2nd ACM/SPEC International Conference on Performance engineering*, vol. 36, pp. 295–306, ACM, 2011.
- [195] G. Dhiman, K. Mihic, and T. Rosing, “A system for online power prediction in virtualized environments using gaussian mixture models,” in *Proceedings of the 47th ACM/IEEE Design Automation Conference (DAC)*, pp. 807–812, IEEE, 2010.
- [196] S. Akoush, R. Sohan, A. Rice, A. W. Moore, and A. Hopper, “Predicting the performance of virtual machine migration,” in *Proceedings of the IEEE International Symposium on Modeling, Analysis & Simulation of Computer and Telecommunication Systems (MASCOTS)*, pp. 37–46, IEEE, 2010.

REFERENCES

- [197] J. L. Berral, Í. Goiri, R. Nou, F. Julià, J. Guitart, R. Gavaldà, and J. Torres, “Towards energy-aware scheduling in data centers using machine learning,” in *Proceedings of the 1st International Conference on energy-Efficient Computing and Networking*, pp. 215–224, ACM, 2010.
- [198] J. Rao, X. Bu, C. Z. Xu, L. Wang, and G. Yin, “Vconf: a reinforcement learning approach to virtual machines auto-configuration,” in *Proceedings of the 6th international conference on Autonomic computing*, pp. 137–146, ACM, 2009.
- [199] H. W. Choi, H. Kwak, A. Sohn, and K. Chung, “Autonomous learning for efficient resource utilization of dynamic vm migration,” in *Proceedings of the 22nd annual international conference on Supercomputing*, pp. 185–194, ACM, 2008.
- [200] Q. Tang, T. Mukherjee, S. K. Gupta, and P. Cayton, “Sensor-based fast thermal evaluation model for energy efficient high-performance datacenters,” in *Proceedings of the Fourth International Conference on Intelligent Sensing and Information Processing*, pp. 203–208, IEEE, 2006.
- [201] N. Ganesh and A. Kannan, “Capacity predictor with varying pricing scheme and interoperability decision in a bursting situation for cloud computing environemnt,” <https://pdfs.semanticscholar.org/19de/833e2a898ff5b49de1c5e7631d00fcd3962f.pdf>, 2006.
- [202] G. Kousiouris, T. Cucinotta, and T. Varvarigou, “The effects of scheduling, workload type and consolidation scenarios on virtual machine performance and their prediction through optimized artificial neural networks,” *Journal of Systems and Software*, vol. 84, no. 8, pp. 1270–1291, 2011.
- [203] A. A. Bankole and S. A. Ajila, “Cloud client prediction models for cloud resource provisioning in a multitier web application environment,” in *Proceedings of IEEE 7th International Symposium on Service Oriented System Engineering (SOSE)*, pp. 156–161, 2013.
- [204] A. V. Do, J. Chen, C. Wang, Y. C. Lee, A. Y. Zomaya, and B. B. Zhou, “Profiling applications for virtual machine placement in clouds,” in *Proceedings*

REFERENCES

- of the IEEE International Conference on Cloud Computing (CLOUD)*, pp. 660–667, 2011.
- [205] Z. Gong, X. Gu, and J. Wilkes, “Press: Predictive elastic resource scaling for cloud systems,” in *Proceedings of the International Conference on Network and Service Management (CNSM)*, pp. 9–16, 2010.
- [206] T. Wood, L. Cherkasova, K. Ozonat, and P. Shenoy, “Profiling and modeling resource usage of virtualized applications,” in *Proceedings of the 9th ACM/IFIP/USENIX International Conference on Middleware*, pp. 366–387, 2008.
- [207] G. Du, H. He, and F. Meng, “Performance modeling based on artificial neural network in virtualized environments,” *Sensors & Transducers*, vol. 153, no. 6, p. 37, 2013.
- [208] N. Elprince, “Autonomous resource provision in virtual data centers,” in *Proceedings of the IFIP/IEEE International Symposium on Integrated Network Management*, pp. 1365–1371, 2013.
- [209] E. Keller, S. Ghorbani, M. Caesar, and J. Rexford, “Live migration of an entire network (and its hosts),” in *Proceedings of the 11th ACM Workshop on Hot Topics in Networks*, pp. 109–114, 2012.
- [210] G. Dasgupta, A. Sharma, A. Verma, A. Neogi, and R. Kothari, “Workload management for power efficiency in virtualized data centers,” *Communications of the ACM*, vol. 54, no. 7, pp. 131–141, 2011.
- [211] L. D. Gray, A. Kumar, and H. H. Li, “Workload characterization of the specpower.ssj2008 benchmark,” in *Proceedings of the fSPEC International Performance Evaluation Workshop*, pp. 262–282, Springer, 2008.
- [212] J. Kennedy and R. C. Eberhart, “A discrete binary version of the particle swarm algorithm,” in *IEEE International Conference on Systems, Man, and Cybernetics*, vol. 5, pp. 4104–4108, IEEE, 1997.
- [213] C. L. Chen, S. Y. Huang, Y. R. Tzeng, and C. L. Chen, “A revised discrete particle swarm optimization algorithm for permutation flow-shop scheduling problem,” *Soft Computing*, vol. 18, no. 11, pp. 2271–2282, 2014.

REFERENCES

- [214] A. Gandelli, F. Grimaccia, M. Mussetta, P. Pirinoli, and R. E. Zich, “Development and validation of different hybridization strategies between ga and pso,” in *Proceedings of the IEEE Congress on Evolutionary Computation*, pp. 2782–2787, IEEE, 2007.
- [215] Y.-T. Kao and E. Zahara, “A hybrid genetic algorithm and particle swarm optimization for multimodal functions,” *Applied Soft Computing*, vol. 8, no. 2, pp. 849–857, 2008.
- [216] A. A. A. Esmín, G. Lambert-Torres, and G. B. Alvarenga, “Hybrid evolutionary algorithm based on pso and ga mutation,” in *Proceedings of the Sixth International Conference on Hybrid Intelligent Systems*, pp. 57–57, IEEE, 2006.
- [217] X. Shi, Y. Liang, H. Lee, C. Lu, and L. Wang, “An improved ga and a novel pso-ga-based hybrid algorithm,” *Information Processing Letters*, vol. 93, no. 5, pp. 255–261, 2005.
- [218] A. Gálvez and A. Iglesias, “A new iterative mutually coupled hybrid ga-pso approach for curve fitting in manufacturing,” *Applied Soft Computing*, vol. 13, no. 3, pp. 1491–1504, 2013.
- [219] X. Shi, Y. Lu, C. Zhou, H. Lee, W. Lin, and Y. Liang, “Hybrid evolutionary algorithms based on pso and ga,” in *Proceedings of the Congress on Evolutionary Computation, CEC’03.*, vol. 4, pp. 2393–2399, IEEE, 2003.
- [220] A. Carlisle and G. Dozier, “An off-the-shelf pso,” in *Proceedings of the workshop on Particle Swarm Optimization*, 2001.
- [221] J. F. Schutte, “Particle swarms in sizing and global optimization,” *Masters Thesis, University of Pretoria, Department of Mechanical and Aeronautical Engineering*, 2001.
- [222] P. C. Mahalanobis, “On the generalized distance in statistics,” in *Proceedings of the statistical laboratory, National Institute of Science of India*, 1936.
- [223] G. J. McLachlan, “Mahalanobis distance,” *Resonance*, vol. 4, no. 6, pp. 20–26, 1999.

REFERENCES

- [224] R. De Maesschalck, D. Jouan-Rimbaud, and D. L. Massart, “The mahalanobis distance,” *Chemometrics and intelligent laboratory systems*, vol. 50, no. 1, pp. 1–18, 2000.
- [225] P. Svärd, B. Hudzia, S. Walsh, J. Tordsson, and E. Elmroth, “Principles and performance characteristics of algorithms for live vm migration,” *ACM SIGOPS Operating Systems Review*, vol. 49, no. 1, pp. 142–155, 2015.
- [226] R. N. Calheiros, R. Ranjan, A. Beloglazov, C. A. De Rose, and R. Buyya, “Cloudsim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms,” *Software: Practice and Experience*, vol. 41, no. 1, pp. 23–50, 2011.
- [227] K. E. Parsopoulos and M. N. Vrahatis, “Particle swarm optimization and intelligence: Advances and applications,” *IGI Global*, 2010.
- [228] F. Alecu, “Software programs, from sequential to parallel,” *Oeconomics of Knowledge*, vol. 1, no. 2, p. 17, 2009.
- [229] G. M. Amdahl, “Validity of the single processor approach to achieving large scale computing capabilities,” in *Proceedings of spring joint computer conference*, pp. 483–485, 1967.
- [230] A. Ghosh and K. Pal, “Self-organization for object extraction using a multilayer neural network and fuzziness measures,” *IEEE TRANSACTIONS ON FUZZY SYSTEMS*, vol. 1, no. 1, 1993.
- [231] G. Carrera and J. Aires-de Sousa, “Estimation of melting points of pyridinium bromide ionic liquids with decision trees and neural networks,” *Green Chemistry*, vol. 7, no. 1, pp. 20–27, 2005.
- [232] J. Schmidhuber, “Deep learning in neural networks: An overview,” *Neural networks*, vol. 61, pp. 85–117, 2015.
- [233] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

REFERENCES

- [234] M. Riedmiller and H. Braun, “A direct adaptive method for faster backpropagation learning: The rprop algorithm,” in *Proceedings of IEEE International Conference on Neural Networks*, pp. 586–591, 1993.
- [235] A. Shahsavand, F. D. Fard, and F. Sotoudeh, “Application of artificial neural networks for simulation of experimental co₂ absorption data in a packed column,” *Journal of Natural Gas Science and Engineering*, vol. 3, no. 3, pp. 518–529, 2011.
- [236] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, “Efficient backprop,” in *Neural networks: Tricks of the trade, LNCS*, vol. 1524, pp. 9–50, Springer, 1998.
- [237] B. Recht, C. Re, S. Wright, and F. Niu, “Hogwild: A lock-free approach to parallelizing stochastic gradient descent,” in *Advances in Neural Information Processing Systems*, pp. 693–701, Curran Associates, Inc., 2011.
- [238] L. Al Shalabi, Z. Shaaban, and B. Kasasbeh, “Data mining: A preprocessing engine,” *Journal of Computer Science*, vol. 2, no. 9, pp. 735–739, 2006.
- [239] R. M. Sirkin, *Statistics for the social sciences*. Sage Publications, 2005.
- [240] R. H. Randles, “Wilcoxon signed rank test,” *Encyclopedia of statistical sciences (1988)*, Wiley, 1988.
- [241] C. Gough, I. Steiner, and W. Saunders, *Energy Efficient Servers: Blueprints for Data Center Optimization*. Apress, Berkeley, CA, 2015.
- [242] E. Oró, V. Depoorter, A. Garcia, and J. Salom, “Energy efficiency and renewable energy integration in data centres. strategies and modelling review,” *Renewable and Sustainable Energy Reviews*, vol. 42, pp. 429–445, 2015.
- [243] A. Greenberg, J. Hamilton, D. A. Maltz, and P. Patel, “The cost of a cloud: Research problems in data center networks,” *SIGCOMM Computer Communication Review*, vol. 39, no. 1, pp. 68–73, 2008.
- [244] A. M. Ferreira and B. Pernici, “Managing the complex data center environment: an integrated energy-aware framework,” *Computing*, pp. 1–41, 2014.

REFERENCES

- [245] T. Daim, J. Justice, M. Krampits, M. Letts, G. Subramanian, and M. Thirumalai, “Data center metrics: an energy efficiency model for information technology managers,” *Management of Environmental Quality: An International Journal*, vol. 20, no. 6, pp. 712–731, 2009.
- [246] V. Avelar, D. Azevedo, A. French, and E. N. Power, “Pue: A comprehensive examination of the metric,” *The Green Grid, White paper*, vol. 49, 2012.
- [247] C. B. Dan Azevedo, J. P. Michael Patterson, and R. Tipley, “Carbon usage effectiveness (cue): a green grid data center sustainability metric,” *The Green Grid, White Paper*, no. 32, 2010.
- [248] M. Patterson, D. Azevedo, C. Belady, and J. Pouchet, “Water usage effectiveness (wue): a green grid data center sustainability metric,” *The Green Grid, White Paper*, vol. 35, 2011.
- [249] M. Banks, E. Benjamin, T. Calderwood, R. G. Llera, and J. Pflueger, “Electronics disposal efficiency (ede): An it recycling metric for enterprises and data centers,” *Green Grid, White Paper*, no. 53, 2012.
- [250] A. Capozzoli, G. Serale, L. Liuzzo, and M. Chinnici, “Thermal metrics for data centers: A critical review,” *Energy Procedia*, vol. 62, pp. 391–400, 2014.
- [251] R. Tozer and M. Salim, “Data center air management metrics-practical approach,” in *Proceedings of the 12th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, pp. 1–8, IEEE, 2010.
- [252] I. Munteanu, V. Debusschere, S. Bergeon, and S. Bacha, “Efficiency metrics for qualification of datacenters in terms of useful workload,” in *Proceedings of the PowerTech conference (POWERTECH)*, pp. 1–6, IEEE, 2013.
- [253] B. Schaeppi, T. Bogner, A. Schloesser, L. Stobbe, and M. D. De Asuncao, “Metrics for energy efficiency assessment in data centers and server rooms,” in *Proceedings of the Electronics Goes Green (EGG)*, pp. 1–6, IEEE, 2012.
- [254] C. Fiandrino, D. Kliazovich, P. Bouvry, and A. Y. Zomaya, “Performance and energy efficiency metrics for communication systems of cloud computing

REFERENCES

- data centers,” *IEEE Transactions on Cloud Computing*, vol. 5, no. 4, pp. 738–750, 2017.
- [255] L. Sisó, J. Salom, M. Jarus, A. Oleksiak, and T. Zilio, “Energy and heat-aware metrics for data centers: Metrics analysis in the framework of coolemall project,” in *Proceedings of the third International Conference on Cloud and Green Computing (CGC)*, pp. 428–434, 2013.
- [256] A. I. Aravanis, A. Voulkidis, J. Salom, J. Townley, V. Georgiadou, A. Oleksiak, M. R. Porto, F. Roudet, and T. Zahariadis, “Metrics for assessing flexibility and sustainability of next generation data centers,” in *Proceedings of IEEE Globecom Workshops*, pp. 1–6, Dec 2015.
- [257] A. Capozzoli, M. Chinnici, M. Perino, and G. Serale, “Review on performance metrics for energy efficiency in data center: The role of thermal management,” in *Proceedings of Third International Workshop on Energy Efficient Data Centers, Cambridge, UK*, pp. 135–151, 2015.
- [258] L. Sis, J. Salom, M. Jarus, A. Oleksiak, and T. Zilio, “Energy and heat-aware metrics for data centers: Metrics analysis in the framework of coolemall project,” in *Proceedings of the third International Conference on Cloud and Green Computing (CGC)*, pp. 428–434, 2013.
- [259] D. Chen, B. Pernici, E. Henis, R. I. Kat, D. Sotnikov, C. Cappiello, A. M. Ferreira, M. Vitali, T. Jiang, and J. Liu, “Usage centric green performance indicators,” *ACM SIGMETRICS Performance Evaluation Review*, vol. 39, no. 3, pp. 92–96, 2011.
- [260] L. Wang and S. U. Khan, “Review of performance metrics for green data centers: a taxonomy study,” *The Journal of Supercomputing*, vol. 63, no. 3, pp. 639–656, 2013.
- [261] M. Wiboonrat, “An empirical study on data center system failure diagnosis,” in *Proceedings of the 3rd International Conference on Internet Monitoring and Protection*, pp. 103–108, IEEE, 2008.

REFERENCES

- [262] ASHRAE, “Thermal guidelines for data processing environments-expanded data center classes and usage guidance,” *Technical report, ASHRAE technical committee (TC) 9.9*, vol. 9, 2011.
- [263] “Design considerations for datacom equipment centers,” *Technical report, American Society of Heating Refrigeration and Air Conditioning Engineers*, 2005.
- [264] M. U. S. Khan and S. U. Khan, “Smart data center,” in *Handbook on Data Centers*, pp. 247–262, Springer, 2015.
- [265] S. Murugesan and G. R. Gangadharan, “Green IT: An overview,” in *Harnessing green IT: Principles and practices*, ch. 1, pp. 1–21, Wiley, 2013.
- [266] R. Basmadjian, P. Bouvry, G. Da Costa, L. Gyarmati, D. Kliazovich, S. Lafond, L. Lefevre, H. De, J.-M. P. Meer, and R. Pries, “Green data centers,” *Large-scale Distributed Systems and Energy Efficiency: A Holistic View. John Wiley & Sons.*, pp. 159–196, 2015.
- [267] S. Murugesan and G. R. Gangadharan Eds., “Harnessing green it: Principles and practices,” *Wiley - IEEE Computer Society Press*, 2013.
- [268] D. Abts, M. R. Marty, P. M. Wells, P. Klausler, and H. Liu, “Energy proportional datacenter networks,” *ACM SIGARCH Computer Architecture News*, vol. 38, no. 3, pp. 338–347, 2010.
- [269] S. Chahal, I. K. Anandarao, S. C. Planner, I. C. Peters, I. I. Director, S. Healy, I. N. Wayman, S. Engineer, and I. S. Owen, “Implementing cloud storage metrics to improve it efficiency and capacity management,” *Intel IT, IT best practices, Cloud Computing and IT Efficiency*, 2011.
- [270] R. K. Yin, *Applications of case study research*. Sage, 2011.
- [271] R. K. Yin, “The case study as a serious research strategy,” *Knowledge*, vol. 3, no. 1, pp. 97–114, 1981.
- [272] T. C. Ferreto, M. A. Netto, R. N. Calheiros, and C. A. De Rose, “Server consolidation with migration control for virtualized data centers,” *Future Generation Computer Systems*, vol. 27, no. 8, pp. 1027–1034, 2011.

REFERENCES

- [273] R. W. Ahmad, A. Gani, S. H. A. Hamid, M. Shiraz, A. Yousafzai, and F. Xia, “A survey on virtual machine migration and server consolidation frameworks for cloud data centers,” *Journal of Network and Computer Applications*, vol. 52, pp. 11–25, 2015.
- [274] G. Magklis, G. Semeraro, D. H. Albonesi, S. G. Dropsho, S. Dwarkadas, and M. L. Scott, “Dynamic frequency and voltage scaling for a multiple-clock-domain microprocessor,” *Micro, IEEE*, vol. 23, no. 6, pp. 62–68, 2003.
- [275] N. Ahuja, “Datacenter power savings through high ambient datacenter operation: Cfd modeling study,” in *Proceedings of the 28th Annual IEEE Semiconductor Thermal Measurement and Management Symposium (SEMI-THERM)*, pp. 104–107, IEEE, 2012.
- [276] J. Tu, L. Lu, M. Chen, and R. K. Sitaraman, “Dynamic provisioning in next-generation data centers with on-site power production,” in *Proceedings of the fourth international conference on Future energy systems*, pp. 137–148, ACM, 2013.
- [277] P. Gilbert, K. Ramakrishnan, and R. Diersen, “From data center metrics to data center analytics: How to unlock the full business value of dcim,” *CA technologies: white paper*, 2013.
- [278] M. H. Beitelmal and C. D. Patel, “Model-based approach for optimizing a data center centralized cooling system,” *Hewlett-Packard (HP) Lab Technical Report*, 2006.
- [279] T. Hartman, “Loop chiller plant design dramatically lowers chilled water costs,” tech. rep., Hartman Co., Marysville, WA (US), 1999.
- [280] D. I. Stewart, A. Buratti, and P. Debagnard, “Data center operational efficiency best practices,” *Reaserch Report, IBM Global Technology Services*, 2012.
- [281] J. Colby and J. Herzing, “The why, what, and how of data center humidification,” tech. rep., Engineered Systems, www.digital.bnppmedia.com, 2016.
- [282] T. Evans, “Humidification strategies for data centers and network rooms,” *APC distributors, White Paper, 58*, 2004.

REFERENCES

- [283] Q. Gu, P. Lago, H. Muccini, and S. Potenza, “A categorization of green practices used by dutch data centers,” *Procedia Computer Science*, vol. 19, pp. 770–776, 2013.
- [284] D. Drutskoy, E. Keller, and J. Rexford, “Scalable network virtualization in software-defined networks,” *IEEE Internet Computing*, vol. 17, no. 2, pp. 20–27, 2013.
- [285] S. Berger, R. Cáceres, D. Pendarakis, R. Sailer, E. Valdez, R. Perez, W. Schildhauer, and D. Srinivasan, “Tvdc: managing security in the trusted virtual datacenter,” *ACM SIGOPS Operating Systems Review*, vol. 42, no. 1, pp. 40–47, 2008.
- [286] V. Dinesh Reddy, B. Setz, G. Rao, G. Gangadharan, and M. Aiello, “Metrics for sustainable data centers,” *IEEE Transactions on Sustainable Computing*, vol. 2, no. 3, pp. 290–303, 2017.
- [287] J. Kennedy, “The particle swarm: social adaptation of knowledge,” in *Proceedings of the IEEE International Conference on Evolutionary Computation*, pp. 303–308, IEEE, 1997.
- [288] Y. Del Valle, G. K. Venayagamoorthy, S. Mohagheghi, J.-C. Hernandez, and R. G. Harley, “Particle swarm optimization: basic concepts, variants and applications in power systems,” *IEEE Transactions on evolutionary computation*, vol. 12, no. 2, pp. 171–195, 2008.
- [289] J. Kennedy, “The particle swarm: social adaptation of knowledge,” in *Proceedings of IEEE International Conference on Evolutionary Computation*, pp. 303–308, IEEE, 1997.
- [290] C. R. Reeves, “Using genetic algorithms with small populations,” in *Proceedings of the Fifth International Conference on Genetic Algorithms*, pp. 92–99, Morgan Kaufmann, 1993.
- [291] S. Tsutsui and A. Ghosh, “Genetic algorithms with a robust solution searching scheme,” *IEEE transactions on Evolutionary Computation*, vol. 1, no. 3, pp. 201–208, 1997.

REFERENCES

- [292] J. M. Kaplan, W. Forrest, and N. Kindler, “Revolutionizing data center energy efficiency,” tech. rep., McKinsey & Company, 2008.
- [293] C. L. Belady and C. G. Malone, “Metrics and an infrastructure model to evaluate data center efficiency,” in *Proceedings of the ASME InterPACK Conference*, pp. 751–755, American Society of Mechanical Engineers, 2007.
- [294] M. Blackburn, “The green grid data center compute efficiency metric: Dcce,” *The Green Grid, White Paper*, no. 34, 2010.
- [295] L. H. SeGO, A. Márquez, A. Rawson, T. Cader, K. Fox, W. I. Gustafson Jr, and C. J. Mundy, “Implementing the data center energy productivity metric,” *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, vol. 8, no. 4, p. 30, 2012.
- [296] K. G. Brill, “Data center energy efficiency and productivity,” *White Paper, The Uptime Institute*, <http://www.uptimeinstitute.org/whitepapers>, 2007.
- [297] J. Haas, M. Monroe, J. Pflueger, J. Pouchet, P. Snelling, A. Rawson, and F. Rawson, “Proxy proposals for measuring data center productivity,” *The Green Grid*, 2009.
- [298] L. H. SeGO, A. Márquez, A. Rawson, T. Cader, K. Fox, W. I. Gustafson, Jr., and C. J. Mundy, “Implementing the data center energy productivity metric,” *Journal of Emerging Technology in Computing System*, vol. 8, pp. 1–22, nov 2012.
- [299] “What is PUE/DCiE? how to calculate, what to measure,” www.42u.com/measurement/pue-dcie.htm.
- [300] N. Rasmussen, “Guidelines for specification of data center power density,” *APC, White paper, 120*, 2005.
- [301] T. Committee, “Green grid metrics describing data center power efficiency,” *Green Grid Industry Consortium, White paper, 7*, 2007.
- [302] L. Newcombe, Z. Limbuwala, P. Latham, and V. Smith, “Data center fixed to variable energy ratio metric (dc-fver),” *Technical Report, BCS-The Chartered Institute for IT*, 2012.

REFERENCES

- [303] J. R. Stanley, K. Brill, and J. Koomey, “Four metrics define data center greenness,” *White paper, Uptime Institute*, 2007.
- [304] G. I. Promotion Council, “New data center energy efficiency evaluation index. dppe (datacenter performance per energy),” *Measurement Guidelines (Ver 2.05)*, March, 2012.
- [305] T. Wilde, A. Auweter, M. Patterson, H. Shoukourian, H. Huber, A. Bode, D. Labrenz, and C. Cavazzoni, “Dwpe, a new data center energy-efficiency metric bridging the gap between infrastructure and workload,” in *Proceedings of the International Conference on High Performance Computing Simulation (HPCS)*, pp. 893–901, July 2014.
- [306] M. K. Patterson, S. W. Poole, C.-H. Hsu, D. Maxwell, W. Tschudi, H. Coles, D. J. Martinez, and N. Bates, “TUE, a new energy-efficiency metric applied at ornls jaguar,” in *Proceedings of International Supercomputing Conference*, pp. 372–382, Springer, 2013.
- [307] B. Lajevardi, K. R. Haapala, and J. F. Junker, “An energy efficiency metric for data center assessment,” in *Proceedings of the Industrial and Systems Engineering Research Conference*, 2014.
- [308] Y. Joshi and P. Kumar, *Energy efficient thermal management of data centers*. Springer Science & Business Media, 2012.
- [309] E. Jaureguiualzo, “Pue: The green grid metric for evaluating the energy efficiency in data center-measurement method using the power demand,” in *Proceedings of IEEE 33rd International Conference on Telecommunications Energy (INTELEC)*, pp. 1–8, 2011.
- [310] Google, “Efficiency: How we do it, measuring and improving our energy use,” <http://www.google.com/about/datacenters/efficiency>.
- [311] R. Tipley, T. Cader, J. Froedge, R. Tozer, and R. Wofford, “Pue scalability metric and statistical analysis,” *The Green Grid, White paper*, 49, 2009.
- [312] L. A. Barroso, “The price of performance,” *ACM Queue*, vol. 3, no. 7, pp. 48–53, 2005.

REFERENCES

- [313] W. chun Feng and K. Cameron, “The green500 list: Encouraging sustainable supercomputing,” *Computer*, vol. 40, pp. 50–55, Dec 2007.
- [314] “Swap (space, watts and performance) metric,” www.sun.com/servers/coolthreads/swap/, 2011.
- [315] M. K. Patterson, S. W. Poole, C.-H. Hsu, D. Maxwell, W. Tschudi, H. Coles, D. J. Martinez, and N. Bates, “TUE, a new energy-efficiency metric applied at ornls jaguar,” in *Proceedings of International Supercomputing Conference*, pp. 372–382, Springer, 2013.
- [316] P. A. Inc, “Par⁴- an energy efficiency rating,” *Underwriter’s Laboratories*, www.par4.org, 2011.
- [317] K.-P. Lee and H.-L. Chen, “Analysis of energy saving potential of air-side free cooling for data centers in worldwide climate zones,” *Energy and Buildings*, vol. 64, pp. 103–112, 2013.
- [318] E. Pakbaznia and M. Pedram, “Minimizing data center cooling and server power costs,” in *Proceedings of the 2009 ACM/IEEE international symposium on Low power electronics and design*, pp. 145–150, ACM, 2009.
- [319] H. Sun, P. Stolf, J. M. Pierson, and G. Da Costa, “Energy-efficient and thermal-aware resource management for heterogeneous datacenters,” *Sustainable Computing: Informatics and Systems*, vol. 4, pp. 292–306, 2014.
- [320] G. Varsamopoulos, M. Jonas, J. Ferguson, J. Banerjee, S. K. Gupta, and I. Lab, “Using transient thermal models to predict cyberphysical phenomena in data centers,” *Sustainable Computing: Informatics and Systems*, vol. 3, no. 3, pp. 132–147, 2013.
- [321] P. Mathew, “Self-benchmarking guide for data centers: Metrics, benchmarks, actions,” *Technical Report, LBNL-3393E, Lawrence Berkeley National Laboratory*, 2010.
- [322] O. VanGeet, “Best practices guide for energy-efficient data center design,” *Technical Report, US Department of Energy Federal Energy Management Program*, 2011.

REFERENCES

- [323] J. W. Vangilder and S. K. Shrivastava, “Real-time prediction of rack-cooling performance,” *ASHRAE transactions*, pp. 151–162, 2006.
- [324] TRANE, “free cooling using water economizers,” *TRANE-engineers newsletter*, vol. 37, no. 3, 2008.
- [325] M. Patterson, B. Tschudi, O. Vangeet, J. Cooley, and D. Azevedo, “Ere: A metric for measuring the benefit of reuse energy from a data center,” *The Green Grid, White paper, 29*, 2010.
- [326] Emerson, “Recycling ratios: The next step for data center sustainability,” *White Paper, www.emersonnetworkpower.com*, 2011.
- [327] R. Sharma, A. Shah, C. Bash, T. Christian, and C. Patel, “Water efficiency management in datacenters: Metrics and methodology,” in *Proceedings of the IEEE International Symposium on Sustainable Systems and Technology*, pp. 1–6, IEEE, 2009.
- [328] B. Subramaniam and W.-c. Feng, “The green index: A metric for evaluating system-wide energy efficiency in hpc systems,” in *Proceedings of the 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum, (USA)*, pp. 1007–1013, IEEE Computer Society, 2012.
- [329] A. Cook, “Technology carbon efficiency,” *CS Technology, White Paper*, 2007.
- [330] M. Bana, A. Docca, and S. Devis, “From compromised to optimized-10 million saved in one data center,” *Future Facilities Ltd., White paper FFL-004*, 2014.
- [331] D. Wong and M. Annavaram, “Knightshift: Scaling the energy proportionality wall through server-level heterogeneity,” in *Proceedings of the 45th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 119–130, 2012.
- [332] F. Ryckbosch, S. Polfiet, and L. Eeckhout, “Trends in server energy proportionality,” *Computer*, vol. 44, pp. 69–72, 2011.

REFERENCES

- [333] C. Bekas and A. Curioni, “A new energy aware performance metric,” *Computer Science-Research and Development*, vol. 25, no. 3-4, pp. 187–195, 2010.
- [334] G. Varsamopoulos, Z. Abbasi, and S. K. Gupta, “Trends and effects of energy proportionality on server provisioning in data centers,” in *Proceedings of the International Conference on High Performance Computing (HiPC)*, pp. 1–11, 2010.
- [335] V. Villebonnet, G. Da Costa, L. Lefevre, J.-M. Pierson, and P. Stolf, “big, medium, little: Reaching energy proportionality with heterogeneous computing scheduler,” *Parallel Processing Letters*, vol. 25, p. 1541006, 2015.
- [336] C. Belady and M. Patterson, “The green grid productivity indicator,” *The Green Grid, White paper, 15*, 2008.
- [337] R. L. Sawyer, “Making large ups systems more efficient,” *Elektron Journal-South African Institute of Electrical Engineers*, vol. 23, no. 6, p. 65, 2006.
- [338] L. Giuntini, “Modeling ups efficiency as a function of load,” in *Proceedings of the International Conference on Power Engineering, Energy and Electrical Drives (POWERENG)*, pp. 1–6, IEEE, 2011.
- [339] N. Rasmussen, “Understanding power factor, crest factor, and surge factor,” *Schneider Electric’s Data center Science, White Paper*, vol. 1, 2006.
- [340] J. Neudorfer and F. J. Ohlhorst, “Data center efficiency metrics and methods,” *White Paper, TechTarget Inc.*, 2010.
- [341] J. W. Vangilder and S. K. Shrivastava, “Capture index: An airflow-based rack cooling performance metric.,” *ASHRAE transactions*, vol. 113, no. 1, 2007.
- [342] O. Sarood, E. Meneses, and L. V. Kale, “A cool way of improving the reliability of hpc machines,” in *Proceedings of International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–12, 2013.
- [343] M. Norota, H. Hayama, M. Enai, T. Mori, and M. Kishita, “Research on efficiency of air conditioning system for data-center,” in *Proceedings of the*

REFERENCES

- 25th International Telecommunications Energy Conference. INTELEC '03.*, pp. 147–151, Oct 2003.
- [344] R. Tozer, C. Kurkjian, and M. Salim, “Air management metrics in data centers.,” *ASHRAE transactions*, vol. 115, no. 1, 2009.
- [345] A. Satman and N. Yalcinkaya, “Heating and cooling degree-hours for turkey,” *Energy*, vol. 24, no. 10, pp. 833–840, 1999.
- [346] M. Herrlin, “Rack cooling effectiveness in data centers and telecom central offices:the rack cooling index (rci),” *ASHRAE transaction*, vol. 111, no. 2, 2005.
- [347] M. K. Herrlin, “Airflow and cooling performance of data centers: two performance metrics,” *ASHRAE transactions*, vol. 114, no. 2, 2008.
- [348] R. K. Sharma, C. E. Bash, and C. D. Patel, “Dimensionless parameters for evaluation of thermal design and performance of large-scale data centers,” in *Proceedings of the 8th ASME/AIAA Joint Thermophysics and Heat Transfer Conference*, pp. 1–1, 2002.
- [349] M. K. Herrlin, “Improved data center energy efficiency and thermal performance by advanced airflow analysis,” in *Proceedings of the Digital Power Forum*, pp. 10–12, 2007.
- [350] X. Qian, Z. Li, and Z. Li, “A thermal environmental analysis method for data centers,” *International Journal of Heat and Mass Transfer*, vol. 62, pp. 579–585, 2013.
- [351] V. Rodoplu and T. Meng, “Bits-per-joule capacity of energy-limited wireless networks,” *IEEE Transactions on Wireless Communications*, vol. 6, pp. 857–865, March 2007.
- [352] R. S. Couto, M. E. M. Campista, and L. H. M. Costa, “A reliability analysis of datacenter topologies,” in *Proceedings of the Global Communications Conference (GLOBECOM)*, pp. 1890–1895, IEEE, 2012.

REFERENCES

- [353] A. Alimian, B. Nordman, and D. Kharitonov, “Network and telecom equipment-energy and performance assessment,” *ECR initiative Draft 3.0.1*, www.ecrinitiative.org/pdfs.
- [354] G. Force, “Harmonizing global metrics for data center energy efficiency,” tech. rep., Global Taskforce Reaches Agreement Regarding Data Center Productivity, 2014.
- [355] ATIS, “Energy efficiency for telecommunication equipment: Methodology for measurement and reporting and transport requirements,” *ATIS*, <https://www.atis.org/docstore/product.aspx?id=28263>, 2014.
- [356] G. Schulz, *The green and virtual data center, Chapter 5*. CRC Press, 2009.
- [357] E. L. Miller and R. H. Katz, “Input/output behavior of supercomputing applications,” in *Proceedings of the 1991 ACM/IEEE conference on Supercomputing*, pp. 567–576, ACM, 1991.
- [358] S. Al-Haj and E. Al-Shaer, “Measuring firewall security,” in *Proceedings of the 4th Symposium on Configuration Analytics and Automation (SAFECONFIG)*, pp. 1–4, Oct 2011.
- [359] J. Snyder, “Firewalls in the data center: Main strategies and metrics,” *Opus One, White paper*.
- [360] M. Arregoces and M. Portolani, “Performance metrics of data center devices,” in *Data center fundamentals*, Cisco, 2003.
- [361] W. Boyer and M. McQueen, “Ideal based cyber security technical metrics for control systems,” in *Proceedings of the International Workshop on Critical Information Infrastructures Security*, pp. 246–260, Springer, 2007.
- [362] D. Newman, “Benchmarking terminology for firewall performance,” 1999.
- [363] R. L. Villars, R. Perry, J. Daly, and J. Scaramella, “Measuring the business value of converged infrastructure in the data center,” *International Data Corporation (IDC), White Paper sponsored by HP*, 2011.

REFERENCES

- [364] “IT financial metrics primer: Eleven essential metrics for optimizing the business value of it,” *APPTIO, White paper*, www.apptio.com/resource-center.
- [365] M. Wara, “Is the global carbon market working?,” *Nature*, vol. 445, no. 7128, pp. 595–596, 2007.
- [366] V. McMorrow, “Cost of unplanned data center outages since 2010,” *Emerson Network Power*, www.emersonnetworkpower.com/en-US/About/NewsRoom/NewsReleases/Pages/Emerson-Ponemon-Cost-Unplanned-Data-Center-Outages.aspx, 2013.
- [367] M. Wiboonrat, “An empirical study on data center system failure diagnosis,” in *Proceedings of the Third International Conference on Internet Monitoring and Protection*, pp. 103–108, IEEE, 2008.
- [368] W. Torell and V. Avelar, “Mean time between failure: Explanation and standards,” *Schneider Electric - Data Center Science, White paper*, 78, 2004.
- [369] “Calculate the ROI of data center investments,” *Forrester, White Paper*, www.forrester.com/report/Calculate+The+ROI+Of+Data+Center+Investments/-/E-RES101661, 2013.
- [370] N. Rasmussen, “Determining total cost of ownership for data center and network room infrastructure,” *Schneider Electric, White Paper*, 8, 2011.

REFERENCES

List of Publications

- [1] V. Dinesh Reddy, B. Setz, G. Subrahmanya V.R.K. Rao, G.R. Gangadharan, Marco Aiello. **Metrics for Sustainable Data Centers**, *IEEE Transactions on Sustainable Computing*, Vol. 2, No. 3, pp. 290-303, DOI:10.1109/TSUSC.2017.2701883, 2017.
- [2] V. Dinesh Reddy, G. R. Gangadharan, and G. Subrahmanya V. R. K. Rao. **Energy-aware virtual machine allocation and selection in cloud data centers**, *Soft Computing*, Springer, ISSN 1433-7479, DOI: 10.1007/s00500-017-2905-z, 2017.
- [3] V. Dinesh Reddy, B. Setz, G.S.V.R.K. Rao, G. R. Gangadharan, and M. Aiello. **Best Practices for Sustainable Data Center**, *IEEE IT Professional* (In Press), 2017.
- [4] V. Dinesh Reddy and G. R. Gangadharan. **Towards an Internet of Things framework for financial services sector**, *Proceedings of the 3rd International Conference on Recent Advances in Information Technology (RAIT)*, Dhanbad, pp. 177-181, DOI: 10.1109/RAIT.2016.7507897, IEEE, 2016
- [5] Chandra Shekhar Verma, V. Dinesh Reddy, G. R. Gangadharan and Atul Negi. **Energy Efficient Virtual Machine Placement in Cloud Data Centers Using Modified Intelligent Water Drop Algorithm**, *Proceedings of the 13th International Conference on Signal Image Technology & Internet Based Systems (SITIS)*, Jaipur, pp. 13-20, DOI: 10.1109/SITIS.2017.14, IEEE, 2017
- [6] V. Dinesh Reddy and G.R. Gangadharan. **Forecasting Data Center Energy Demand: A Deep Learning Approach**, *Sustainable Computing, Informatics and Systems*, Elsevier (Journal Paper Under Review).

LIST OF PUBLICATIONS

- [7] V. Dinesh Reddy and G.R. Gangadharan, M. Aiello. **Energy efficient resource management in cloud data centers using a hybrid evolutionary algorithm**, *Soft Computing*, Springer (Journal Paper Under Review).

Appendix A

Overview of Techniques Used

A.1 Particle Swarm Optimization (PSO)

PSO is a stochastic optimization algorithm based on social simulation models [287, 288]. PSO is a popular method to find optimum of a numerical function defined on a continuous domain. PSO is a collective, anarchic, iterative method with emphasis on cooperation. Each particle in the swarm is able to communicate to some other in the position and quality of the best site it knows. The algorithm is initialized with a population of search points (particle (X_i)) that moves stochastically in the search space. Each particle moves in a multi dimensional space according to its own velocity, particle's best performance, and the best performance of its best informant. For each iteration, PSO updates velocity and position of each particle according to Equation A.1 and Equation A.2. After updation PSO calculates the fitness of the particles and flies in the search space towards the local and global best solutions in a navigated way.

$$v_{i+1} = \omega v_i + \phi_1 \beta_1 (p_i - x_i) + \phi_2 \beta_2 (p_g - x_i) \quad (\text{A.1})$$

$$x_{i+1} = x_i + v_{i+1} \quad (\text{A.2})$$

PSO has stochastic mechanism that makes the particles to exploit some specific areas and rise up from the local optima [289]. PSO has ability to keep track the particle that has best value in the population ((\hat{L}_j)) and this is global best. When a particle takes part of the population, best experience or position of the one particle is called local best ((L_i)). The particle velocity is updated using these best

A.2 Discrete Binary Version of PSO

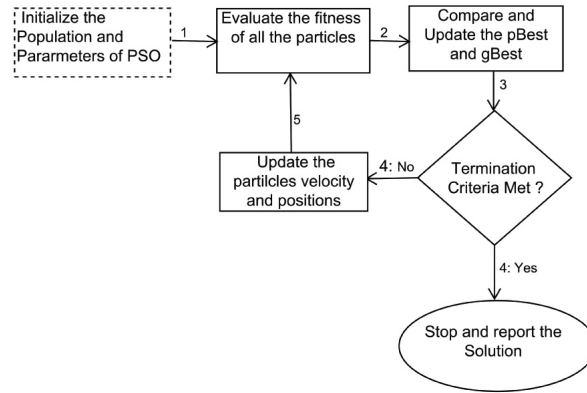


Figure A.1: Flow of Particle Swarm Optimization

values according to the Equation A.3. Further, each particle updates its position using the updated velocity, and its current position as shown in Equation A.4.

$$V_i(t + 1) = wV_i(t) + k_1r_1(t)(L_i(t) - X_i(t)) + k_2r_2(t)(\hat{L}_j(t) - X_i(t)) \quad (\text{A.3})$$

$$X_i(t + 1) = X_i(t) + V_i(t + 1) \quad (\text{A.4})$$

Where k_1 and k_2 are the learning factors which determines the convergence properties of the algorithm, w is the inertia weight coefficient that determines how the previous velocity of the particle influences the velocity in the next iteration, and r_1, r_2 are the random number between $(0,1)$.

A.2 Discrete Binary Version of PSO

The particle swarm works by adjusting trajectories through manipulation of each coordinate of a particle. However, many optimization problems are set in a space featuring discrete, qualitative distinctions between variables and between levels of variables. In the binary version of the PSO, the trajectories are changes in the probability that a coordinate will take on binary value $(0 \text{ or } 1)$.

A.3 Genetic Algorithm (GA)

A genetic algorithm is an adaptive heuristic search algorithm which can be applied to both constrained and unconstrained optimization problems [290, 291].

A.3 Genetic Algorithm (GA)

Genetic algorithms are built based on the transformative process of natural selection, mating, and reproduction that mimics biological evolution. A GA simulates this process with natural populations of individuals that evolve according to the principles of natural selection and survival of the fittest. Each individual of the population is coded to make a chromosome that represents a candidate solution for the given optimization problem.

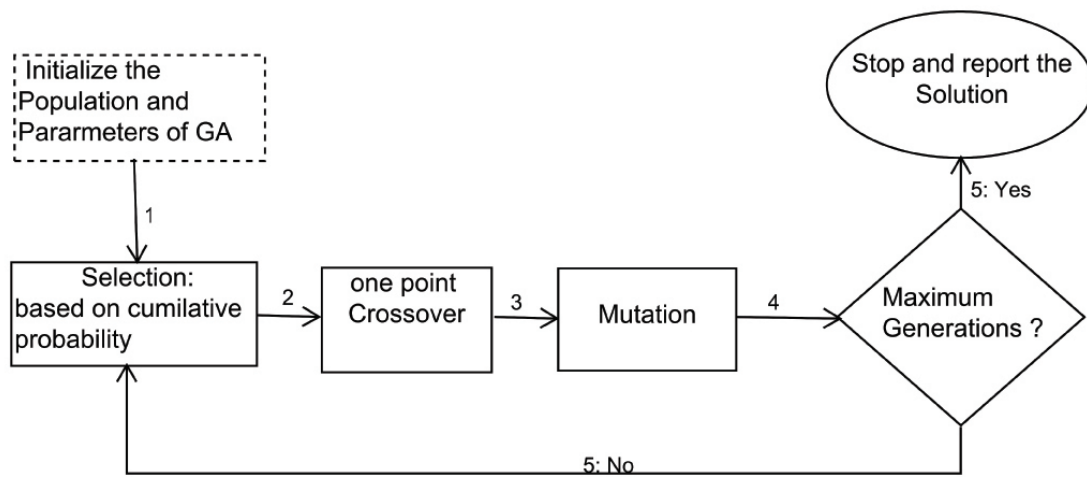


Figure A.2: Flow of Genetic Algorithm

Individuals who are more successful in adapting to their environment are selected based on their cumulative probability calculated based on fitness. The selected individuals reproduce the off-springs by exchanging pieces of their genetic information (characteristics of parents) depending on the cross over rate, whilst individuals who are less fit will be eliminated. This phase is known as crossover. Mutation rate is basically a measure of the likeness that random elements of a candidate solution are flipped to get a new solution.

To improve the offspring solutions Mutation operator is applied by altering some genes in the strings depending on the mutation rate. We made swap mutation where we randomly select the two positions on the chromosome and swap those values. This selection-crossover-mutation cycle repeats until a satisfactory solution is found as shown in Figure A.2.

A.4 Coordination of the particles for Modified Discrete Particle Swarm Optimization (MDPSO) approach

This section present the change in the fitness values of particles for each iteration using MDPSO. We simulated a data center comprising 100 heterogeneous physical machines and 300 virtual machines in the said experimental environment with the following initial parameters:

Population size = 40,

Inertia weight coefficients: $k_1 = 3$ and $k_2 = 2$.

These values are chosen after several experiments and the way these weights coordinate the searching process after each iteration is presented here. We presented the change in the fitness value of each particle, starting from the first iteration to termination with an interval of 20 in Figure A.3.

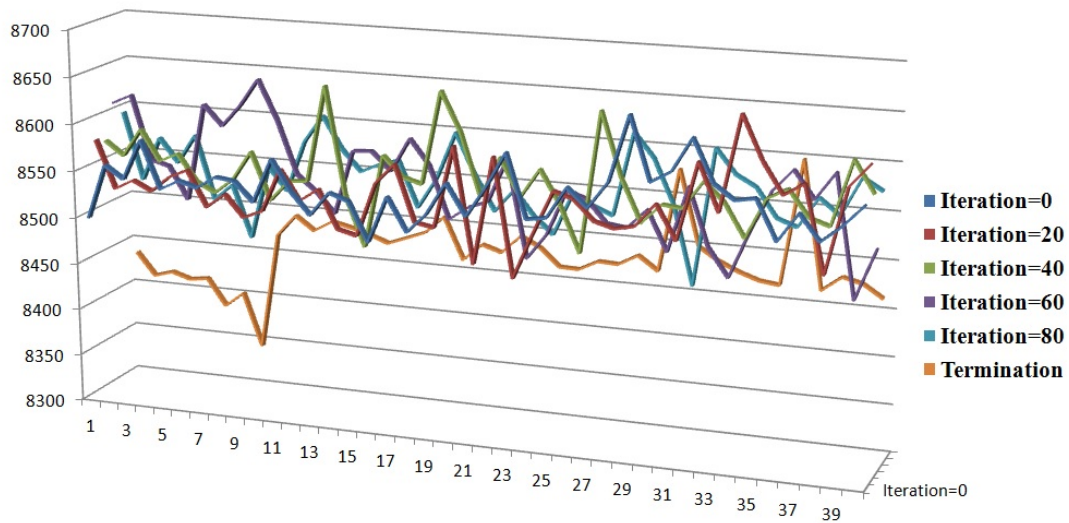


Figure A.3: Particle coordination

Appendix B

Data Center Metrics Definitions

This appendix presents the definitions and details of data center metrics in all dimensions of operations such as energy efficiency, cooling, greenness, performance, thermal and air management, network, security, storage, and financial impact.

B.1 Energy Efficiency Metrics

1. **Adaptability Power Curve (APC)** : The Adaptability Power Curve (APC) metric measures the data center's adaptability with regards to the existing energy usage pattern [256]. The APC is defined as follows:

$$APC = 1 - \frac{\sum_{i=1}^n |K_{APC} E_{Pi} - E_{DCi}|}{\sum_{i=1}^n E_{DCi}}$$

$$K_{APC} = \frac{\sum_{i=1}^n E_{DCi}}{\sum_{i=1}^n E_{Pi}}$$

Where E_{DCi} is the DC energy consumption in kWh, E_{Pi} is the planned energy in kWh, i indicates the time period, n indicates the sample size, and K_{APC} is the adjustment factor for normalizing E_{DCi} and E_{Pi} .

2. **Corporate Average Data Center Efficiency (CADE)** : CADE measures the performance of both the IT equipment and Facilities as a percentage [292]. CADE compares the performance (in terms of overall energy efficiency) of set of data centers of an organization. It measures the return

on investment (ROI) on a green computing initiative. The higher CADE value, the more energy efficient the data center. It accounts aggregate IT load, facility capacity, CPU usage and the aggregate energy devoured by the facility. However, it cannot be useful when basic measurement have not been calculated by the site operators. Also it does not consider how well servers, storage, and communication equipment are utilized.

$$CADE = Facility\ Efficiency \times Asset\ Efficiency$$

where Facility Efficiency and Asset Efficiency are calculated as follows.

Facility Efficiency (FE)

$FE = Facility\ Energy\ Efficiency \times Facility\ Utilization$ where Facility energy efficiency is the ratio of IT load to total power consumed by data center and Facility Utilization is the ratio of Actual IT load over data center capacity.

Asset Efficiency (AE)

$$AE = IT\ Utilization \times IT\ Energy\ Efficiency$$

where IT Asset utilization is measured as average CPU utilization.

3. **Compute Power Efficiency (CPE)** : CPE evaluates the overall efficiency of the data center considering the distribution losses, idle equipment power consumption and other overheads involved [293]. CPE estimates the productivity of a data center as a function of power used and is calculated as follows :

$$CPE = \frac{ITEU}{PUE} = \frac{ITEU * IT\ Equipment\ Power}{Total\ Facility\ Power}$$

where ITEU is IT Equipment Utilization.

4. **DCAdapt (DCA)** : The DCA metric determines the shift of the energy curve after the data center has adapted its operation to a predefined operational mode. It compares the energy profiles before and after the modification of the data center operation mode [256]. DCA is defined as follows:

$$DCA = 1 - \frac{\sum_{i=1}^n |K_{DCA} E_{DCReal\ i} - E_{DCBaseline\ i}|}{\sum_{i=1}^n E_{DCBaseline\ i}}$$

$$K_{DCA} = \frac{\sum_{i=1}^n E_{DCBaseline\ i}}{\sum_{i=1}^n E_{DCReal\ i}}$$

Where $E_{DCReal\ i}$ is the energy consumption in kWh after changing operational modes, $E_{DCBaseline\ i}$ is the energy consumption in kWh before changing operation modes, i indicates the time period, n indicates the sample size and K_{DCA} is the adjustment factor for normalizing the $E_{DCReal\ i}$ and $E_{DCBaseline\ i}$ energy curves.

5. **DCcE and ScE** : Server compute efficiency (ScE) determines the proportion of work done by server to provide primary services which does not include the services like anti-virus, network management, virtualization management, disk defragmentation etc. Aggregation of ScE of all servers with in the data center gives the compute efficiency of the entire data center (known as data center compute efficiency (DCcE)) [294].

ScE and DCcE are time-based metrics. ScE measures the time server used for primary services. The ScE is measured over frequent time periods and calculated as the ratio of sum of the samples where the server used for primary services (p) over total number of measured samples in the same time period (n).

$$ScE = \frac{\sum_{i=1}^n p_i}{n} \times 100$$

DCcE is the average of ScE's of all servers during the same time period.

$$DCcE = \frac{\sum_{j=1}^k ScE_j}{n}$$

where k is the total number of servers in the data center. ScE helps data center managers to improve overall energy use by monitoring the servers which are not running live applications for long stretches. DCcE is used to find inefficiencies lie within a specific data center so that site operators can address these issues there by increasing efficiency over time by right sizing server population.

6. **Data Center Energy Productivity (DCeP)** : DCeP measures to what extent the energy consumed by the data center is used to perform the useful work [295]. DCeP quantifies the work performed (W) which is either useful to the end customer or owner of the data center over the resources consumed by the data center to produce this work (E_{Total}). This can be expressed as follows :

$$DCeP = \frac{\text{Useful Work performed (W)}}{\text{Energy Consumption by Data Center (}E_{Total}\text{)}}$$

To calculate the said useful work and energy consumption we have a user defined “assessment window” (time limit). Useful Computational Unit (UCU) is a measure of useful work done. W is counted as a weighted sum of UCUs derived by different applications during the user defined assessment window.

The work of a data center is unique to its operator. For example, useful work for a search engine may be number of on-line searches performed, while an on line store may use number of sales. This factor makes it very difficult to define a universal productivity metric [296, 297]. The work produced by different applications(W) is defined as

$$W = \sum_{j=1}^{N_a} \sum_{i=1}^{M_j} V_j U_j(t_{ij}, T_{ij}) C_{ij}$$

where

N_a - Total number of applications,

M_j - UCUs produced by j^{th} application during user defined time limit,

V_j - the relative value of a UCU produced by the j^{th} application and

$$C_{ij} = \left\{ \begin{array}{l} 1 \quad \text{if the } i^{th} \text{ UCU of } j^{st} \text{ application} \\ \quad \text{finished within the time limit.} \\ 0 \quad \text{otherwise.} \end{array} \right\}$$

To account for the value of timely completion, defined a utility function $U_j(t_{ij}, T_{ij})$ which is time-based. For application j, t_{ij} is the elapsed time from start to finish of the UCU, and T_{ij} is the exact time by when the UCU of application must be finished [298].

7. **Data Center infrastructure Efficiency (DCiE)** : DCiE is the inverse of PUE and measures the extent of IT power usage to the total power consumed by the data center [7]. It is defined as the ratio of what IT uses versus the total power available. The DCiE value will always be less than 1. The smaller the number the less efficient our utilization is. DCiE can vary from one to zero: $0(worse) \leq DCiE \leq 1(better)$. DCiE is a performance improvement metric gives better way to calculate output power to input power.

Both PUE and DCiE can be used to compare peer data centers and find chances to enhance a data center’s operational effectiveness. Standard values of PUE and DCiE are given in Table B.1 [299].

PUE	Utilization	DCiE(%)
3.0	Very Poor	33
2.5	Poor	40
2.0	Average	50
1.5	High	67
1.2	Very High	83

Table B.1: Efficiency level of PUE and DCiE

8. **Data Center Lighting Density (DCLD)** : DCLD measures the data center power consumption for lighting with respect to area. It is measured in Watt/ft² and defined as

$$DCLD = \frac{\text{Power consumption for data center lighting}}{\text{Total data center space}} \left[\frac{kW}{ft^2} \right]$$

9. **Data Center Power Density (DCPD)** : DCPD specifies variations of density within a data center. The average power density per rack is 5.94kW and the peak value is 7.7kW as reported by Data Center User’s Group [300]. Higher power density can be achieved using converged infrastructure.

10. **Data Center Performance Efficiency (DCPE)** : Data Center Performance Efficiency (DCPE) metric [301] defines the performance efficiency as

the output to input ratio. Inputs for a data center are power and data. Outputs are heat, data and useful work done by the data center.

$$DCPE = \frac{Useful\ Work}{Total\ Facility\ Power}$$

DCPE is the key focus for the industry as it is very difficult to define Useful Work.

11. **Data Center Fixed to Variable Energy Ratio metric (DC-FVER) :** DC-FVER is a way to measure useful work and makes data center operators to focus on operating waste in the software, IT and Mechanical & Electrical infrastructure. DC-FVER measures what percentage of data center energy consumption is variable. This gives us the picture of energy to useful work and to a fixed burden. The DC-FVER is calculated as follows [302] :

$$DC-FVER = \frac{Variable\ energy + Fixed\ Energy}{Variable\ energy}$$

DC-FVER is useful in reducing fixed energy consumption. But measuring the output will vary from operator to operator. DC-FVER can also be measured for the part of the data center allowing full control. For example DC-FVER(IT) is for IT equipment and DC-FVER (Utility) for total data center. For a data center DC-FVER value 10 indicates 90% energy consumption is fixed. For an idle data center DC-FVER value 1 indicates 0% energy consumption is fixed.

12. **Deployed Hardware Utilization Efficiency (DH-UE) :** DH-UE helps increasing servers and storage energy efficiency by knowing the actual number of servers needed to run the peak load considering overhead in virtualization [303].

$$DH - UE_{(servers)} = \frac{M}{Total\ number\ of\ servers\ deployed}$$

where M is the minimum number of servers required to compute peak load. It is also ascertained as the sum of the most astounding of peak loads every server encounters and any overhead incurred in virtualization.

13. **Deployed Hardware Utilization Ratio (DH-UR)** : Idle servers consume energy as much as the servers at peak load. DH-UR estimates how much fraction of deployed servers (or storage) are not in use [303]. But it is difficult to find whether a box is running any useful work. It is defined for servers and storage as follows :

$$DH-UR_{(servers)} = \frac{\text{Number of servers which are live executing some applications}}{\text{Number of servers actually deployed}}$$

$$DH-UR_{(storage)} = \frac{\text{Amount of frequently accessed data within the last 90 days}}{\text{Total storage deployed}}$$

14. **Datacentre Performance Per Energy (DPPE)** : DPPE measures overall data center energy efficiency by computing data center throughput per unit of non-green energy [304]. DPPE defined in the following equation is an integrated metric expressed as a function of four sub metrics, which are ITEU, ITEE, PUE, and GEC.

$$DPPE = ITEU \times ITEE \times \frac{1}{PUE} \times \frac{1}{1 - GEC}$$

$$DPPE = \frac{\text{Actual usage of IT equipment}}{\text{grid energy usage}}$$

DPPE gives the energy savings of different elements in data center. It identifies the each energy saving activity of the data center by considering technologies of cloud computing and source of energy. DPPE can be used for data centers powered by renewable or other sources.

15. **Data center Workload Power Efficiency (DWPE)** : DWPE is a metric to calculate energy efficiency of High-Performance Computing (HPC) system including data centers [305]. DWPE defines the energy efficiency of running a given workload on a specific HPC system in a specific data center. DWPE is calculated by determining the energy efficiency of a specific workload (WPE) and dividing it by the overhead for operating a given system in a certain data center (sPUE).

$$DWPE = \frac{WPE}{sPUE}$$

It considers various aspects including influence of the system management software, and the HPC workloads.

16. **Energy ExpenseS (EES)** : The Energy ExpenseS (EES) metric quantifies the change in the data center expenses after upgrading equipment or the introduction of flexibility mechanisms [256]. The metric is defined as follows:

$$EES = \frac{\sum_{i=1}^n [(E_{DCi} * Cost_{ei} + Eoth_{DCi} * Cost_{othi})_{bas}]}{\sum_{i=1}^n (E_{DCi} * Cost_{ei} + Eoth_{DCi} * Cost_{othi})_{bas}} - \frac{\sum_{i=1}^n [(E_{DCi} * Cost_{ei} + Eoth_{DCi} * Cost_{othi})_{cur}]}{\sum_{i=1}^n (E_{DCi} * Cost_{ei} + Eoth_{DCi} * Cost_{othi})_{bas}}$$

Where E_{DCi} is the total electricity consumed by the data center in primary energy terms, $Cost_{ei}$ is the electricity cost per kWh, $Eoth_{DCi}$ is the total primary energy produced from other sources, $Cost_{othi}$ is the electricity cost per kWh from other sources, i is the time period, *bas* is baseline and *cur* is current.

17. **Energy Wasted Ratio (EWR)** : The Energy Wasted Ratio gives an indication of the actual amount of energy that can be saved [255]. The wasted energy is the integral of wasted power over a certain interval. The wasted power is the difference between the actual power and the ideal power, proportional to the load.

$$EWR = \frac{E_{DC-related\ to\ useless\ workload}[Wh]}{E_{DC}[Wh]}$$

Where $E_{DC-related\ to\ useless\ workload}$ is defined as:

$$\int_{t_0}^{t_1} (P(t) - load(t) * P_{max}) dt = [Wh]$$

18. **IT Hardware Power Overhead Multiplier (H-POM)** : For a piece of hardware, H-POM indicates the wasted power which is not used for computing due to conversion losses and other overheads [303]. For a data center, it

is defined as follows.

$$H-POM = \frac{\text{Total Hardware load of data center at the plug}(AC)}{\text{Total Hardware compute load}(DC)}$$

19. **IT-Power Usage Effectiveness (ITUE)** : ITUE describes the overhead of internal infrastructure such as fans, voltage regulators and power supplies for IT equipment [306]. PUE does not consider upgradation of IT equipment, cooling or power conversion loss. ITUE looks similar to PUE but measure for the IT equipment where PUE is for a site. ITUE is defined as total IT energy used for cooling, power distribution and computing components over total computational energy.

$$ITUE = \frac{\text{Total Energy into IT Equipment}}{\text{Total Energy delivered to Compute Components}}$$

where the compute components include CPU, memory, and storage etc.

20. **Operating System Workload Efficiency (OSWE)** : OSWE measures the efficiency and elasticity of the data center. It is defined as the ratio of count of all OS instances which includes OS inside virtual machines($Count_{OS}$) to the total facilities power (P_{DC}) at the time of assessment [297].

$$OSWE = \frac{Count_{OS}}{P_{DC}}$$

This metric can be used to improve the capacity planning of a data center. But the difficulty lies in counting the OS instances.

21. **Power Density Efficiency (PDE)** : PDE is a metric to link both the thermal and energy efficiency of the data centers [307]. PDE evaluates the impact on energy efficiency by the physical changes made in a rack. PDE can be calculated as follows :

$$PDE = \left[1 + \epsilon \left(\frac{P_{inf}}{P_{IT}} \right) \right]^{-1}$$

where P_{inf} is the power draw of the supporting infrastructure, mainly the cooling system, P_{IT} is the power consumed by the IT equipment in the racks,

v_r is the total volume of the racks, v_s is the total volume of the IT equipment inside the racks, and ϵ is the ratio of rack volume to IT equipment volume (v_r/v_s). The PDE metric enables evaluation of the impact of physical changes and improvements of IT equipment within the racks, but it does not reflect thermal and air flow management inefficiencies within the room [308].

22. **Primary Energy Savings (PEsavings) :** The Primary Energy Savings (PEsavings) metric quantifies the change in the data center energy profile after upgrading equipment or the introduction of flexibility mechanisms and the energetic, economic and environmental upgrade of the data center behaviour [256]. The metric is defined as follows:

$$PEsavings_{DC} = \frac{\sum_{i=1}^n [(E_{DCi} + E_{othDCi})_{bas} - (E_{DCi} + E_{othDCi})_{cur}]}{\sum_{i=1}^n (E_{DCi} + E_{othDCi})_{bas}}$$

Where E_{DCi} is the total electricity consumed by the data center in primary energy terms, E_{othDCi} is the total primary energy produced from other sources, i is the time period, bas is baseline and cur is current.

23. **Power Usage Effectiveness Level 1-4 (PUE₁₋₄) :** Power Usage Effectiveness helps data center professionals to decide energy efficiency of their facility and to screen the effect of their productivity endeavors [246]. PUE is defined as a ratio of the amount of power entering a data center to the power delivered to run the computing equipment which includes server, network equipment, storage etc.

$$PUE_{1-3} = \frac{Total\ Facility\ Energy}{IT\ Equipment\ Energy}$$

PUE is divided into multiple different levels. Level 1-3 use different sensor inputs to determine the IT Equipment Energy. Level 1 measures the energy from Uninterruptible Power Supply (UPS) outputs. Level 2 measures the Power Distribution Unit (PDU) outputs. Level 3 uses the energy measured directly at the IT equipment. The measurements intervals also differ for each

level: monthly / weekly for level 1, daily / hourly for level 2 and continuous for level 3.

Total data center energy measures different types of energy purchased at utility hand off to the data center [309]. For example Google uses the following equation to calculate PUE considering all sources overhead.

$$PUE = \frac{ESIS + EITS + ETX + EHV + ELV + EF}{EITS - ECRAC - EUPS - ELV + ENet}$$

Details of the above equation can be found in [310]. An additional level, PUE₄ has been defined by Siso et. al [255].

$$PUE_4 = \frac{E_{DC}[Wh]}{E_{IT} - E_{fan-rack} - E_{PSU}[Wh]}$$

The PUE₄ focuses on the actual IT work, as it excludes the consumption of fans in the rack and also excludes the consumption of PSUs.

PUE will yield a factor of greater than one. A PUE close to 1 indicates all of the data center source energy is used for computing. The large value of PUE indicates inefficient utilization of energy. While comparing different facilities PUE does not consider the location and climate changes, which is a real drawback.

24. Power Usage Effectiveness Scalability (PUE_{scalability}) :

The Power Usage Effectiveness Scalability (PUE_{scalability}) is another metric defined by the Green Grid Consortium [311]. The metric provides information about the ability to scale the total facility power and whether the infrastructure supports this. PUE_{scalability} is defined as follows:

$$PUE_{scalability} = \frac{m_{actual}}{m_{PUE}}$$

Where m_{actual} is determined by the linear relationship of IT power (P_{IT}) to facility power (P_{DC}) such that

$$P_{DC} = m_{actual}P_{IT} + b,$$

m_{PUE} is the total facility energy usage divided by the IT energy usage.

25. **partial PUE (pPUE)** : pPUE is a metric for monitoring and managing the power for portion of the data center facility [246]. pPUE is the ratio of total energy used by the components within a zone over the Total IT equipment energy inside that zone. In a mixed-use or co-location data centers, some necessary information needed to calculate PUE may not be available. In such cases pPUE metric will be useful to evaluate a zone of the facility, such as IT infrastructure or HVAC system. pPUE is not an alternative to PUE. It prevents the incorrect use of PUE.

$$pPUE = \frac{\textit{Total energy delivered to a zone}}{\textit{Energy used by IT equipment inside that zone}}$$

Zone can be either physical or logical. If a zone does not contain any IT Equipment which generally treated as overhead, pPUE is undefined. If pPUE of one zone is bad, we can improve the energy efficiency by working on that zone's components.

26. **Performance per Watt (PpW)** : PpW measures the energy efficiency of the computer hardware used in a high performance computing system for a specific workload [312]. High PpW indicates that the system is highly energy efficient. The metric covers HPC system hardware (in most cases only the compute-subsystems), software and HPC applications. A well known example of this metric is the FLOPS/watt metric used by the Green500 List [313]. In case of the FLOPS/watt metric, PpW indirectly measures the operations performed for each joule of energy consumed.

$$PpW = \frac{\textit{Performance}}{\textit{Power}}$$

27. **Site Infrastructure Power Overhead Multiplier (SI-POM)** : SI-POM estimates how much power consumed in overhead (distribution losses incurred from transformer to minor building loads) [303]. It is a dimensionless ratio of data center power consumption at utility meter (PCUM) over Total AC power consumption for hardware at the plug (IT-AC power).

$$SI-POM = \frac{\textit{PCUM}}{\textit{IT-AC power}}$$

We can increase SI-POM by increasing energy efficiency of the PDU, UPS and other components.

28. **SWaP: Space, Watts and Performance** : SWaP gives performance of rack-optimized server deployments considering together space, energy, and performance [314]. It is used to compare the performance of different server deployments. The SWaP is calculated as follows:

$$SWaP = \frac{Performance}{Space \times Power Consumption}$$

where Performance is taken from industry-standard benchmarks, Space is the height of the server in rack units, and Power consumption is watts consumed by the system. The higher your SWaP numbers, the less data center space and power need to do a computing job. It is used by data center operators to choose among different servers which deliver the optimum performance for their needs and works in both power constrained and unconstrained situations.

29. **Total-Power Usage Effectiveness (TUE)** : TUE combines PUE and ITUE for a total efficiency picture. TUE is defined as the ratio of total energy use and the specific energy used in the computational components [315].

$$TUE = PUE \times ITUE$$

TUE compares various HPC sites to describe the total energy use from the utility to the silicon.

30. **PAR⁴** : PAR⁴ measures IT equipment power consumption at different load conditions such as idle, loaded, peak and switched off [316]. PAR⁴ enables us to compare IT equipment of different manufacturers, models and generations in terms of energy efficiency. PAR⁴ is used for forecasting the comparative efficiency for each year using a variation of Moore's Law. This metric allows us to classify equipment based on its past and future energy efficiency as well as by IT equipment architecture.

B.2 Cooling Metrics

In this section we describe the metrics that measure the efficiency of different cooling solutions such as chillers, CRAC, CRAH, Economizers.

1. **Air Economizer Utilization Factor (AEUF)** : If the outside air is cold enough, we can use this air to mix properly with the exhaust air in such a way that mixed air will fall in ASHRAE recommended temperature and humidity ranges for the equipment [317]. AEUF tells us how many hours in a year air economizer provides "free" cooling (with out using compressor based cooling). It is defined as follows :

$$AEUF = \frac{\text{free colling hours by Air Economizer}}{24 \times \text{days in a year}} \times 100$$

But for effective use of air side economizers, we must filter the outside air to avoid dust and particulate matter into data center and should do proper mixing of outside and return air depending on the local climate.

2. **Coefficient of Performance Ensemble (CoP)** : CoP of the data center is the ratio of total heat extracted by cooling infrastructure over the work needed to cool the air [318]. To calculate the COP of the data center we consider the heat produced by CRAC blowers and work input to humidify and dehumidify air in data center. If COP is high it needs less work to cool the air which indicates CRAC units are efficient.

$$COP = \frac{\text{Heat extracted by Air Conditioners}}{\text{Net - Work Input of cooling System}}$$

Further, COP is used in analytical models to compute the total power consumed (by the computing part and by the cooling part), predict the cyber-physical phenomena, and to propose thermal aware mapping policies [319, 320]

3. **Data Center Cooling System Efficiency (DCCSE)** : DCCSE indicates the efficiency of HVAC system in terms of power used per unit of cooling output. It is defined as follows :

$$DCSE = \frac{\text{Average cooling system power usage} \left[\frac{kW}{\text{tons}} \right]}{\text{Average cooling load in the data center}}$$

As per LBNL database, the accepted and better values for DCCSE are 0.8 kW/ton and 0.6 kW/ton respectively [321].

4. **Data center Cooling System Sizing Factor (DCSSF)** : This is a dimensionless metric which is the ratio of the total cooling capacity (TC) in tons not considering the back up to the peak chiller load (PL) over one year in tons [321]. This metrics depends on the site location.

$$DCSSF = \frac{TC(w/o\ backup)}{PL}$$

A high DCSSF indicates the necessity to rightsize the capacity of cooling plant and need to increase the partial load efficiency. Using a Variable Frequency Drives in chillers and a flexible design may improve chiller efficiency at part-load.

5. **Energy Efficiency Ratio (EER)** : To determine the energy efficiency of the cooling system the Energy Efficiency Ratio (EER) is used [255]. The EER is defined as follows:

$$EER = \frac{Q_{cooling}[W_{th}]}{P_{cooling}[W_{el}]}$$

Where $Q_{cooling}$ is the heat removed by the cooling system, and $P_{Cooling}$ is the electrical power used by the cooling system.

6. **HVAC System Effectiveness (HSE)** : HVAC system of a data center includes the CRAC, ventilation, a central air handling system, and minor lighting load. The HSE metric is the ratio of the IT equipment energy over the total energy consumption of HVAC system. Higher HSE value indicates that HVAC system is more effective to the IT load [322]. HSE of 0.7 indicates standard value, 1.4 indicates efficient and 2.5 indicates more efficient.

$$HSE = \frac{IT\ Electrical\ Use}{(HVAC + Fuel + Steam + Chilled\ Water)\ Energy\ use}$$

7. **Recirculation Index (RI)** : RI is indicative of an amount of recirculated air to inlet airflow of at least one rack. It is used to compare the cooling performance for raised floor and cold-aisle clusters. The Recirculation Index represents extent of inlet airflow which is not directly from the perforated tiles for a rack [323].

$$RI = \frac{Air\ flow\ not\ from\ the\ perforated\ tiles}{Total\ inlet\ air\ flow\ of\ atleast\ one\ rack}$$

It is used to implement automated method for measuring the cooling performance of a cluster of racks, where a cluster includes two rows separated by a hot aisle.

8. **Water Economizer Utilization Factor (WEUF)** : Supply air of a cooling system is cooled indirectly with water instead of mechanical cooling using water economizer. WEUF characterizes the extent to which water economizer is used for full cooling over a year [324]. It is defined as follows :

$$WEUF = \frac{\text{Water Economizer Hours (full cooling)}}{24 \times \text{days in a year}} \times 100$$

B.3 Greenness Metrics

Green Data centers are emerging, based on the design and operations of data centers in a more efficient and eco-friendly manner. Green initiatives and practices not only reduces GHG emission but also achieves measurable benefits. In this section we explore the metrics which measures the greenness of a data center.

1. **CO₂Savings** : The CO₂Savings metric quantifies the change in the data center CO₂ profile after upgrading equipment or the introduction of flexibility mechanisms [256]. The metric is defined as follows:

$$CO_2Savings =$$

$$\frac{\sum_{i=1}^n [(CO2_{ei} + CO2_{othi})_{bas} - (CO2_{ei} + CO2_{othi})_{cur}]}{\sum_{i=1}^n (CO2_{ei} + CO2_{othi})_{bas}}$$

Where $CO2_{ei}$ is the total CO₂ emissions released by the data center's consumed energy, $CO2_{othi}$ is the total CO₂ emissions released by the energy produced from other resources, i is the time period, bas is baseline and cur is current.

2. **Carbon Usage Effectiveness (CUE)** : CUE is a source based sustainability metric to better manage the environmental aspects of data centers which indicates CO₂ footprint in the daily operations of data centers. For

data centers using the GRID power, CUE is defined as ratio of total CO₂ emission for energy consumption of the data center annually over data center IT equipment energy [247].

$$\begin{aligned}
 CUE &= \frac{\text{co}_2 \text{ emission caused by}}{\text{the Total data center energy}} \\
 &= \frac{\text{Annual kWh IT load}}{\text{Annual kWh IT load}} \\
 &= CEF \times PUE \quad \left[\frac{KgCO_2}{kWh} \right]
 \end{aligned} \tag{B.1}$$

where CEF is carbon emission factor of the site. Unlike PUE, CUE is measured in $KgCO_2/kWh$. If a data center using 100% Green Energy, it will have a CUE of zero.

3. **Electronics Disposal Efficiency (EDE)** : EDE measures to what extent disposal of discarded IT electronics and electrical equipment (IT EEE) is efficient and evaluates progress in disposal of e-waste over time. This metric helps to know how well an organization responsibly manage e-waste that reaches end of current use (EOCU) or end of life (EOL) [249].

$$EDE = \frac{\text{Total weight}_{\text{responsibly disposed}}}{\text{Total weight}_{\text{all decommissioned}}}$$

where responsible disposed(numerator) is how IT EEE at its EOCU or EOL well managed through certified entities and all decommissioned(denominator) is the sum of weights of whole system/component reused, recycled, and wasted.

4. **Energy Reuse Effectiveness (ERE)** : ERE measures the reused energy outside the data center or in other locations of a site. If the energy is reused it reduces the energy data center need to buy [325].

$$ERE = \frac{\text{Total Energy} - \text{Reuse}}{\text{IT Equipment Energy}}$$

Total Energy is the sum of energy consumed by the cooling system, lost of energy during power distribution, energy utilized for data center lighting and energy used by all the IT equipment. ERE ranges from 0 to ∞ . If ERE reaches zero, then 100% of the energy is reused in other locations.

5. **Energy Reuse Factor (ERF)** : ERF is the ratio of energy reused in the site and total energy. ERF ranges from 0 (indicates no reuse) to 1 (indicates 100% reuse). ERF can be defined as follows [325] :

$$ERF = \frac{\text{Energy Reused in a data center}}{\text{Total Energy}}$$

we can relate ERF, ERE and PUE in the following way :

$$ERE = (1 - ERF) \times PUE$$

6. **Green Energy Coefficient (GEC)** : GEC quantifies the extent of renewable energy used in a data center. GEC is used to promote the use of Green Energy. GEC calculates the fraction of renewable energy used by the data center over the overall energy delivered to the data center [304].

$$GEC = \frac{\text{Renewable energy used in kWh}}{\text{Total power consumption in kWh}}$$

7. **Grid Utilization Factor (GUF)** : The Grid Utilization Factor (GUF) metric indicates the percentage of time that a data center has to use energy from the grid because the locally generated resources are not sufficient [256]. GUF is defined as:

$$GUF = \frac{\sum_{n=1}^N f(n)}{N}$$

$$f(n) = \begin{cases} 1, & \bar{n}e(n) < 0 \\ 0, & \bar{n}e(n) \geq 0 \end{cases}$$

Where $f(t)$ is a step function indicating if the renewable resources can cover the energy demand at time t , n is the measurement index, N is the number of measurements, $\bar{n}e(n)$ is the mean value of the net exported electricity at sampling period Δt between measurements n and $n - 1$.

8. **Material Recycling Ratio (MRR)** : MRR evaluates how much material is recycled or reclaimed or re-purposed producing a product or service [326].

It is defined as follows :

$$MRR = \frac{\text{total material : } (recycled + reclaimed + repurposed)}{\left(\begin{array}{c} Total : \\ inbound material \end{array} \right) - \left(\begin{array}{c} Total : outbound \\ product or service \end{array} \right)}$$

In the said equation, both numerator and denominator are measured in Mass (lbs/kg). 100% MRR indicates that there is no wastage. MRR (lifecycle), MRR (building), MRR (operations) and MRR (e-Waste) are used for detailed reporting of recycled materials.

9. **Water Usage Energy (ω)** : Water and energy are interconnected in either the way they are used or treated. Water Usage Energy captures the energy impact of water usage which is the ratio of total energy consumption of water treatment and distribution to the site over the power consumption of IT equipment [327]. It is defined as follows :

$$\omega = \frac{[(E_d + E_n)]}{E_{IT}} \times 10^3$$

where E_n, E_d are embedded energy in indirect, direct water usage respectively and E_{IT} is energy consumption of IT equipment. This metric helps us to manage water efficiency which impact the energy efficiency of the data center and also used to compare water efficiency across data centers. This metrics needs seasonal benchmarking to capture the effect of regional and seasonal water availability.

10. **The Green Index (TGI)** : Different benchmark suites uses different metrics. Choosing and Combining these performance outputs to formulate a single numbered green metric that stresses different components of a data center is a difficult task. TGI metric (Equation B.2 or B.3 or B.4) is one of this kind which evaluates the energy efficiency of servers from the various benchmark tests [328]. We can also use different weights such as time (t_i), energy (e_i) and power (p_i) consumed by each benchmark respectively.

$$TGI \text{ using } W_{t_i} = \frac{\sum_{i=0}^n w_{t_i} * EE_i}{EE_{Ref_i}} \quad (B.2)$$

$$\propto \frac{1}{p_i}$$

$$\begin{aligned}
 TGI \text{ using } W_{e_i} &= \frac{\sum_{i=0}^n w_{e_i} * EE_i}{EE_{Ref_i}} \\
 &\propto \frac{1}{e_i}
 \end{aligned} \tag{B.3}$$

$$\begin{aligned}
 TGI \text{ using } W_{p_i} &= \frac{\sum_{i=0}^n w_{p_i} * EE_i}{EE_{Ref_i}} \\
 &\propto \frac{1}{p_i}
 \end{aligned} \tag{B.4}$$

where EE_i is the ratio of Performance to Power Consumed, when computing i^{th} standard of a benchmark suite on supercomputers and EE_{Ref_i} is the ratio of Performance to Power Consumption when executing i^{th} standard of a reference system on a supercomputers.

11. **Technology Carbon Efficiency (TCE)** : TCE measures the carbon impact (in pounds) for every Kwh delivered to IT equipment. This carbon footprint depends on the energy sources used to produce electricity [329]. It is calculated as follows :

$$TCE = \frac{\text{Total Facility Power}}{\text{IT Equipment Power}} \times ECER$$

where ECER is Electricity Carbon Emission Rate. Less TCE indicates more green the data center is. TCE is used to compare overall environmental impact of facilities considering PUE rating and energy source. Placing 1 MW of additional IT load on a site having less TCE would save million pounds of CO₂ emission.

12. **Water Usage Effectiveness (WUE)** : WUE allows us to understand the effect water consumption has on the local electric grid. WUE at high level can be calculated as water usage of a data center per year over energy consumption of IT equipment (Kwh) [327]. WUE is used to optimize the water use of an operational data center. Water usage includes water used for on-site for power generation, cooling system and water evaporated in a data

B.4 Performance / Productivity Metrics

center and its sub systems. Based on the usage of water at site and source, The Green Grid consortium defined the following metrics [248] :

$$WUE = \frac{\text{Annual Water Usage}}{ITEE} \quad [\text{Liters/kWh}] \quad (\text{B.5})$$

$$\begin{aligned} WUE_{\text{source}} &= \frac{ASEWU + ASWU}{ITEE} \\ &= [EWIF \times PUE] + \frac{ASWU}{ITEE} \end{aligned} \quad (\text{B.6})$$

where ASEWU is Source Energy Water Usage per Year, ASWU is Site Water Usage per year, EWIF is energy water intensity factor and ITEE is IT Equipment Energy. Using WUE, we can determine ways to increase a data center efficiency and sustainability and can compare with similar data centers. By using WUE in conjunction with PUE and CUE metrics, an organization can reduce energy use.

B.4 Performance / Productivity Metrics

This section presents the metrics, which gives us infrastructure utilization, performance, and productivity of a data center.

1. **ACE Performance Score** : ACE is a score calculated using three conflicting indicators (Availability, Capacity, and Efficiency) to know the impact of design, decisions and configurations on data center performance [330]. ACE score measures the operational flexibility of a data center and how much a data center compromised compared to its design intent. It is based on the Computational Fluid Dynamics models of the site and integrations that can be made to data center in future.

ACE is very useful for taking effective decisions during operation to improve the performance to highest level. Even though it is a good metric we have some challenges to calculate ACE. We may not get the correct base scores without having a good model of the data center and integration to an existing model may be challenging without the right tools.

2. **CPU Usage** : CPU Usage measures to what extent the allocated CPU is utilized, defined as

$$CPU\ Usage = \frac{Used\ CPU}{Allocated\ CPU}$$

3. **Data Center Productivity (DCP)** : DCP is the ratio of the useful work that a data center produces over the quantity of any resource that it consumes to produce the work. DCP tallies the consumption of a data center-related resource, against the data center's output.

$$DCP = \frac{useful\ computing\ work}{total\ facility\ power}$$

But there is no uniquely agreed definition on "useful computing work", because each vendor tries to define it with their own preferences.

4. **Data Center Energy Efficiency and Productivity Index (DEEPI)** : DEEPI is the result of multiplying IT Productivity per Embedded Watt (IT-PEW) and Site Infrastructure Energy Efficiency Ratio (SI-EER) [296]. DEEPI indicates delivered IT Productivity per Watt of energy delivered to site infrastructure. This metric can be compared with other peer IT organizations. DEEPI helps us to find the improvements to be made and best practices that can be used.

$$DEEPI = IT-PEW \times SI-EER$$

where the details of IT-PEW and SI-EER are given as follows :

- (a) IT Productivity Per Embedded Watt (IT-PEW) measures the electricity consumption of the data centers in terms of productivity. IT-PEW is the composite metric which considers Data architecture, Reliability, Hardware specifications, Technology refresh and Archiving inactive data in different stages.

$$IT-PEW = \frac{IT_{productivity}}{Embedded\ Watt}$$

IT-PEW allows hardware manufacturer to compare their different model using benchmarks, also it allows to compare the manufacturers of different products.

- (b) Site Infrastructure Energy Efficiency Ratio (SI-EER) measures to what extent data center operators are managing the efficiency of site infrastructure systems. SI-EER is calculated as follows :

$$SI-EER = \frac{\textit{power-in}}{\textit{conditioned power-out}}$$

where Power-in is measured at the utility electric meter and power-out is the power used to run the IT equipment for computing. From the data of 85 clients it is observed that SI-EER of 2.5 is Acceptable, 2.0 is better and 1.6 is best achieved under some specific conditions [296].

5. **Dynamic Range (DR)** : This metric is commonly used in the literature as an approximation for energy proportionality. It is calculated as the ratio of difference between peak power and idle power to the peak power [331].

$$DR = \frac{(\textit{Peak Power} - \textit{Idle Power})}{\textit{Peak Power}} * 100$$

where Peak Power is the power consumption at 100% utilization and Idle Power is the power consumption at 0% utilization. An energy proportional computing systems will have DR of 100%. Drawback with DR is that it does not account the variations in power usage across different utilizations. DR is a poor measurement of the servers actual proportionality.

6. **Energy Proportionality (EP)** : EP quantifies a server's energy proportionality that derives a relationship between actual server energy consumption and ideal energy-proportional server [332]. This metric is used to quantify how closely a server's energy proportionality approaches perfect scaling when utilization varies from 0% to 100%. Figure B.1 shows an energy proportional system when EP=0.5, where "area A" is the region between the actual and the ideal power consumption curves. "area B" is the region under the ideal curve. Given the "area A" and "area B", EP is calculated as follows :

$$EP = 1 - \frac{\textit{area A}}{\textit{area B}}$$

An ideal energy proportional server will have EP=1 and power consumption will be proportional to its load where as the server power consumption is constant when EP=0 even with varying loads.

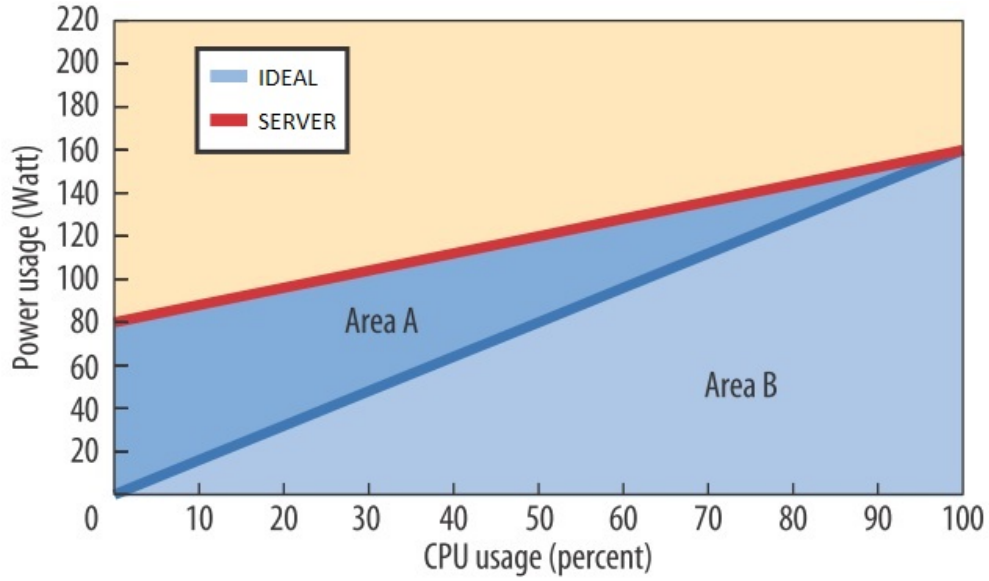


Figure B.1: Energy Proportional System, adopted from [332]

7. **FLOPS per Watt (FpW)** : Floating Point Operations Per Second per watt is the benchmark measurement for rating the server performance [333].

$$\begin{aligned}
 FLOPS/watt &= \frac{\text{Floating Point Operations / Second}}{\text{Joules / Second}} \\
 &= \text{Floating Point Operations / Joule}
 \end{aligned}
 \tag{B.7}$$

8. **Idle-to-peak Power Ratio (IPR)** : In order to measure how close to the origin, the power curve starts (i.e range), Varsamopoulos et al. [334] proposed the metric Idle-to-peak Power Ratio , which is defined as the ratio of the power consumption at 0% utilization (P_{idle}) over the power consumption at 100% utilization (P_{peak}). It is calculated as follows:

$$IPR = \frac{P_{idle}}{P_{peak}}$$

IPR is a normalized metric. So, it can be used for comparing the systems with different power consumption magnitude. This metric works well with the systems where the range of idle to peak power is high. IPR ranges from 0 to 1 where low value indicates more energy proportional system.

9. **Linear Deviation (LD)** : Linear Deviation is a measure of the energy proportionality curves linearity [331]. Linear Deviation is calculated as follows:

$$LD = \frac{Area_{actual}}{Area_{linear}} - 1$$

where $Area_{actual}$ and $Area_{linear}$ are the areas under the server's actual and linear energy proportionality curves, respectively. A server is considered linearly energy proportional if $LD = 0$, super-linearly energy proportional if $LD > 0$, and sub-linearly energy proportional if $LD < 0$.

10. **Linear Deviation Ratio (LDR)** : In order to measure how close the power curve to linear curve, we can use the metric linear deviation ratio (LDR) [334]. LDR is the maximum of the following ratio after absolute-value comparison.

$$LDR = \max \left| \frac{(P_{Actual} - P_{Ideal}) \text{ at } 0\% \text{ utilization}}{(P_{Actual} - P_{Ideal}) \text{ at } 100\% \text{ utilization}} \right|$$

where P_{Actual} , P_{Ideal} are the actual power consumption and hypothetical linear power consumption respectively. Detailed calculations of P_{Actual} and P_{Ideal} are given in [334]. Lower LDR values denote a more linear system. If $LDR < 0$, it denotes that the actual power curve is under the straight line. Positive LDR value denotes a power curve that is over the straight line. Further LDR is a normalized metric, so it can be used for direct comparison.

11. **Proportionality Gap (PG)** : This metric gives the deviation between the servers actual and the ideal energy proportionality at individual utilization levels [331]. This metric is useful when we need to know the disproportionality of servers in detailed granularity and to pinpoint the reason for dis-proportionality. PG at utilization level $u\%$ is given as follows:

$$PG_{x\%} = \frac{Power_{actual@x\%} - Power_{ideal@x\%}}{Power_{peak}}$$

PG is 0 for for an ideal energy proportional server $\forall x$. For both super-linear and sub-linearly proportional systems, PG is very large at 0% utilization of the server.

Metrics like PG, LD, EP etc. are useful when we want to improve energy proportionality in heterogeneous data centers [335].

12. **SWaP: Space, Watts and Performance** : SWaP gives performance of rack-optimized server deployments considering together space, energy, and performance [314]. It is used to compare the performance of different server deployments. The SWaP is calculated as follows:

$$SWaP = \frac{Performance}{Space \times Power Consumption}$$

where Performance is taken from industry-standard benchmarks, Space is the height of the server in rack units, and Power consumption is watts consumed by the system. The higher your SWaP numbers, the less data center space and power need to do a computing job. It is used by data center operators to choose among different servers which deliver the optimum performance for their needs and works in both power constrained and unconstrained situations.

13. **Data Center Utilization (U_{DC})** U_{DC} : calculates the proportion between IT equipment power consumption relative to the total capacity of the data center. U_{DC} of 91% demonstrates the most extreme achievable as believed by the industry. 50% and 45% are the peak and average utilization of a data center during the period of monitoring respectively [336].
14. **Server Utilization (U_{server})** U_{server} : measures utilization of the processor in comparison to its maximum ability. Maximum ability is the performance of the processor in the highest frequency state. U_{server} of 80% indicates the most extreme achievable that the industry believes in. 45% is the peak and 15% is the average utilization of servers during the period of monitoring [336].
15. **Uninterruptible Power Supply (UPS) losses** : Efficient UPS will deliver all the power received without any loss. But UPS itself will consume some power and there are some losses associated with UPS namely "square-law", "proportional" and "no-load" losses. The only means of comparing the efficiencies of UPSs is to evaluate these losses across all load levels [337].

No-load losses - If there is no load on UPS, then total input will be used by the UPS. This is called as no-load losses which is also known as tare, parallel and constant. These losses occur in powering transformers, logic circuits,

network cards etc are independent of load. This accounts more than 40% of all UPS losses. UPS efficiency can be much improved by decreasing these losses.

Proportional losses - With increasing load switching losses and the conduction losses will vary as huge amount of power will be processed throughout supply chain. All these varying losses in power path will put them to Proportional losses.

Square-law losses Electrical current flowing through the components of UPS increases with the load. This produces " I^2R " losses as the squandered power in the form of heat are proportional to the square of the electrical current. Square-law losses contribute up to 4% at higher loads [338].

16. **Uninterruptible Power Supply-Crest Factor (UCF):** CF calculates the proportion between the quick crest current required by the load (X_{PEAK}) and the Root Mean Square current (X_{RMS}) [339]. The crest factor of 1.4 is common in most IT and electrical equipment. A UPS must be sized properly to serve the peak load otherwise output will be distorted. Modern data centers with the power factor corrected components will eliminate the crest factor problem.

$$Crest\ factor = \frac{X_{PEAK}}{X_{RMS}}$$

17. **Uninterruptible Power Supply-Power Factor (UPF) :** Power factor is used to calculate the actual power in AC systems. It is the ratio of the actual power to the apparent power [339].

$$Power\ Factor = \frac{Actual\ Power}{Apparent\ Power}$$

Apparent power includes AC current without delivering energy. Hence it is larger than the actual power.

18. **Uninterruptible Power Supply-Power Factor Corrected (UPFC) :** PFC represents the ability of a non-linear reactive UPS to improve its power factor and reduce reactive power [340]. A power supply coming with "0.86" PFC indicates that UPS has to be given with 700VA to deliver 600 watts.

But it is better to choose a power supply with over 0.9 power factor corrected. Components which are power factor corrected will eliminate the crest factor problem.

19. **UPS energy efficiency** : UPS efficiency is the ratio of the UPS output power over the input power. UPS energy efficiency is calculated in KW. It is defined as follows :

$$UPS\ efficiency = \frac{UPS\ Output\ Power(kW)}{UPS\ Input\ Power(kW)} \times 100$$

The UPS efficiency varies depending on its load factor.

20. **Uninterruptible Power Supply-Surge Factor (USF)** : With a switch-mode power supply, it needs extra power to start data center components such as compressors of chillers, motors and some times disk drives. Surge factor indicates the ability of the UPS to handle these loads. This situation exists for a few seconds or more for special loads. For large disk arrays, it is desirable to have surge factor of 1.5 [339].

B.5 Thermal and Air Management Metrics

Cooling has been the major issue consuming nearly one third of the data centers energy consumption. High performance computing servers are bringing in new thermal and power challenges for data center operators. Data center operators must ensure minimum amount of energy for cooling which can be done through efficient air movement and using free cooling. In this section we present metrics for understanding the air flow in data centers, environmental conditions such as temperature, humidity, heat etc.

1. **Airflow Efficiency** : Airflow efficiency is measured in terms of power required to move the air inside a data center. This assesses overall efficiency of moving air gently through out the data center from cooling units to vents and takes into account facility layout and fan efficiency, measured in W/cfm

$$Air\ flow\ Efficiency = \frac{Total\ fan\ power}{Total\ fan\ airflow}$$

B.5 Thermal and Air Management Metrics

where the total fan power includes supply and return and total fan air flow includes supply and exhaust. The standard values of airflow efficiency are given in Table B.2 [321].

Standard	Better	Good
0.6 W/cfm	0.3 W/cfm	0.1 W/cfm

Table B.2: Standard Values of Airflow Efficiency

- Capture Index (CI) :** Unlike temperature-based metrics, Capture Index (CI) is solely a function of airflow. CI is defined in terms of local cooling resources (airflow from tiles, local coolers, local extract vents, etc.) and it measures the robustness and scalability of a grouping of equipment [341]. This is of two kinds: Cold & Hot. “Cold Aisle CI” is the fraction of cool air delivered to the rack which originates from CRAC units or perforated floor tiles. “Hot Aisle CI” assesses the fraction of hot air scavenged by coolers or return vents from the racks. CI assesses the efficiency of supplying cool air or capturing hot air to/from the rack on a 0% (bad) to 100% (good) scale.
- Data Center Temperature (DC) :** Temperature raise in a data center above the threshold has a negative impact on the IT equipment causing for reduction in reliability, lifetime of components. Also operating IT equipment at high temperatures ($> 30^{\circ}\text{C}$) for long time may cause unplanned downtime [342]. ASHRAE thermal recommendations for class 1 data center operations are given in Table B.3 [262].

Recommended	Allowable
18 – 27 ⁰ C	15 – 32 ⁰ C

Table B.3: AHSRAE thermal recommendations for Class 1 data centers

- Dew Point (DP) :** The dew point is a temperature, at which water vapor in the air condensates into liquid water. It is a measurement of the true amount of moisture in the air stream. For a data center, server room or communications room, this means that it is the temperature at which water

droplets inside your equipment may start to appear. In a sensible mode of heat transfer (processes staying above the dew point temperature) air temperature increases, but the amount of moisture in the air remains the same therefore dew point temperature remains the same. ASHRAE recommended 17°C maximum Dew Point in tightly controlled environment for data center operations [262].

5. **Heat Flux (HF) :** Heat Flux is most frequently used metric in a data center which includes area of layout and is measured in W/m^2 . It calculates the heat load to a given footprint area.

$$\text{Heat Flux} = \frac{\text{Total heat load}}{\text{The footprint area of the layout under consideration}}$$

By including volume of the facility we can extend the heat flux to a volumetric heat load (W/m^3) .

6. **Imbalance of Temperature (IoT) :**

The Imbalance of Temperature (IoT) allows for the evaluation of the quality of cooling [255]. It can be defined as the IoT of a node (server, blade), or the IoT of a rack (node-group).

$$IoT_{NG} = \frac{T_{CPU,max} - T_{CPU,min}}{T_{CPU,maxref}}$$

Where $T_{CPU,max}$ and $T_{CPU,min}$ are the maximum and minimum temperature reached by the CPUs in the node-group at a given time-stamp. And $T_{CPU,maxref}$ is the value of reference for the maximum acceptable temperature. The maximum acceptable temperature is selected as 100°C.

7. **Mahalanobis Generalized Distance (D^2) :** Mahalanobis generalized distance is a statistical metric. D^2 is a new way to characterize the non uniformity in rack heat load. Mahalanobis Generalized Distance is calculated by taking into account the distribution of a population in the Euclidean distance between two given points [343]. This metric has a large value when the heating non-uniformity covers a localized region and a small value when

B.5 Thermal and Air Management Metrics

it covers a broad region. In the data center context, D^2 can be calculated as follows :

$$D^2 = \begin{pmatrix} X - x_{ave} & Y - y_{ave} \end{pmatrix} S^{-1} \begin{pmatrix} X - x_{ave} \\ Y - y_{ave} \end{pmatrix}$$

where S is the covariance matrix of variances.(X,Y) and (x_{ave}, y_{ave}) are central coordinates of the CRAC units and racks respectively. High D^2 indicates the significant variations in heat loads across the rack, where a small D^2 indicates uniformity.

8. **CRAC Flow (M_c)** : M_c is the total air supplied by the cooling units in data center [344]. Generally, data centers have air flow rate more than what is needed by the servers, due to redundant CRAC units.
9. **Negative Pressure Flow (M_n)** : In a raised floor data center, if the under floor air velocity is high, then Venturi effect takes place. With this effect (Bernoulli law of fluid dynamics), air is drawn down to the floor void via grille [251]. Negative pressure flow is negligible in real time but is found near floor grilles near CRAC units, tile edges etc. This may cause insufficient air to meet the local cooling demand. This can be observed by placing a sheet of paper above the floor grilles to see whether it is drawn down or not.
10. **Bypass Air Flow (M_{bp})** : Not all the air produced by CRAC units reaches front intake. Some air returns directly to the CRAC units without passing through the IT equipment is Bypass air [344]. We can reduce Bypass air flow by sealing air gaps in the raised floor, cable cut-outs within cabinets. By relocating floor grilles so that they supply where it is needed we can reduce bypass air flow.
11. **Recirculation Air Flow (M_r)** : M_r is the air that is discharged from hot air aisle, which returns and mixes with air from cool air aisle that enters servers to cool them. We can use filler panels and internal air dams to prevent recirculation.

Apart from the above we can see the other components like Hall air flow(M_h), Floor air flow(M_f) and server air flow (M_s). Given the mass flow rates we can write the following mass equations.

$$M_f = M_n + M_c = M_b + M_h$$

B.5 Thermal and Air Management Metrics

$$M_s = M_h + M_r$$

Based on the said definitions and equations, the following ratios are defined [251] :

- **Negative pressure Ratio** is the ratio of Negative pressure flow (M_n) to the CRAC Flow (M_c).

$$NPR = \frac{M_n}{M_c}$$

- **Bypass Ratio** is the ratio of Bypass air flow (M_b) to the floor air flow (M_f).

$$BPR = \frac{M_b}{M_f}$$

- **Recirculation Ratio** is the ratio of recirculation flow (M_r) to the Server air flow (M_s).

$$RR = \frac{M_r}{M_s}$$

- **Balance Ratio** is the ratio of CRAC flow (M_c) to the Server air flow (M_s).

$$BR = \frac{M_c}{M_s}$$

12. **Degree-Days (DD)** : Degree-Days quantifies cooling and heating demands where a day's average temperature is above and below a base temperature which are "Cooling Degree Days (CDD)" and "Heating Degree Days (HDD)" respectively. Base temperature is generally 65°F or 18°C. Detailed definition and calculation of HDD and CDD can be found in [345].

13. **Relative Humidity** : Water vapor or moisture content in the air is known as humidity. Environment inside a data center plays vital role in improving the availability of IT equipment. Too high or too little humidity will reduce the performance and equipment downtime [282].

We can measure humidity using relative humidity calculated as the ratio of the amount of water in the air (absolute humidity) over the maximum vapor

B.5 Thermal and Air Management Metrics

(saturation humidity) that air can hold at given temperature. Amount of water in air is measured in grams per cubic meter (g/m^3).

$$Relative\ Humidity = \frac{actual\ vapor\ density}{saturation\ vapor\ density}$$

According to ASHRAE guidelines, Relative Humidity in the range of 20% to 80% is allowable and a Relative Humidity of 60% will give the best performance of servers [262].

14. **Rack Cooling Index (RCI) :** The Rack Cooling Index (RCI) is a cooling performance metric for analyzing the thermal environment in data centers and it is room health indicator. For continuous operation, servers and other electronic devices depend on intake temperature which RCI deals with. It is well suited as a design specification for new data centers. If the intake temperature exceeds the maximum recommended, we get “over-temperature” condition. If the intake temperature drop below the level of minimum recommended, then ”under-temperature” condition exists [346]. RCI has high and low limits: RCI_{Hi} (Equation B.8) and RCI_{Lo} (Equation B.9) based on temperature distribution along the rack height.

$$RCI_{Hi} = \left[\frac{Total\ Over-Temp}{Max\ Allowable\ Over-Temp} \right] * 100[\%] \quad (B.8)$$

$$RCI_{Lo} = \left[\frac{Total\ Under-Temp}{Max\ Allowable\ Under-Temp} \right] * 100[\%] \quad (B.9)$$

If both RCI_{Hi} , RCI_{Lo} are equal to 100% then ambient conditions are within the recommended range. But lower RCI_{Hi} or RCI_{Lo} shows more prominent likelihood that the data center equipment is encountering temperature above/below the allowable Maximum/Minimum respectively. Table B.4 gives ratings of RCI based on different analyses [347]. RCI is used to eval-

Poor	Acceptable	Good	Ideal
$\leq 90\%$	91% - 95%	$\geq 96\%$	100%

Table B.4: Compliance Of RCI

uate and report the effectiveness of cooling solutions if combined with CFD

B.5 Thermal and Air Management Metrics

modeling. It is used as a standardized way for specifying thermal quality which helps in marketing our cooling solutions.

15. **Return Heat Index (RHI) and Supply Heat Index (SHI)** : With inadequate air management systems and rack layouts that allow mixing of hot and cold air streams will increase the energy consumption of a data center. The level of separation of cold and hot air streams can be measured by the supply and return heat indices [348]. SHI is calculated as the ratio of sensible heat gained in the cold aisle to the heat gained at the rack (Equation B.10) which is a dimensionless measure. The return heat index (RHI) is defined as the ratio of heat extracted by the cooling system to the heat gained at the rack exit (Equation B.11).

$$SHI = \frac{\delta Q}{Q + \delta Q} \quad (B.10)$$

$$\begin{aligned} RHI &= \frac{Q}{Q + \delta Q} \\ &= \frac{\text{Total heat extraction by the CRAC Unit}}{\text{Total Enthalpy rise at the rack exhaust}} \end{aligned} \quad (B.11)$$

where Q is the total heat scavenged by the local heat extractors and δQ is heat gained by the air due to infiltration in cold aisle. SHI is function of rack inlet, outlet temperatures and CRAC outlet temperature. SHI varies between 0 to 1. Lower SHI indicates the better. For different layouts SHI and RHI values are shown in Table B.5 [348].

Layout	SHI	RHI
Room return	0.21	0.81
Ceiling return	0.2	0.83

Table B.5: SHI and RHI for different infrastructures

16. **Return Temperature Index (RTI)** : Return Temperature Index (RTI) is a measure of the energy performance of the equipment room air-management system [349]. The RTI is defined as follows :

$$RTI = \frac{RAT - SAT}{\Delta T_{Equip}} * 100[\%]$$

where RAT is return air temperature, SAT is supply air temperature and ΔT_{Equip} is increase across equipment. RTI below 100% indicates By-Pass air flow and RTI above 100% indicates Recirculation. By-pass air and Recirculation effects are detrimental to the thermal and energy performance.

The RTI is a measure of how well air management system controls by-pass or recirculation air. We can use RTI to evaluate the acuteness of energy penalty when RCI_{Hi} is improved. So the combined RCI and RTI metrics gives the overall performance of the cooling system [347].

17. **β -Index** : β -Index measures the extent of increase in inlet temperature due to recirculation [350]. Even with significant local hot spots, SHI and RHI can show the favorable values. To handle this local inefficiency, we can use β -index which is the ratio of the temperature differentials as given by

$$\beta = \Delta T_{inlet} / \Delta T_{rack}$$

where ΔT_{inlet} is the temperature gained by cooled air while flowing from chilled-air entry to rack inlet air, and ΔT_{rack} is the temperature differential between rack outlet and inlet air.

The value of β varies from rack to rack and also for different locations in front of the rack. In general, we consider average value for β . $\beta = 0$ indicates that there is no effect of recirculated hot air. $\beta = 1$ indicates that T_{inlet} is equal to Average rack outlet temperature. There is a possibility of β greater than 1 which indicates the existence of a local self-heating loop inside a rack.

B.6 Network Metrics

This section presents the metrics used to monitor efficiency of communication network of a data center. A Network topology can be described as an undirected graph G . This graph G contains a set of vertices (V) represent the nodes and set of edges (E) represent the links. The diameter of G , given by D_{max} . Let G' is the sub graph after removing different elements from the network. The diameter of the largest connected component of G' is D'_{max} . This information is used to calculate

path stretch and diameter stretch. Along with network efficiency metrics, these metrics are discussed as follows :

1. **Bits per Joule Capacity (BJC):** This metric is used to analyze the performance of energy-limited wireless sensor made up of tiny nodes and ad hoc networks. BJC is the maximum network delivery efficiency per Joule of energy in terms number of bits [351]. Under a fixed network size and the given error constraint, for each source destination pair in a network, only finite number of bits can be delivered.

Let $\epsilon \in (0, 1]$ and $f : R_+^{N \times N} \rightarrow R_+$. The bits-per-Joule capacity of K_G is defined as follows :

$$C_J(K_G; \epsilon, f) \stackrel{\text{def}}{=} \sup \{b/ b \text{ is } (\epsilon, f)\text{-achievable}\}$$

where the supremum is the least upper bound taken over the set of families of encoders, decoders and schedules on G . The definitions of the terms K_G , $R_+^{N \times N}$, R_+ and (ϵ, f) -achievable are in [351].

2. **Communication Network Energy Efficiency (CNEE):** CNEE measures the efficiency of a packet delivery process in a data center network. CNEE measures the energy spent for the successful message delivery by the network to the computing servers. CNEE is measured in joules/bit [254].

$$CNEE = \frac{\text{Power Consumed by Network Equipment}}{\text{Effective Network Throughput Capacity}}$$

3. **Diameter Stretch (DS) :**Diameter Stretch is the ratio of diameter of the largest connected component of G' over the diameter of the original graph G [352].

$$DS = \frac{D'_{max}}{D_{max}}$$

4. **Energy Consumption Rating Variable Load (ECR-VL) :** Technology level of a network system can be measured by normalizing its energy consumption to the highest sustained throughput recorded. But these systems in the field exhibits short term bursts. For this reason, we use variable

load metric ECR-VL which gives the network efficiency and differentiates the energy efficiency under various loads. It is measured in Watts per Gbps.

$$ECR-VL = \frac{w.E_{100} + x.E_{50} + y.E_{30} + z.E_{10} + \alpha.E_i}{w.TP_f + x.TP_{50} + y.TP_{30} + z.TP_{10}}$$

where TP_f is maximum throughput (Gbps) achieved in the measurement cycle, $TP_{50} = TP_f * 0.5$, $TP_{30} = TP_f * 0.3$, $TP_{10} = TP_f * 0.1$ and w,x,y,z,α are weight coefficients such that $(w+x+y+z+\alpha) = 1$. Other details are given in [353].

5. **Network Power Usage Effectiveness (NPUE) :** NPUE is the ratio of the total power consumed by the IT equipment (P_{IT}) over power consumed by network equipment ($P_{Network}$). NPUE measures the power used to operate data center communication system which is the part of IT equipment [254].

$$NPUE = \frac{P_{IT}}{P_{Network}}$$

6. **Network Traffic per Kilowatt-Hour :** Network traffic (bits) per kilowatt-hour is calculated as the ratio of outbound bits over data center energy. Information can be easily obtained in a data center to calculate this metric and is correlated to work. But it depends on the data center type and does not focus on useful work done [354].

$$Network\ traffic/KwH = \frac{Outbound\ bits}{data\ center\ energy}$$

7. **Path Stretch (PS) :** This metric quantifies the increase on the average path length. It is the ratio between all server pairs after removing vertices or edges from G (L'_{avg}) in relation to the average path length on G (L_{avg}) [352].

$$PS = \frac{L'_{avg}}{L_{avg}}$$

Most routing algorithms used in data centers considers shortest path between each pair of servers. This metrics also consider the same to evaluate the path quality on the network.

8. **Maximum Relative Size (RS_{max})** : This metric gives the reliability in complex networks. RS_{max} is the ratio of relative size of the largest connected component over the sum of existent servers in each connected component [352].

$$RS_{max} = \frac{\max_{1 \leq j \leq n} |S_j|}{\sum_{j=1}^n |S_j|}$$

where n is the total number of connected components in the resulting graph G' and $|S_j|$ gives the count of servers in each connected component j .

9. **Telecommunications Energy Efficiency Ratio (TEER)** : TEER assesses the energy efficiency of individual equipment and network configurations calculated as the ratio of useful work over power. TEER specifies equipment classes and is used for benchmarking similar equipment. TEER is defined as follows [355] :

$$TEER = \frac{\text{Useful work}}{\text{Power}}$$

where Useful work is based on the equipment function and Power is dependent on the equipment measurement.

10. **Network Utilization ($U_{network}$)** : $U_{network}$ indicates the "network utilization" which is the percentage of bandwidth used relative to the bandwidth capacity in the data center [336]. The maximum achievable value that industry believes for $U_{network}$ is 80%. Peak and average network utilization are 30% and 10% respectively.

B.7 Storage Metrics

It is important to monitor and notify the measurements that boost efficiency to meet storage requirements of a data center. In storage perspective, data center operators need to get the information on metrics such as throughput, availability etc, which are discussed as follows :

1. **Capacity** : Capacity measures the energy consumed within the storage facility and it is computed as follows [356] :

$$\text{Capacity} = \text{Capacity storage} / \text{Watt}$$

The Capacity Metric (GB/Watt) represents the energy efficiency of storing the data of the user's applications. The metric is defined as the ratio of space used by files written and stored on the storage system (GB) to average power of the storage system under typical usage over a period of time (Watts).

2. **Low-cost Storage Percentage (LSP)** : This metric depends on the criticality of the data and presents the fraction of data that is stored on the lowest-cost and high-capacity drives. If the data stored is less critical then we can move the same data to a lower cost storage without disturbing business operations [269].
3. **Memory Usage** Memory Usage refers to the utilization of the main memory computed as follows :

$$\text{Memory Usage} = \frac{\text{Used Memory by a server/application}}{\text{Allocated Main Memory}}$$

4. **Overall Storage Efficiency (OSE)** : OSE is the ratio of customer stored data compared to the raw storage capacity. This is used to understand to what extent the raw storage capacity utilized. But measuring customer stored data is difficult due to data duplication and the user's view differ from the storage frame view [269].
5. **Response Time (RT)** : Response time describes the time to complete a single read or write operation, measured in milliseconds (fraction of a millisecond for flash storage). Ideal latency value would be zero, indicating that the application will be served without any delay for read/write operations.
6. **Slot Utilization (SU)** : This metric measures the efficiency of storage as the ratio of storage frame slots that are filled with hard drives over the total available storage frame slots. Efficient utilization of frame slot indicates the minimization of storage cost of data center.
7. **Throughput** : Throughput indicates energy efficiency of I/O operations (i. e. , data read and write). Throughput measures of speed at which the storage system delivers data [357]. Throughput is measured in two ways: I/O rate and data rate measured in accesses/second and bytes/second respectively. I/O throughput computed considering the number of operations

as:

I/O throughput = I/O operations per second/Watt

Data Transfer throughput (D) considers the amount of data involved computed as follows :

D = Mega Bytes moved per second/Watt

The I/O throughput is used for applications like transaction processing where data transfer is small. Data transfer is used for applications where the amount of each request will be huge.

8. **Storage Usage ($U_{storage}$)** : This measures the percentage of storage used relative to the overall storage capacity within the data center. $U_{storage}$ of 70% is believed as the maximum achievable, 40% is peak and 35% is average storage utilization [336].

B.8 Security Metrics

Security is one of the major concerns for business operations. As the data center houses the owner's core assets and clients data, they must be safe guard against physical as well as software threats. A firewall must have the capacity to handle the quickly advancing, network intensive service environment of the data center. Quality of the firewall policy is the key to any cyber defense perimeters. For security metrics, number of blocked attacks or intrusions detected will become the logical starting points. This section describes the physical and IT security aspects of the data centers and some basic measurements of complexity and performance of firewalls, intrusion detection and prevention system.

1. **Average Comparisons Per Rule (ACPR)**: ACPR is an important metric that measures the performance of the firewall. It measures the average number of comparisons required to match a rule in the firewall policy [358]. If comparisons are more it will affect the performance of a firewall. So, performance of firewalls depend on the order of policy rules. It is suggested to move high frequent matched rules to the top that makes firewalls to trigger these rules very fast. ACPR will be minimum if frequent rules appear early

in the policy else indicates the network administrators to reorder the rules.

$$ACPR(F) = \sum_{i=1}^n i * f_i$$

where f_i is the hit count for rule r_i and n indicates the number of rules.

2. **Accessibility Surface (AS)** : This metric is used to quantify the firewall rule / policy. AS of a firewall is the sum of accessibility surface areas of all interfaces (l) in a firewall (F) [358].

$$AS(F) = \sum_{i=1}^l IAS(F_i)$$

Where IAS is accessibility surface of an interface described in 12. Higher value of AS indicates that the policy is permissive thus policy need to be refined or analyze the traffic using intrusion detection systems.

3. **Application Transaction Rate (ATR)** : Firewalls must peer deep into the application layer to secure the traffic to detect the attacks that move from the network to application layer. Capability of the firewall to secure discrete application-layer transactions and gateways contained in an open connection is known as Application Transaction Rate [359].
4. **Concurrent Connections (CC)** : Concurrent connections measures the firewall ability to handle the growing information processing capabilities. CC measures the maximum number of open connections which are point-to-point and through firewall device. This number reflects maximum information points firewall can support [358].
5. **Connection Establishment Rate (CER)** : For a TCP/IP session, dozens of connection will be established across the organization's firewall. CER measures the speed of firewalls to establish connections and the full three-way handshake for a TCP/IP session [360].
6. **Connection Tear down Rate (CTR)** : CTR represents the rate at which firewalls can obliterate the connections and free resources to be utilized for other activity [360].

7. **Defense Depth (DeD)** : It is the minimum number of independent single machine compromises required for a successful network attack. It is useful for identifying the firewall policy or system configuration which can be defeated by a single point of failure [359].
8. **Detection Performance (DeP)** : DeP is used to measure the effectiveness of the detection mechanisms employed in an organization. The metric is defined as detection probability multiplied by the number of true alarms during tests [361].

$$Detection\ Performance = P_d * (1 - P_{fa})$$

where P_d is Probability of detection in test cases and P_{fa} is the probability of false alarm.

9. **Data Transmission Exposure (DTE)** : DTE indicates the unencrypted data transmission volume [361]. This metric counts the unencrypted communication channel pairs and TCP-port pairs in use.
10. **Firewall complexity (FC)** : Firewalls serve as the first line of defence against threats. However, the protection depends on how good the firewalls configured. Complexity of a firewall policy F, is defined as the average of complexity of F for each destination address [358].

$$FC = \frac{\sum_{i=1}^n Complexity(D_i, F)}{total\ number\ of\ destinations}$$

where Complexity (D_i, F) ranges from 0 to 1 and gives the complexity of a policy to a particular destination D_i defined as follows ://

$$Complexity(D, F) = \frac{\begin{array}{l} positive\ rules \\ allowing\ traffic\ to\ D_i \end{array}}{\begin{array}{l} total\ rules\ for \\ controlling\ traffic\ to\ D_i \end{array}}$$

If complexity is 0, then all the traffic to destination D_i is blocked. If the value is 1 means that all traffic to D_i is allowed. This metric indicates the extent of intra-policy conflicts.

11. **HTTP Transfer Rate** : This measures the transaction rate that the firewall handles per unit of time. An entire HTTP transaction includes connection request, transferring objects, and closing connections [360].
12. **Interface Accessibility Surface (IAS)** : For a given firewall F, if the policy size for an interface (p) is 'n' then accessibility surface is sum of the areas covered by each single rule defined for that interface [358].

$$IAS(F_p) = \sum_{i=1}^n RA(i)$$

where RA(i) is Rule Area for rule_i given in 17.

13. **IP Fragmentation Handling (IPFH)** : IPFH is the ability of firewall to bring together the fragments before any rule is applied [360].
14. **IP Throughput** : It is the ability of the firewall to process bits from interface to interface [362]. It accounts the amount of bits or packets processes per second from one interface to other.
15. **Illegal Traffic Handling (ITH)** : ITH accounts the ability of firewall to concurrently handle both legal and illegal traffic. Illegal traffic can be either dropped or denied and is explicitly specified in rule [360].
16. **Latency** : Latency is the delay time of network traffic in the firewall [359]. This is accountable while calculating total delay of the network traffic. Latency is measured in milliseconds under steady state load near the firewall's limit.
17. **Rule Area (RA)** : Area of rule 'r' is the count of all source/destination addresses combinations where redundant combination are counted as one [358]. RA can be calculated if and only if we know the trustfulness of source and damage impact of destination. Trust and impact are the values between 0 and 1 which gives the trustfulness of an origin and the damage effect under attack respectively.

$$RA(r) = \sum_{i=1}^s \sum_{j=1}^d T_i P_j$$

where s and d are the number of source and destination addresses covered by rule r respectively. T_i is the trust of each source and P_j is the impact for each destination address.

18. **Reachability Count (RC)** : RC is the number of access points from a specific origin [361]. Decreasing the reachability count reduces the security hazards. This metric requires complete network configuration and ports access information and defined as follows:

$$\text{Reachability count} = N_s + N_o + N_p$$

where N_s is the number of ports that reply to traffic from a source, N_o is the Number of machines that have two-way connection-oriented sessions to the point of origin and N_p is the number of paths that have physical access to secure parts like storage media drives. The provenance for physical access may be outside the fence or from a partially controlled area inside the data center.

19. **Rogue Change Days (RCD)** : A rogue change is a system configuration change that is not informed to security experts. Rogue change days reflects the number of days that these unspecified changes are unknown to the security experts [361]. It is defined as the product of total rogue changes and number of days these changes are not identified by security experts. This metric is used to know the security impacting changes but does not give the impact of these changes.
20. **Vulnerability Exposure (T)** : This metric gives the number of days vulnerabilities are open. It is the sum of exposure time interval of each known and unpatched vulnerabilities until they are discovered locally or by public [361]. High vulnerability days indicates the greater risk. Total vulnerability days are defined as follows :

$$T = \sum_{i=1}^{nv} (t - T_i)$$

where t is the current date, nv is total known and unpatched vulnerabilities, and T_i is the date i^{th} vulnerability is discovered.

B.9 Financial Impact Metrics

This section presents the financial impact metrics of a data center. We explore various metrics which accounts cost associated with designing and operation of data center, financial impact of data center outage and return on investments on management tools and technologies for sustainable data center.

1. **Business Value of Converged Infrastructure (BVCI) :** Converged infrastructure (CI) is an approach to data center management that relies on a specific vendor and the vendor's partners to provide pre-configured bundles of hardware and software in a single chassis. As the scope of virtual server deployment expands asset utilization, employee productivity increases and reduces the capital costs. A recent research on 22 companies indicates that substantial business benefits are associated with higher convergence and asset sharing [363]. We can observe the effect of converged infrastructure through measures like IT cost per unit of workload, faster deployment, MTBF, MTTR, asset types in maintenance and average time to provision server and storage.

2. **Capital Expenditure and Operational Expenditure of data center :**

Capital Expenditure(CapEx) are funds used by an organization to acquire or upgrade physical and non-consumable assets which will be depreciated over time. To be simple, CapEx are single payments in exchange for goods or services. For a data center CapEx include the purchase of data center servers, land, buildings, cooling equipment, network infrastructure, and software. All these assets will be depreciated over a number of years in the data center. CapEx is often used by the organization to undertake new projects or investments with the intent of substantial return on investment [364].

Operational Expenditure (OpEx) refers to the day-to-day costs of operation and are recurring. They represent the cost of keeping the company operational and include costs of technical and commercial operations, administration, etc. Data center operating costs include a range of utilities including electricity and water needed just for the physical data center in-

frastructure, rentals for leased infrastructure, personnel wages, contractors, peripheral licenses and taxes that apply to the data center. CapEx and OpEx are interconnected issues. For example, a data center infrastructure management technology that enables automated maintenance and provisioning tasks may have higher acquisition costs but will be cheaper to operate [364].

3. **Carbon Credit (CCr)** : Carbon Credit measures the extent to which an organization emitting less CO_2 . Carbon Credit is a tradable greenhouse gas emission reduction unit. This gives us the offset credits that an organization holds which can be sold or to be bought to offset CO_2 emissions [356]. Rules and conditions to measure this metric vary based on the national regulations [365].
4. **Data Center Outage** : Recent studies have shown that the average data center downtime is 1 to 2 hours, which is slightly more than \$7,900 per minute [366]. Following are few terms associated with the data center outage :

- **Availability (A)**: It is the probability that a data center will be operating at certain time. It is the function of reliability and maintainability. The system availability is the ratio of data center operating time over the total time [367]:

$$A = \frac{MTBF}{MTBF + MTTR} \quad (B.12)$$

Where MTBF and MTTR are described below.

- **Mean Time Between Failures (MTBF)** : MTBF represents the average exposure time between consecutive failures or outages of a data center. The MTBF is usually expressed as years per failure [368]. MTBF can be determined as

$$MTBF = \frac{\text{total surviving hours}}{\text{number of failures or outages}}$$

- Mean Time To Failure (MTTF) : MTTF gives the expected time to failure for a non-repairable component. MTTF is the length of time hardware or other devices can perform reasonably. MTTF is used to evaluate the reliability of device or system [367].
- Mean Time To Repair (MTTR) : In case of failure MTTR represents the average time it takes to repair the component or fix the problem. For a system, it is the sum of average time to fix the failure by repair staff and the average time a machine spends in the queue waiting to be repaired. It is expressed in hours [367].
- Reliability (λ): Reliability is a measure of the ability of equipment performing satisfactory in the specified environment and operational loads, during a specified time [367]. We can calculate reliability using MTTR and MTTF as follows :

$$\lambda = \frac{1}{MTBF} \quad \mu = \frac{1}{MTTR}$$

λ gives Component failure rate and measured as faults per hour. μ gives component repair rate and measured as number of repairs per hour.

5. **Return On Investment (ROI) :** ROI is a metric typically used to compare profitability and the efficiency of different investments. ROI measures the amount of returns relative to the investments cost. If an enterprise has immediate objectives then ROI can be measured in different aspect of meeting one or more of these objectives not looking into cost savings or profits [369]. ROI is expressed as percentage of gain or loss defined as follows:

$$ROI = \frac{\text{Net profit}}{\text{cost of investment}} \times 100$$

A high ROI implies that gains are up to the stand and on par with investment cost. This is the most used profitability ratio metric because of its flexibility. But ROI calculation can be manipulated, so results may vary between users and organizations. Further, basic ROI calculation does not take time into consideration.

Data centers are expensive to build and maintain. Apart from sales and marketing a service, improving energy efficiency and reliability is the key to maximize the return on investment in data center technologies. With best data center infrastructure management solutions, operators can have high ROI in terms of improved energy efficiency, productivity, availability, and manageability.

6. **Total Cost of Ownership (TCO)** : TCO represents the sum of capital expenditure and the cost over time to operations and maintain the data center (operating expenditure). A TCO metric is cost/sever when the specifications and number of servers are known. Also we have cost/kW when the details of the servers to be installed are unknown. Here kW is the power available to the servers not the power into the site. The use of cost/sq. ft would not be valid as the high-density data center will have a greater cost/sq. ft. TCO can be expressed in a per-rack basis which normalizes the measurement of TCO but it requires a significant amount of data including capital, engineering, installation and operating cost data of various elements [370].

Appendix C

Data Center Assessment Checklist

We have compiled a handy check list of some key questions and metrics to run through while evaluating the data center as part of Chapter 6. It is not a definitive list, but it covers the main aspects that operators need to think about. This checklist is a living document of recommended actions to increase energy efficiency in data centers. Designed for data center owners and operators, this appendix provides actionable guidance to both prioritize and implement energy saving measures in data centers. More specifically, individuals can use relevant actions into an action plan or into the recommendations section of an energy assessment report. We divided this list into eight sections that represent data center subsystems and other areas that deserve attention:

- Energy Efficiency (EE)
- Thermal and Air Management (TA)
- Cooling Plant (CP)
- Overall Performance and Distribution Chain (OPD)
- Greenness (G)
- Network (N)
- Storage (S)

- Security (S)

Each section starts with a set of questions that captures the present status of the data center and measures the performance presented in Table C.1 to Table C.8.

Notes:

1. We assess the different features of a data center and various metrics to measure the performance of the data center along with the data required for each metric.
2. There are three priority levels for metrics: 1 indicates that the metric is very important and must be measured; 2 indicates that it is an important metric and must be measured if data is available; 3 indicates that it need to be measured, only if easily available.

Table C.2: Key Questions and Metrics for Thermal and Air Management

2. Thermal and Air Management						Notes
TA.Q1	What is the temperature set point of the cooling system ?					
TA.Q2	What is a typical return temperature?					
TA.Q3	What is a typical (average) supply temperature?					
TA.Q4	What is the IT equipment intake temperature on an average?					
TA.Q5	Does the data center need humidity control?					
TA.Q6	Do you have automatic humidification controls?					
TA.Q7	What type of humidifier do you have?					
TA.Q8	What is the prevalent humidification set point? (% RH)					
TA.Q9	Recommended and allowed IT equipment intake temperatures?					
TA.Q10	Recommended and allowed IT equipment intake humidity?					
TA.Q11	Whether the air temperature humidity sensors installed ?					
TA.Q12	Do you have the centralized CRAC/CRAH units ?					
TA.Q13	For all the cooling equipment report nameplate data to identify capacity and design conditions.					
TA.Q14	How many CRAC/CRAH/AHUs are there that are standby units?					
TA.Q15	Is there any supplemental cooling?					
TA.Q16	Do you use water side economizer or air economizer ?					
	Air Supply Path					
TA.Q17	Plenum height?					
TA.Q18	Plenum static pressure?					
TA.Q.19	What percent of cool air is wasted due to floor leakages?					
	Air Return Path					
TA.Q.20	Plenum height?					
TA.Q 21	Plenum static pressure?					
TA.Q.22	What percent of cool air is wasted due to floor leakages?					
TA.Q.23	Do you handle these leaks in raised floor and other places?					
TA.Q.24	Raised floor exist?					
TA.Q.25	Drop ceiling exist?					
TA.Q.26	Where are cables and pipes located?					
TA.Q.27	Are the cable penetrations sealed?					
TA.Q.28	Is the cable build-up in plenum more than 33% of the plenum height?					
TA.Q.29	Is there a cable-mining (allow proper pressure distribution) program in place?					
TA.Q.30	Degree to which hot and cold aisles are currently fully enclosed?					
TA.Q.31	Do you practice minimization of bypass air and recirculated air at the racks?					
TA.Q.32	Supply Air: Where are the overhead diffusers or perforated floor tiles placed?					

TA.Q.33	Location and configuration of the CRAC/CRAH units?								
TA.Q.34	Is there a fan speed control mechanism in place?								
TA.Q.35	Do some areas of the data center have very high load density?								
TA.Q.36	Is the HVAC system optimized to ensure correct airflow rates?								
Metric ID	Metric Name	Unit	Data Required	Priority	Value	Notes			
TA.M1	Return Heat Index	-	Total heat extraction by the CRAC Unit, Total Enthalpy rise at the rack exhaust	1					
TA.M2	Supply Heat Index	-	Enthalpy rise due to infiltration in cold aisle, Total Enthalpy rise at the rack exhaust	1					
TA.M3	Rack Cooling Index	-	Total Over-Temp, Max Allowable Over-Temp	2					
TA.M4	β – Index	-	Rack outlet air temperature, Rack inlet air temperature						
TA.M5	Average increase in Rack Temperature	F	Rack Intake and outtake temperatures	2					
TA.M6	Airflow Efficiency	W/cfm	CRAC Power / AHU Power, CRAC Airflow / AHU Airflow	2					
TA.M7	Return Temperature Index	-	Temperature increase across equipment, Supply and Return air temperatures	2					
TA.M8	System Pressure Drop	in w.g.	Return Side and Supply Side Pressure Drops	3					

Table C.4: Key Questions and Metrics for Overall Performance and Distribution Chain

5. Overall Performance and Distribution Chain							Notes
OPD.Q1	What is the current usage factor? (% of space?)						
OPD.Q2	What is the average age at which you replace your servers?						
OPD.Q3	Do you purchase or lease IT equipment?						
OPD.Q4	Do you have power management enabled on your servers?						
OPD.Q5	Do you use nameplate ratings when provisioning power for new IT equipment?						
OPD.Q6	Do you return old servers to the vendor who supplied them when they are replaced by new servers?						
OPD.Q7	What percentage of your servers are ≥ 5 years old?						
OPD.Q8	Do you measure the effectiveness of delivering a service ?						
OPD.Q9	What is the average utilization of servers?						
OPD.Q10	Is load balancing in place?						
OPD.Q11	UPS Technology Type						
OPD.Q12	load factor of active UPS module (average) ?						
OPD.Q13	UPS Redundancy Configuration ?						
OPD.Q14	UPS Input Power Factor ?						
OPD.Q15	Is there a standby generator?						
OPD.Q16	Standby generator power configuration						
OPD.Q17	Are there PDUs with built-in transformers?						
OPD.Q18	Types and constitution of MV and LV transformer(s) ?						
OPD.Q19	What is the load balance between the phases?						
OPD.Q20	Average Load Factor per Active PDUs / Transformers						
OPD.Q21	What is the lighting power density?						
OPD.Q22	Is there any automatic lights controlling mechanism in place?						
OPD.Q23	What type of lamps are used?						
OPD.Q24	What type of ballasts are used?						

Metric ID	Metric Name	Unit	Data Required	Priority	Value	Notes
OPD.M1	IT Rack Power Density	kW/rack	Peak Power of the IT Equipment, Number of Racks	2		
OPD.M2	ACE Performance Score	-	Availability, Capacity, and Efficiency	1		
OPD.M3	Data Center Energy Efficiency and Productivity Index	Productivity/ watt	IT productivity, power used for IT equipment	1		
OPD.M4	UPS Load Factor	-	UPS Input kW, UPS rating	1		
OPD.M5	UPS Power Energy Efficiency	%	UPS Input kW, UPS Output kW	1		
OPD.M6	PDU Efficiency	%	PDU Input and Output	1		
OPD.M7	IT Peak Power Density	W/sf	Peak Power of the IT Equipment, Floor Area of the data center	1		
OPD.M8	Data Center Lighting Density	kW/sq. ft	Power consumption for data center lighting, Total data center space	2		

Table C.6: Key Questions and Metrics for Network

7. Network				Notes		
N.Q1	Redundancy of core and edge network infrastructure					
N.Q2	Are you providing redundant private VLANs with high speed up to 1000 mbps					
N.Q3	on-site meet-me-room ?					
N.Q4	Dark fiber available ?					
N.Q5	network drop availability using HSRP or VRRP					
N.Q6	Do you have 24x7 support contracts with network gear vendors					
N.Q7	private data transport between data centers					
N.Q8	wireless service available					
N.Q9	How many direct connections to different network providers are available?					
N.Q10	100% Network uptime service level agreement					
N.Q11	Zero-tolerance spamming policy					
N.Q12	Gather application flow characteristics at the edge of the data center network					
N.Q13	BGP routing for customer owned IP addresses ?					
N.Q14	Identify the applications process name and network port number mapping ?					
N.Q15	What type of cross-connects are used?					
N.Q16	Carrier-neutral access to provider of your choice ?					
N.Q17	How do you eliminate downtime as a result of a fiber cut ?					
Metric ID	Metric Name	Unit	Data Required	Priority	Value	Notes
N.M1	Communication Network Energy Efficiency		Power Consumed by Network Equipment, Effective Network Throughput Capacity	1		
N.M2	Network Power Usage Effectiveness		IT equipment power consumption, Power consumed by network equipment	1		
N.M3	Network Utilization		Bandwidth used, Available bandwidth capacity	1		

Appendix D

Fact Sheet for Publications


Table D.1: Fact sheet for publication [1]

Title	Metrics for Sustainable Data Centers
Authors	V. Dinesh Reddy, Brian Setz, G. Subrahmanya V. R. K. Rao, G. R. Gangadharan, and M. Aiello
Publication	IEEE Transactions on Sustainable Computing , Vol. 2, No. 3, Pages 290 - 303, 2017
ISBN/ISSN	2377-3782
DOI	10.1109/TSUSC.2017.2701883
Status	Published
Publisher	IEEE
Publication Type	Journal

Metrics for Sustainable Data Centers

[View Document](#)

884
Full
Text Views

 Open Access

Related Articles

A wireless, passive carbon nanotube-based gas sensor

Characterization of micromachined spiked biopotential electrodes

[View All](#)

5 Author(s) [V. Dinesh Reddy](#) ; [Brian Setz](#) ; [G. Subrahmanya V. R. K. Rao](#) ; [G. R. Gangadharan](#) ; [Marco Aiello](#) [View All Authors](#)

[Abstract](#) [Authors](#) [Figures](#) [References](#) [Citations](#) [Keywords](#) [Metrics](#) [Media](#)

Abstract:

There are a multitude of metrics available to analyze individual key performance indicators of data centers. In order to predict growth or set effective goals, it is important to choose the correct metric and be aware of their expressivity and potential limitations. As cloud based services and the use of ICT infrastructure are growing globally, continuous monitoring and measuring of data center facilities are becoming essential to ensure effective and efficient operations. In this work, we explore the diverse metrics that are currently available to measure numerous data center infrastructure components. We propose a taxonomy of metrics based on core data center dimensions. Based on our observations, we argue for the design of new metrics considering factors such as age, location, and data center typology (e.g., co-location center), thus assisting in the strategic data center design and operations processes.

Published in: IEEE Transactions on Sustainable Computing (Volume: 2, Issue: 3, July-Sept. 1 2017)

Figure D.1: Publication [1]

Table D.2: Fact sheet for publication [2]

Title	Energy-aware virtual machine allocation and selection in cloud data centers
Authors	V. Dinesh Reddy, G. R. Gangadharan, and G. Subrahmanya V. R. K. Rao
Publication	Soft Computing, Pages 1-16, 2017
ISBN/ISSN	1432-7643
DOI	https://doi.org/10.1007/s00500-017-2905-z
Status	Published
Publisher	Springer
Publication Type	Journal

The screenshot shows the Springer Link interface for the article. At the top, there is a search bar and navigation links for Home, Contact us, and Login. The article title is prominently displayed, along with the journal name 'Soft Computing' and page range 'pp 1-16'. The authors' names are listed below the title. A 'Methodologies and Application' badge is visible, along with the first online date '06 November 2017' and a download count of 112. The abstract section is partially visible, starting with 'Data centers evolve constantly in size, complexity, and power consumption...'. On the right side, there is a sidebar with a 'Log in to check access' prompt, a 'Buy (PDF)' button for EUR 34.95, a list of benefits (unlimited access, instant download, local sales tax), a 'Subscribe to Journal' button, and a 'Cite article' dropdown menu.

Figure D.2: Publication [2]

Table D.3: Fact sheet for publication [3]

Title	Best Practices for Sustainable Data Centers
Authors	V. Dinesh Reddy, Brian Setz, G. Subrahmanya V. R. K. Rao, G. R. Gangadharan, and M. Aiello
Publication	IEEE IT Professional, In Press, 2017
ISBN/ISSN	1520-9202
DOI	
Status	Accepted (in Press)
Publisher	IEEE
Publication Type	Journal



Dinesh Reddy V <dineshvemula@gmail.com>

IT Professional, ITPro-2017-03-0029.R1 , Decision: Accept

1 message

IT Professional <onbehalfof+san+computer.org@manuscriptcentral.com>

Mon, May 8, 2017 at 11:43 AM

Reply-To: san@computer.org

To: dineshvemula@gmail.com, briansetz@gmail.com, subrahmanyavrk.rao@cognizant.com, geeyaar@gmail.com,

aiellom@ieee.org

Cc: itpro-ma@computer.org, rdeuelgallegos@gmail.com

IT Professional, ITPro-2017-03-0029.R1
manuscript type: General Interest
"Best Practices for Sustainable Data Centers"

08-May-2017

Dear Dr. G R Gangadharan,

Congratulations! Your above referenced manuscript has been officially accepted for publication in a future issue of IT Professional, subject to final editing for English style, clarity, and organization. The editor's comment below.

Before submitting your manuscript as outlined below, if you would like include the figure on Airflow management, which you said you couldn't include because of page length, you may include that figure in the appropriate place and upload the updated manuscript. It is optional, however.

Once the production cycle begins, the editor will assign your manuscript to a professional editor. The editor will contact you directly to discuss any recommended changes that will enhance the presentation of the article. The final editing of the manuscript will be a collaborative process in which you and the IT Professional staff work together to achieve a concise, well-worded article. Please note that IT Professional reserves the right to change the title of any paper that is accepted for publication.

Kind Regards,

Mr. Andy Morton on behalf of Prof. San Murugesan, EIC
IT Professional
itpro-ma@computer.org

Editor Comments:

The authors have successfully addressed the recommendations of the reviewers. A figure on airflow would have been helpful if the page limit could be stretched a bit.

My recommendation is to accept the revised paper.

Figure D.3: Publication [3]

Table D.4: Fact sheet for publication [4]

Title	Towards an Internet of Things framework for financial services sector
Authors	Vemula Dineshreddy, and G. R. Gangadharan
Publication	In Proceedings of the 3rd International Conference on Recent Advances in Information Technology (RAIT) RAIT, Dhanbad, India, Pages 177-181, 2016
ISBN/ISSN	978-1-4799-8579-1
DOI	10.1109/RAIT.2016.75078977
Status	Published
Publisher	IEEE
Publication Type	Conference

The screenshot displays the IEEE Xplore Digital Library interface. At the top, there are navigation links for 'Institutional Sign In', 'IEEE', 'Browse', 'My Settings', 'Get Help', and 'Subscribe'. A search bar is present with the text 'Enter keywords or short phrases (searches metadata only by default)'. Below the search bar is an advertisement for 'Need Full-Text access to IEEE Xplore for your organization?' with a 'REQUEST A FREE TRIAL >' button. The main content area shows the title 'Towards an "Internet of Things" framework for financial services sector' with a 'Sign In or Purchase to View Full Text' button and a '284 Full Text Views' indicator. The authors are listed as 'Vemula Dineshreddy ; G. R. Gangadharan'. Below the title, there are tabs for 'Abstract', 'Authors', 'Figures', 'References', 'Citations', 'Keywords', 'Metrics', and 'Media'. The 'Abstract' tab is selected, showing the following text: 'The ability to apply state-of-the-art Internet of Things (IoT) technology to extract customer insights through analytics by shaping the information into consumables for other connected systems is creating a lot of opportunities for banking and financial services. This paper presents an architecture based on Internet of Things for banking and finance sector by managing, mobile, household devices, wearable sensors and other sensing devices for various applications including retail banking, insurance, and investments. We have presented a case study of different banking applications flow with IoT-intelligence by analyzing users' data. In addition, we have a mapping of the proposed architecture onto the various applications of banks and financial services.' Below the abstract, there is a 'Published In:' section listing 'Recent Advances in Information Technology (RAIT), 2016 3rd International Conference on'. At the bottom, there are two columns of metadata: 'Date of Conference: 3-5 March 2016', 'Date Added to IEEE Xplore: 11 July 2016', '▼ ISBN Information: Electronic ISBN: 978-1-4799-8579-1, Print on Demand(PoD) ISBN: 978-1-4799-8580-7', 'INSPEC Accession Number: 16140907', 'DOI: 10.1109/RAIT.2016.7507897', 'Publisher: IEEE', and 'Conference Location: Dhanbad, India'.

Figure D.4: Publication [4]

Table D.5: Fact sheet for publication [5]

Title	Energy Efficient Virtual Machine Placement in Cloud Data Centers Using Modified Intelligent Water Drop Algorithm
Authors	Chandra Shekhar Verma, V. Dinesh Reddy, G.R. Gangadharan, Atul Negi
Publication	In Proceedings of the 13th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS) SITIS, Jaipur, India, Pages 13-20, 2017
ISBN/ISSN	978-1-5386-4283-2
DOI	10.1109/SITIS.2017.14
Status	Published
Publisher	IEEE
Publication Type	Conference

The screenshot shows a digital publication page. At the top left, the title "Energy Efficient Virtual Machine Placement in Cloud Data Centers Using Modified Intelligent Water Drop Algorithm" is displayed in a large, bold font. Below the title is a blue button with the text "Sign In or Purchase to View Full Text". To the right of the title is a "Related Articles" section with a light orange background, containing two article titles and a "View All" link. Below the title and button is a light blue bar with the author count "4" and the names "Chandra Shekhar Verma ; V. Dinesh Reddy ; G.R. Gangadharan ; Atul Negi", along with a "View All Authors" link. Below this bar is a navigation menu with tabs for "Abstract", "Authors", "Figures", "References", "Citations", "Keywords", "Metrics", and "Media". The "Abstract" tab is selected. The abstract text follows, starting with "Abstract:" and describing the paper's focus on energy-efficient VM placement in cloud data centers. At the bottom, the "Published in:" information is provided.

Energy Efficient Virtual Machine Placement in Cloud Data Centers Using Modified Intelligent Water Drop Algorithm

[Sign In or Purchase to View Full Text](#)

Related Articles

- [Dynamic load-balancing for BSP Time Warp](#)
- [On aspect-orientation in distributed real-time dependable systems](#)

[View All](#)

4 Author(s) | [View All Authors](#)

Chandra Shekhar Verma ; V. Dinesh Reddy ; G.R. Gangadharan ; Atul Negi

Abstract | Authors | Figures | References | Citations | Keywords | Metrics | Media

Abstract:

Cloud Computing is an emerging distributed computing paradigm for the dynamic provisioning of computing services on demand over the internet. Due to heavy demand of various IT services over the cloud, energy consumption by data centers is growing significantly worldwide. The intense use of data centers leads to high energy consumptions, excessive CO2 emission and increase in the operating cost of the data centers. Although many virtual machine (VM) placement approaches have been proposed to improve the resource utilization and energy efficiency, most of these works assume a homogeneous environment in the data centers. However, the physical server configurations in heterogeneous data centers lead to varying energy consumption characteristics. In this paper, we model and implement a modified Intelligent Water Drop algorithm (MIWD) algorithm for dynamic provisioning of virtual machines on hosts in homogeneous and heterogeneous environments such that total energy consumption of a data center in cloud computing environment can be minimized. Experimental results indicate that our proposed MIWD algorithm is giving superior results.

Published in: Signal-Image Technology & Internet-Based Systems (SITIS), 2017 13th International Conference on

Figure D.5: Publication [5]

Synopsis of
**Energy Efficient Data Center Management
Strategies**

Thesis submitted for the degree of

Doctor of Philosophy

in

Computer Science

by

Vemula Dinesh Reddy Reg. No. 14MCPC19

under the guidance of

Dr. G.R. Gangadharan



Institute for Development and Research in Banking Technology

(Established by Reserve Bank of India)

Hyderabad-500057



School of Computer & Information Sciences

University of Hyderabad

Hyderabad-500046

Telangana, India

1 Introduction

The proliferation of Cloud computing has resulted in the establishment of large-scale data centers around the world. These data centers are energy-intensive building types and consume large amount of electrical energy resulting in high operating costs and carbon dioxide (CO₂) emissions to the environment. While the energy consumption of data centers is already significant, the growth of the global cloud-based economy along with the need to power our connected devices and always-on lifestyles increases the required resources even further[1, 2]. The world’s Information and communications technology (ICT) infrastructure is estimated to consume 1,500 TWh of electricity, roughly 10% of global usage. Furthermore, energy consumption is expected to continue slightly increasing in the near future, rising 4% from 2014-2020 [3]. High power consumption generates heat and requires an accompanying cooling system that costs in a range of \$2 to \$5 million per year for traditional data centers [4]. Servers and IT equipment are responsible for 55% of the energy used by the data center followed by 30% for cooling equipment [5]. A source of high energy consumption is not just the amount of computing resources used and power inefficiency of the hardware but also lies in the inefficient usage and dynamic power ranges of servers. A survey by IBM shows that the average resource utilization rate is lower than 20% in data centers. Even completely idle servers still consume up to 70% of their peak power [4]. Therefore, keeping servers underutilized is highly inefficient from the energy consumption perspective.

Green cloud computing requires energy efficient use of data centers with minimum impact on environment [5]. Green computing in the cloud can be achieved by eliminating inefficiencies and waste in the way electricity is delivered to computing resources, and in the way, these resources are utilized to serve application workloads. This can be done by implementing effective policies and algorithms for energy-aware resource management in data centers [6].

The operational efficiency of the data centers assumes central importance. Even small gains in efficiency translate into end-user perceivable cost reductions, providing key competitive advantage. From the perspective of users, the performance of the data center is measured in terms of response time, virtual machine (VM) provisioning time, etc. To serve the customers in a better way, data center providers should follow an optimal VM provisioning within very short time. Traditional algorithms like Best-Fit, First-Fit, and Modified Best-Fit, etc. will not be efficient for large-scale data centers as they take a long time for optimal VM provisioning [7, 8]. The objective of the energy efficient resource provisioning is to find the near-optimal solution that improves the resource utilization and decrease the energy consumption of the data center in an acceptable time.

The operations of a data center are quickly transforming from individual and disconnected tactical activities with a primary historical goal of “high service levels at any cost” to a planned and predictable approach with the modified metric “service at what cost” [9]. Energy efficiency of the data center is influenced by many factors, such as data center layout design and characteristics, ambient weather conditions, rack density, the operation of heating, ventilation, and air conditioning (HVAC) systems and their behavior. This complex connection makes it hard to predict data center energy consumption. With sensor data and information about data center

operations, forecasting energy consumption helps in planning and operations of data centers. Well-planned resource provisioning makes good return on investment and elasticity of computing infrastructures.

2 Problem Statement and Objectives

This thesis focuses on modeling energy-efficient data center management strategies, i.e., ensuring that computing resources are efficiently utilized to serve application workloads to minimize energy consumption, planned provisioning with demand forecasting, and evaluating the current practices and implementing the new practices. In particular, the following research problems are investigated:

- **Energy efficient placement of virtual machines.** To reduce the number of required host servers in a data center, it is necessary to have efficient virtual machine (VM) placement strategies in place. Determining the optimal placement of VMs to improve physical resources utilization and to reduce the energy consumption of the data centers while satisfying the service level of agreement (SLA) is an essential aspect of the data center.
- **Selecting appropriate VMs to migrate.** Determining a set of VMs that should be migrated from an overloaded or under-loaded host has a significant impact on the VM migration time and energy consumption of the data center, and can cause the SLA violation. So, designing a VM selection policy, considering different resources along with CPU utilization plays an important role in improving the energy efficiency of the data centers. The problem consists in determining the best subset of VMs to migrate that will provide the most beneficial system reconfiguration.
- **Forecasting data center energy demand.** Forecasting data center electrical energy demand is very challenging due to data center scales, dynamics in workloads and complexity involved. Developing forecasting models with accurate predictions gives operators enough time to avoid the risk of over-provisioning during non-peak hours, and reduces the risk of under-provisioning. To provide a capacity management process for resource pools in a just in time manner, it is necessary to have an efficient forecast model for data centers to predict and estimate proper energy demand in real-world situations.
- **Analyzing metrics and practices of the data center.** In order to, predict growth or set effective goals, it is important to choose the correct metric and being aware of their expressivity and potential limitations. Understanding and analyzing data center metrics allows the operators to have a better view on possible inefficiencies by focusing on the core parameters. All the data center operators need to document and then automate their existing practices. It is necessary to compare the current approaches in a data center with industry standards and assess whether the practices are still valid and/or optimal. Determine and implement the best practices for data center operations to optimize the workflows and to decrease operating cost in the long term.

To deal with the challenges associated with the said research problems, the following objectives have been delineated:

- Explore the research in the area of energy-efficient resource management in a data center using dynamic VM provisioning and selection to gain a systematic understanding of the existing techniques and approaches.
- Propose novel and optimal methods for dynamic VM placement and selection.
- Develop accurate methods for forecasting data center energy demand that helps in planning operations of the data center.
- Analyze metrics and their current values for the data centers.
- Explore the current practices of the data center and develop the best practices for data center sustainability.

3 Literature Survey

3.1 Virtual Machine Placement and Selection

Energy efficient resource allocation and selection is a challenging issue in data centers. Beloglazov et al. [10] proposed a threshold based modified best fit decreasing algorithm for virtual machines (VM) provisioning on hosts. This algorithm heuristically uses varying threshold level for CPU utilization. It performs live migration of VMs using minimum migration time method. However, dynamic consolidation may lead to more SLA violation. Wang et al. [11] proposed a mixed integer programming approach to solve virtual machine placement problem and developed linear and nonlinear power consumption models. To find the near optimal solution, they used a heuristic based iterative rounding technique. However, this takes more time with increasing virtual machines and servers. Tang et al. [7] proposed a hybrid genetic algorithm for VM placement analyzing network overhead. But it has not considered the cost of VM migration. Buyya et al. [12] proposed a non-power aware policy and DVFS policies for energy saving. For virtual machine placement, Power Aware Best Fit Decreasing (PABFD) algorithm is proposed in [10, 13]. However, for strict SLA, this approach consumes more energy.

Goudarzi et al. [14] used dynamic programming to create multiple copies of VMs and to put them on servers and then used local search to find underutilized servers. In this approach, only the original VM serves the request while all other copies will remain ideal. However, SLA violation is not considered, and multiple copies of VMs create network overhead. Wu et al. [15] proposed a Simulated Annealing (SA) based VM Placement algorithm. It is suitable for static VM consolidation but not for dynamic VM consolidation. Wang et al. [16] proposed efficient VM placement optimization based on particle swarm optimization (PSO) with the local fitness first scheme. This approach did not consider over-utilization of the hosts which may lead to SLA violations. Kumar et al. [17] focused on energy efficient VM allocation using PSO in cloud operations. Their objective is to minimize the total resource wastage in

a data center, but they have not considered the over utilization of the hosts. Jeyarani et al. [18] proposed a Self Adaptive PSO (SAPSO) for VM allocation in data centers. Though SAPSO performs well in terms of energy consumption, the analysis of SLA violations shows that this method does not maintain the quality of service when the VM requests increase.

Virtual machine migration is a major way of reducing unnecessary consumption in a data center [19]. Buyya et al. [12] proposed various strategies such as single threshold (ST), Minimum migration (MM), Highest Potential growth (HPG) and random choice policies. For VM migration they illustrate that ST policy and MM policy consume nearly same energy, but MM policy shows very less VM migrations compared to ST policy. Zhou et al. [20] proposed an adaptive three threshold energy aware algorithm (ATEA) considering resource usage patterns of virtual machines. ATEA may not guarantee the optimal energy consumption of data center. Dai et al. [21] proposed a minimum power and the minimum communication algorithms based on integer programming for energy efficient VM placement and migration respectively. However, it did not consider SLA violation. Wang et al. [22] presented a decentralized double threshold VM selection policy considering the utilization of the physical nodes. However, they have not considered the energy consumption and this strategy may not give the optimal solution always. Bose et al. [23] proposed cloud spider architecture that integrates replication and scheduling methods to minimize the live migration costs across Wide Area Networks. This method requires additional storage requirements. Zhang et al. [24] presented an approximate approach based on bin packing algorithm to migrate virtual machines. They considered the resource utilization and the migration cost to get the optimal solution. Bobroff et al. [25] proposed a dynamic VM consolidation algorithm. They used time series forecasting and bin-packing heuristic for minimizing the physical resources. However, multiple resources are not considered in this approach. Cardoso et al. [26] proposed a PowerExpandMinMax algorithm, for VM consolidation based on min-max and share features of VM technologies. This algorithm does not follow the restrict resource constraints and live VM migration.

VM allocation and selection are NP hard problems [27, 28], inferring that an optimal arrangement cannot be found in deterministic polynomial time. Traditional approaches like First-fit, best-fit, best fit decreasing, etc. are deterministic in nature and do not always guarantee the optimal solution. Bio-inspired approaches like PSO and genetic algorithm (GA) are iterative in nature and provide global optimal solutions in case of energy efficient virtual machine allocation in data centers [29–31]. However, improving the performance, exploration and exploitation capacities of the above-mentioned algorithms is still a challenging issue. Further, there is a need to develop models, best practices and algorithms considering energy consumption of the data centers at various levels to reduce the operational cost.

3.2 Forecasting Techniques for Data Center Management

Gmach [32] et al. use a trace-based approach to capacity management that relies on the characterization of workload demand patterns and predict future demands based on the patterns. They use a three-stage approach to recommend the most

likely pattern for the workload. Aksanli [33] et al. design an adaptive data center job scheduler which utilizes short-term prediction of solar and wind energy production. Abhishek et al. [34] combine online measurements with prediction and resource allocation techniques to provide guarantees to web applications running on shared data centers. Application workloads are modeled using a time-domain description of a generalized processor sharing (GPS) server. Time series analysis techniques are used to update the parameters of this model. Prevost et al. [35] et al. proposed load demand prediction using the neural network and autoregressive linear prediction algorithms and combined with stochastic state transition models to optimal resource allocation by minimizing the energy consumed. Kong et al. [36] proposed a fuzzy prediction method to model the uncertain workload and the vague availability of virtualized server nodes, by using the type-I and type-II fuzzy logic systems. Farahnakian et al. [37] present a CPU usage prediction method based on the linear regression technique and this is integrated with the live migration process to predict over-loaded and under-loaded hosts. The proposed approach approximates the short-time future CPU utilization based on the history of usage in each host. Ricciardi et al. [38] developed a methodology, that exploits load fluctuations and effectively control the system using service-demand matching algorithm and determines the subset of servers that may be powered off in data centers. However, their proposed method did not always guarantee the optimal solution and the authors have not considered the SLA violations.

3.3 Analyzing Metrics for Sustainable Data Centers

The primary step in developing a model to capture the effects of data center is to decide which dimensions are relevant, define the metrics, and populate them[39]. The Green Grid consortium proposed the power usage effectiveness (PUE)[40], which currently is the prevailing metric. Schaeppli et al. explored energy related metrics for IT equipment, data storage and network equipment [41]. Metrics to monitor and control the air flow in a data center are discussed in [42, 43]. Capozzoli et al. reviewed thermal, power and energy consumption metrics [44]. Chen et al. identified and presented usage-centric green performance indicators at various levels such as server and storage [45]. Wang et al. presented a set of performance metrics for a green data center [46]. With numerous metrics available for measuring data center efficiency, there is a need to analyse these metrics to choose the correct metric and being aware of their expressivity and potential limitations.

3.4 Best Practices for Sustainable Data Centers

Strong et al. focused on room-level bypass airflow and proposed necessary changes at three levels of data centers to reduce operating expenses and increase cooling capacity [47]. Hamann et al. [48] developed a set of measurement-based best practices metrics and guidance for improving the energy efficiency of a data center. With the data and available key metrics, they derived insights into the sources of energy inefficiencies. Lau et al. [49] developed a rating system for server power supplies. They explored criteria for lower load conditions and the energy efficiency opportunities in server power supply. Al-Fares et al. [50] described how to leverage commodity Ethernet switches

to support the full aggregate bandwidth of clusters consisting of tens of thousands of elements. Evans et al. [51] proposed humidity control mechanisms and design guidelines for computing infrastructure installations. Further, they proposed practices for achieving desired humidity and described issues with over-humidification.

4 Proposed Contents of the Thesis

This thesis consists of 7 chapters including an introductory and a concluding chapter followed by Appendices. The content of this thesis can be broadly divided into 5 categories: a systematic literature review, novel algorithms for dynamic VM placement and selection, novel algorithms for data center energy demand prediction, analyzing metrics for energy efficient data centers and evaluation of best practices for sustainable data centers. The content of each of these chapters is summarized as follows:

4.1 Chapter 1 : Introduction

Chapter 1 provides the motivation and importance of improving the energy efficiency in data centers. This chapter begins with the discussion of energy efficient data center management strategies that include optimal VM placement and selection, demand forecasting, etc. Further, it presents the problem statement and contributions.

4.2 Chapter 2: Literature Review

Energy efficient techniques for resource management is inevitable to reduce the electricity consumption in cloud data centers. To identify open challenges in the energy efficient virtual machine placement and consolidation and to facilitate further advancements, it is essential to synthesize the research in this area conducted to date. This chapter presents a literature review on the virtual machine placement and consolidation using soft computing and machine learning techniques and presents a taxonomy based on their objectives.

4.3 Chapter 3: Energy aware Virtual Machine Placement and Selection Approaches in Cloud Data Centers

This chapter presents proposed techniques for energy efficient VM placement and selection .

Technique 1: We propose a Modified Discrete Particle Swarm Optimization (MDPSO) approach for optimal energy-aware virtual machine placement that minimizes the power consumption of the physical machine by estimating the increase in the power consumption before a VM is placed onto the server.

Technique 2: Taking the recent advances in multi-core architectures, we develop a parallelized optimization algorithm called “Interactive PSO-GA” (IPSOGA). IPSOGA performs parallel processing of particle swarm optimization (PSO) and genetic

algorithm (GA) using multithreading and shared memory for information exchange to enhance convergence time and global exploration. To enhance the population of each generation, we incorporate the social interaction between the algorithms. This technique helps to balance between improving convergence time and accuracy.

Technique 3: Inspired by the imitating behavior of humans, we developed a swarm based approach for virtual machine placement namely imitation based optimization (IBO). The search for the optimal solution is completed using a particle swarm optimization-like method but that does not contain inertia and velocity components. The particles try to imitate the best solution. This technique generates optimal solution that tend to satisfy the structural information and provides consistency.

Further, we present a novel virtual machine selection method considering the factors such as memory, bandwidth and size of the VM (MBS-VM). This method optimally selects the virtual machines from a under/over utilized server and performs migration to further improve the energy efficiency in a data center.

The first part of this chapter is published in **Soft Computing, Springer**.

4.4 Chapter 4: Machine Learning Approaches for Forecasting Data Center Energy Demand

The energy efficiency of the data center is influenced by many factors, such as data center layout design and characteristics, ambient weather conditions, rack density, the operation of HVAC systems and their behavior. This complex connection makes it hard to predict data center energy consumption. With sensor data and information about devices operations, forecasting energy consumption for data centers helps in planning and operations. This chapter presents promising ideas and results about the data center energy demand prediction using two machine learning approaches.

Technique 1: Multi layer neural networks involve multiple levels of non linearity and they perform hierarchical feature extraction. These models are able to learn useful information of raw data and exhibit high performance. We propose “Multi layer feed forward neural networks” for forecasting energy demand of the data centers. We trained our model with the popular back-propagation algorithm where the weights connecting the layers are updated in an iterated manner.

Technique 2: Although a multilayer backpropagation network with enough neurons can implement just about any function, backpropagation will not always find the correct weights for the optimum solution. We need reinitialize the network and retrain several times to guarantee that you have the best solution. To overcome the said problems and to improve the accuracy of the predictions, we propose a deep learning approach with “parallel stochastic gradient descent” training for forecasting energy demand of the data centers.

4.5 Chapter 5: Metrics for Sustainable Data Centers

Most data centers lack an integrated energy management system that jointly optimizes and controls all its components to reduce the operational cost. There are a multitude of metrics available to analyse energy efficiency of the data centers. In order to, predict growth or set effective goals, it is important to choose the correct metric and being aware of their expressivity and potential limitations. Understanding and analyzing data center metrics allows the operators to have a better view on possible inefficiencies by focusing on the core parameters. We present an analysis of metrics that are commonly used in data centers, starting from the power grid and going all the way up to the service delivery. We propose a classification based on the different core dimensions of data center operations such as energy efficiency, cooling, greenness, performance, thermal and air management, network, security, storage, and financial impact. Furthermore, we derive relationships between metrics, and discuss the advantages and disadvantages of each metric in order to expose the research gaps and illustrate the latest research trends in computing the efficiency of a data center. This proposed work on analysis of metrics is published in **IEEE Transaction on Sustainable Computing**.

4.6 Chapter 6: Best Practices for Sustainable Data Centers

All the data center operators need to compare their current approaches with industry standards and assess whether their practices are still valid and/or optimal. It is thus essential to consider any opportunity to reduce the energy consumption of the data centers, both in design and operations. In this chapter, we have analyzed seven data centers in India and the Netherlands. Based on our findings and industry standards, we propose a set of best practices to improve the energy efficiency of the data centers which spans the categories of Energy Efficiency, Cooling, Air and Thermal management, Greenness, Storage, and Networks. Following some of these best practices, data centers surveyed in our study have achieved 10 – 20% improvements in their energy consumption. The chapter provides efficient alternatives in daily operations of the data centers and costs saving opportunities. The proposed work of this chapter is accepted for publication in **IT Professional, IEEE**.

4.7 Chapter 7: Conclusions and Future Directions

Chapter 7 summarizes the contributions of the thesis and outlines the future directions. This thesis develops the techniques for optimizing the energy consumption in a data center using energy efficient VM placement and selection. Further, this thesis tackles the problem of forecasting data center energy consumption for better planning and operations of data centers. In our future work, we will consider implementing extra constraints on the VM placement to co-allocate VMs on the physical server or to performance or privacy concerns. There is scope for improvement of VM placement algorithms using buffered VM placement requests. Another direction of future research is to exploit VM resource usage pattern for more efficient resource provisioning and higher energy efficiency.

5 Acknowledgements

This thesis received funding from the Netherlands Organization for Scientific Research (NWO) in the framework of the Indo Dutch Science Industry Collaboration programme in relation to project NextGenSmart DC (629.002.102).

6 List of Articles Published During the Candidature

6.1 Journal Papers

- (1) V. Dinesh Reddy, B. Setz, G. Subrahmanya V.R.K. Rao, G.R. Gangadharan, Marco Aiello, “Metrics for Sustainable Data Centers”, *IEEE Transactions on Sustainable Computing*, Vol. 2, No. 3, pp. 290-303, doi:10.1109/TSUSC.2017.2701883, 2017 (**DBLP**).
- (2) V. Dinesh Reddy, G. R. Gangadharan, and G. Subrahmanya V. R. K. Rao, “Energy-aware virtual machine allocation and selection in cloud data centers”, *Soft Computing*, Springer, ISSN 1433-7479, doi: 10.1007/s00500-017-2905-z, 2017 (**DBLP, SCOPUS, and SCI indexed**).
- (3) V. Dinesh Reddy, B. Setz, G.S.V.R.K. Rao, G. R. Gangadharan, and M. Aiello, “Best Practices for Sustainable Data Center”, *IEEE IT Professional (In Press)*, 2017 (**DBLP, SCOPUS, and SCI indexed**).

6.2 Papers in Conference Proceedings

- (4) V. Dineshreddy and G. R. Gangadharan, “Towards an Internet of Things framework for financial services sector,” *3rd International Conference on Recent Advances in Information Technology (RAIT)*, Dhanbad, pp. 177-181, doi: 10.1109/RAIT.2016.7507897, 2016 (**DBLP and SCOPUS indexed**).

6.3 Papers Communicated

- (5) V. Dinesh Reddy and G.R. Gangadharan, “Forecasting Data Center Energy Demand: A Deep Learning Approach,” *Sustainable Computing, Informatics and Systems*, Elsevier (**DBLP, SCOPUS and SCI indexed**).
- (6) V. Dinesh Reddy and G.R. Gangadharan, M. Aiello, “Energy efficient resource management in cloud data centers using a hybrid evolutionary algorithm,” *Soft Computing*, Springer (**DBLP, SCOPUS and SCI indexed**).

REFERENCES

- [1] T. Dandres, N. Vandromme, G. Obrekht, A. Wong, K.K. Nguyen, Y. Lemieux, M. Cheriet, and R. Samson, “Consequences of future data center deployment in Canada on electricity generation and environmental impacts: a 2015–2030 prospective study”, *Journal of Industrial Ecology*, 21(5), pp. 1312–1322, 2017.
- [2] J. Koomey, “Growth in data center electricity use 2005 to 2010”, *A report by Analytical Press, completed at the request of the New York Times*, pp. 9–10, 2011.
- [3] M.P. Mills, “The Cloud Begins With Coal- An overview of the electricity used by the global digital ecosystem”, *Technical report, Digital Power Group*, 2013.
- [4] A.K. Kiani and N. Ansari, “On The Fundamental Energy Trade-offs of Geographical Load Balancing”, *IEEE Communications Magazine*, 55(5), pp. 170–175, 2017.
- [5] S. Murugesan and G.R. Gangadharan, *Harnessing green IT: Principles and practices*, Wiley Publishing, 2012.
- [6] B. Pernici, M. Aiello, J. vom Brocke, B. Donnellan, E. Gelenbe, and M. Kretsis, “What IS Can Do for Environmental Sustainability: A Report from CAiSE’11 Panel on Green and Sustainable IS”, *CAIS*, 30(18), 2012.
- [7] M. Tang and S. Pan, “A hybrid genetic algorithm for the energy-efficient virtual machine placement problem in data centers”, *Neural Processing Letters*, 41(2), pp. 211–221, 2015.
- [8] S.S. Gill, R. Buyya, I. Chana, M. Singh, and A. Abraham, “BULLETT: Particle Swarm Optimization Based Scheduling Technique for Provisioned Cloud Resources”, *Journal of Network and Systems Management*, 26(2), pp. 361–400, 2018.
- [9] H. Geng, *Data center handbook*, John Wiley & Sons, 2014.
- [10] A. Beloglazov and R. Buyya, “Adaptive threshold-based approach for energy-efficient consolidation of virtual machines in cloud data centers.”, *In Proceedings of the 8th International Workshop on Middleware for Grids, Cloud and e-Science*, p. 4, 2010.
- [11] Y. Wang and Y. Xia, “Energy Optimal VM Placement in the Cloud”, *In Proceedings of the IEEE 9th International Conference on the Cloud Computing*, pp. 84–91, 2016.
- [12] R. Buyya, A. Beloglazov, and J. Abawajy, “Energy-efficient management of data center resources for cloud computing: a vision, architectural elements, and open challenges”, *arXiv:1006.0308*, 2010.
- [13] A. Beloglazov and R. Buyya, “Optimal online deterministic algorithms and adaptive heuristics for energy and performance efficient dynamic consolidation of virtual machines in cloud data centers”, *Concurrency and Computation: Practice and Experience*, 24(13), pp. 1397–420, 2012.

- [14] H. Goudarzi and M. Pedram, “Energy-efficient virtual machine replication and placement in a cloud computing system”, *In Proceedings of the IEEE 5th International Conference on Cloud Computing*, pp. 750–757, 2012.
- [15] Y. Wu and M. Tang, “A simulated annealing algorithm for energy efficient virtual machine placement”, *In Proceedings of the IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pp. 1245–1250, 2012.
- [16] S. Wang, Z. Liu, Z. Zheng, Q. Sun, and F. Yang, “Particle swarm optimization for energy-aware virtual machine placement optimization in virtualized data centers”, *In Proceedings of the International Conference on Parallel and Distributed Systems (ICPADS)*, pp. 102–109, 2013.
- [17] D. Kumar and Z. Raza, “A PSO Based VM Resource Scheduling Model for Cloud Computing”, *In Proceedings of IEEE International Conference on Computational Intelligence & Communication Technology (CICT)*, pp. 213–219, IEEE, 2015.
- [18] R. Jeyarani, N. Nagaveni, and R.V. Ram, “Self Adaptive Particle Swarm Optimization for Efficient Virtual Machine Provisioning in Cloud”, *International Journal of Intelligent Information Technologies*, 7(2), pp. 25–44, 2011.
- [19] Y. Yang, B. Mao, H. Jiang, Y. Yang, H. Luo, and S. Wu, “SnapMig: Accelerating VM Live Storage Migration by Leveraging the Existing VM Snapshots in the Cloud”, *IEEE Transactions on Parallel and Distributed Systems*, 2018.
- [20] Z. Zhou, Z. Hu, and K. Li, “Virtual Machine Placement Algorithm for Both Energy-Awareness and SLA Violation Reduction in Cloud Data Centers”, *Scientific Programming*, (5612039), 2016, doi: 10.1155/2016/5612039.
- [21] X. Dai, J. M. Wang, and B. Bensaou, “Energy-efficient virtual machine placement in data centers with heterogeneous requirements”, *In Proceedings of the IEEE 3rd International Conference on Cloud Networking (CloudNet)*, pp. 161–166, 2014.
- [22] X. Wang, X. Liu, L. Fan, and X. Jia, “A decentralized virtual machine migration approach of data centers for cloud computing”, *Mathematical Problems in Engineering*, Hindawi Publishing Corporation, (878542), 2013.
- [23] S.K. Bose, S. Brock, R. Skeoch, and S. Rao, “CloudSpider: Combining replication with scheduling for optimizing live migration of virtual machines across wide area networks”, *In Proceedings of the 11th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing*, pp. 13–22, 2011.
- [24] X. Zhang, K. Li, and Y. Zhang, “Minimum-cost virtual machine migration strategy in data center”, *Concurrency and Computation: Practice and Experience*, 27(17), pp. 5177–5187, 2015.
- [25] N. Bobroff, A. Kochut, and K. Beaty, “Dynamic placement of virtual machines for managing sla violations”, *In Proceedings of the 10th IFIP/IEEE International Symposium on Integrated Network Management*, pp. 119–128, 2007.

- [26] M. Cardoso, M. R. Korupolu, and Singh A. Shares and, “Shares and utilities based power consolidation in virtualized server environments”, *In Proceedings of the IFIP/IEEE International Symposium on Integrated Network Management*, pp. 327–334, 2009.
- [27] R.G. Michael and S.J. David, *Computers and intractability: a guide to the theory of NP-completeness*, W. H. Freeman & Co, New York, 1979.
- [28] F. Hao, M. Kodialam, T.V. Lakshman, and S. Mukherjee, “Online allocation of virtual machines in a distributed cloud”, *IEEE/ACM Transactions on Networking*, 25(1), pp. 238–249, 2017.
- [29] G. Wu, M. Tang, Y. Tian, and W. Li, “Energy-efficient virtual machine placement in data centers by genetic algorithm”, *In proceedings of International Conference on Neural Information Processing*, pp. 315–323, 2012.
- [30] Y.D. Valle, G.K. Venayagamoorthy, S. Mohagheghi, J.C. Hernandez, and R.G. Harley, “Particle swarm optimization: basic concepts, variants and applications in power systems”, *IEEE Transactions on Evolutionary Computation*, *IEEE*, 12 (2), pp. 171–195, 2008.
- [31] A. Askarzadeh, “A Memory-based Genetic Algorithm for Optimization of Power Generation in a Microgrid”, *IEEE Transactions on Sustainable Energy*, 2017.
- [32] D. Gmach, J. Rolia, L. Cherkasova, and A. Kemper, “Capacity management and demand prediction for next generation data centers”, *In Proceedings of the IEEE International Conference on Web Services (ICWS)*, pp. 43–50, 2007.
- [33] B. Aksanli, J. Venkatesh, L. Zhang, and T. Rosing, “Utilizing green energy prediction to schedule mixed batch and service jobs in data centers”, *ACM SIGOPS Operating Systems Review*, 45(3), pp. 53–57, 2012.
- [34] A. Chandra, W. Gong, and P. Shenoy, “Dynamic resource allocation for shared data centers using online measurements”, *ACM SIGMETRICS Performance Evaluation Review*, 31(1), pp. 300–301, 2003.
- [35] J. Prevost John, N. KranthiManoj, K. Brian, and M. Jamshidi, “Prediction of cloud data center networks loads using stochastic and neural models”, *In Proceedings of the 6th International Conference on System of Systems Engineering (SoSE)*, pp. 276–281, 2011.
- [36] X. Kong, C. Lin, Y. Jiang, W. Yan, and X. Chu, “Efficient dynamic task scheduling in virtualized data centers with fuzzy prediction”, *Journal of network and Computer Applications*, 34(4), pp. 1068–1077, 2011.
- [37] F. Farahnakian, P. Liljeberg, and J. Plosila, “LiRCUP: Linear regression based CPU usage prediction algorithm for live migration of virtual machines in data centers”, *In Proceedings of the 39th EUROMICRO Conference on Software Engineering and Advanced Applications (SEAA)*, pp. 357–364, 2013.
- [38] S. Ricciardi, D. Careglio, J. Sole-Pareta, U. Fiore, and F. Palmieri, “Saving energy in data center infrastructures”, *In Proceedings of International Conference on Data Compression, Communications and Processing*, pp. 265–270, 2011.

- [39] T. Daim, J. Justice, M. Krampits, M. Letts, G. Subramanian, and M. Thirumalai, “Data center metrics: an energy efficiency model for information technology managers”, *Management of Environmental Quality: An International Journal*, 20(6), pp. 712–731, 2009.
- [40] C.L. Belady, A. Rawson, J. Pfeuger, and T. Cader, “The Green Grid Data Center Power Efficiency Metrics: PUE and DCiE”, *Green Grid, white paper-06*, 2007.
- [41] B. Schaeppi, T. Bogner, A. Schloesser, L. Stobbe, and M.D. De Asuncao, “Metrics for energy efficiency assessment in data centers and server rooms”, pp. 1–6, 2012.
- [42] A. Capozzoli, G. Serale, L. Liuzzo, and M. Chinnici, “Thermal Metrics for Data Centers: A Critical Review”, *Energy Procedia*, 62, pp. 391–400, 2014.
- [43] R. Tozer and M. Salim, “Data center air management metrics-practical approach”, *In Proceedings of the 12th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, pp. 1–8, 2010.
- [44] A. Capozzoli, M. Chinnici, M. Perino, and G. Serale, “Review on Performance Metrics for Energy Efficiency in Data Center: The Role of Thermal Management”, *In Proceedings of Third International Workshop on Energy Efficient Data Centers, Cambridge, UK*, pp. 135–151, 2015.
- [45] D. Chen, B. Pernici, E. Henis, R.I. Kat, D. Sotnikov, C. Cappiello, A.M. Ferreira, M. Vitali, T. Jiang, and J. Liu, “Usage centric green performance indicators”, *ACM SIGMETRICS Performance Evaluation Review*, 39(3), pp. 92–96, 2011.
- [46] L. Wang and S.U. Khan, “Review of performance metrics for green data centers: a taxonomy study”, *The Journal of Supercomputing*, 63(3), pp. 639–656, 2013.
- [47] P.E. Lars Strong, “Reducing Room-Level Bypass Airflow Creates Opportunities to Improve Cooling Capacity and Operating Costs”, *White Paper, Upsite Technologies*, 2013.
- [48] H.F. Hamann, M. Schappert, M. Iyengar, T. van Kessel, and A. Claassen, “Methods and techniques for measuring and improving data center best practices”, pp. 1146–1152, 2008.
- [49] L. Henry, “Efficient Power Supplies for Data Center and Enterprise Servers”, *Technical report, Southern California Edison Design and Engineering Services*, 2013.
- [50] M. Al-Fares, A. Loukissas, and A. Vahdat, “A scalable, commodity data center network architecture”, *ACM SIGCOMM Computer Communication Review*, 38(4), pp. 63–74, 2008.
- [51] T. Evans, “Humidification strategies for data centers and network rooms”, *APC distributors, White Paper-58*, 2004.