

Energy Management in Data Centers from Design to Operations and Maintenance

Montri Wiboonrat
KMITL Business School
King Mongkut's Institute of Technology Ladkrabang
Bangkok, Thailand
montri.wi@kmitl.ac.th

Abstract—Data centers are the information factories of the digital evolution. Creating, storing, processing, distributing, and analyzing data all need energy. Therefore, data center industry consumes energy more than 2 percent of the global electricity consumption. Energy efficiency need to discuss at the outset of data center design. The root cause of oversizing data center design is the research question because this will affect investment or CAPEX and long-term operating costs or OPEX of data center as long as data center life cycle (DCLC). Data center measurement in power usage effectiveness (PUE) unit helps data center owners and consultants realized on relationship between oversizing data center design and total cost of ownership (TCO). The research results propose modular data center as a solution to handle uncertainty demand of IT equipment, scalability for growth as your need, flexibility in any size of infrastructure, fast deployment because of prefabricated design, and more efficiency by applying energy management platform called data center infrastructure management (DCIM).

Index Terms—energy management, energy efficiency, data center, operations and maintenance, TCO

I. INTRODUCTION

In 2018, the globally electricity demand of data centers was an estimated 198 terawatt hours or around 1% of power consumption in the world [1]. The global traffic has tripled from 2015 to 2019 and expected will be doubled by 2022 [2]. A data center is one of the most investment assets of any IT enterprise. Power requirements for data center have become gigantic. Data center requires power for IT equipment, power distribution systems, and cooling systems. As data center basic design and installation, the amount of power needed for cooling system is approximately equal to that required to operate IT equipment. For high performance computing data center, the power costs are significant. The top 7 biggest information factories by market capitalization in 2020 are Microsoft, Apple, Amazon, Alphabet, Facebook, Alibaba, and Tencent. During and after COVID-19 pandemic, the demand of Internet and 5G traffic skyrockets, the information industry could lead to a high demand in energy use. Each year, data centers use an estimated 200 terawatt hours (TWh), therefore data center's energy usage must be vigilantly management, as seen the demand of energy forecasting in data center industry in Fig. 1.

Before COVID-19, the rising demand for connectivity data, ICT's energy consumption is staying nearly flat, but after COVID-19 the increasing Internet traffic and data loads are countering by increasing efficiencies and utilized server and storage resources. The mechanisms of energy conservation can apply with 2 scenarios; existing data centers and new data centers. For existing data centers, consolidation, virtualization, and utilization are the keys while for new data centers, the design philosophy must be optimal between energy efficiency, system reliability, utilization, and total cost of ownership (TCO).

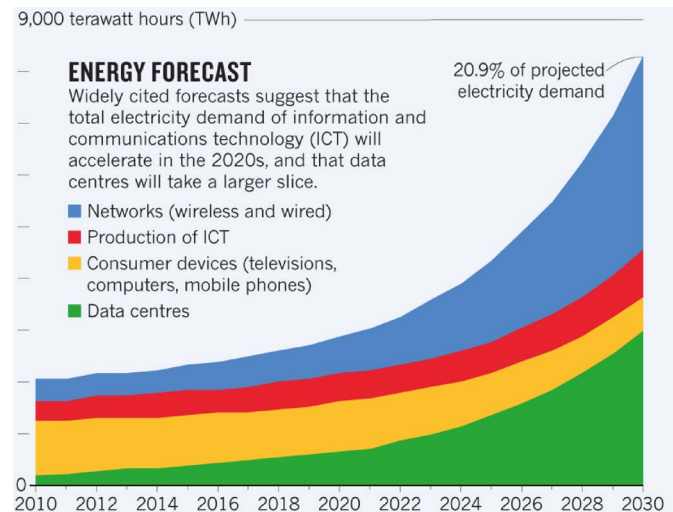


Figure 1. Energy forecast in data center industry [4].

A modern data center still shares many aspects with predecessor architectures, topologies, structures, and services. The purpose of data center is to provide information services. To provide information services, it requires continuous power for continual functioning must be met. The design of the data center must consider the functional requirements of the data center are:

- To provide space (layout design) of servers, storages, and networking devices.
- To provide the power needed to maintain critical operations.

- To provide the cooling and controlled environment within the parameters needed to sustain operations.
- To provide connectivity to other systems both inside and outside the data center.
- To provide data safety and security as regulation and standard requirements

In the design philosophy of this research, these needs must be met and in the most efficient and reliable way possible.

Enterprises must streamline data center operations while delays and costs to invest, design, and build a data center present challenges to many enterprises. The efficiency of the data center system relies entirely on the efficiency of the system design while reliability depends on system topology or redundant systems. Moreover, sustainable operations of data center confide on their design, operating and maintenance conditions, and human intervention. This research paper proposes the best practices for energy conservation of data center design and operations maintenance. The data center design processes are based on optimal solution of investment, reliability, efficiency, scalability, flexibility, and TCO.

II. BACKGROUND

A. Data Center Topologies

The Uptime Institute Tier Standard: Topology is a basic purpose of concept model that defines taxonomic classification system for comparing the functionality, capacity, and expected availability of a particular site infrastructure as design topology.

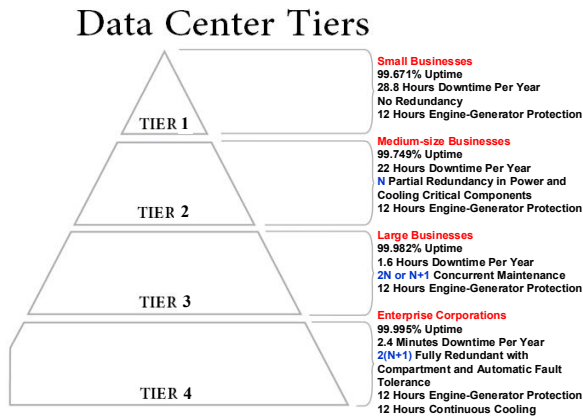


Figure 2. Data center tier classification [6].

This topology describes criteria to distinguish 4 classifications of site infrastructure topology based on increasing levels of redundant capacity components and distribution paths. This standard emphasizes on the definitions of the 4 Tiers and the performance confirmation tests for determining compliance to the definitions. The Tier classifications define the site level infrastructure topology required to sustain data center operations. This topology establishes 4 distinctive definitions of data center site infrastructure Tier classifications [5]: Tier 1: Basic Site Infrastructure; Tier 2: Redundant Site Infrastructure Capacity Components; Tier 3: Concurrently Maintainable Site

Infrastructure; and Tier 4: Fault Tolerant Site Infrastructure, as illustrated in Fig. 2.

B. Energy Efficiency (Power Usage Effectiveness: PUE)

The purpose of PUE is tracking how an individual data center perform over a period. PUE provides a common best practice for comparing and improving power usage in data centers. PUE defines as the ratio of total facility energy draw into data center divided by the total IT equipment energy [10], as illustrated in Table I.

TABLE I. MEASUREMENT LEVEL OF PUE

Where to measure energy point?	Level 1	Level 2	Level 3
How often to measure?	Basic	Intermediate	Advanced
IT Equipment Energy	Required	UPS outputs	PDU outputs
Total Facility Energy	Required	Utility Inputs	Utility Inputs
	Additional recommended measurements	UPS input/outputs	PDU outputs
		Mechanical inputs	inputs/outputs
			Mechanical input
Measurement Intervals	Required	Monthly	Daily
	Additional recommended measurements	Weekly	Hourly
			Continuous

Data collection metering location is defining as either Level 1, Level 2, and Level 3 subject to the definition provides by the Green Grid [10]. The point of measurement of each level represented in Fig. 3.

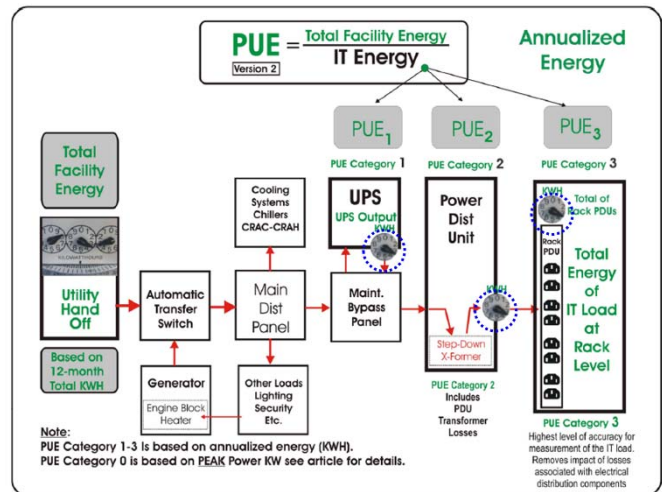


Figure 3. Measurement Poits of PUE.

C. Energy Conversion

Power and cooling limit the number of IT equipment in a data center rather than space. Presumption that all electrical energy used by IT equipment are converting to heat, in real practical situation a few of IT equipment's is converting to noise, vibration, or light. Therefore, 1,000 watts power will generate 3,141 BTUs (British thermal units) of heat [3]. Cooling is normally measured in tons, a unit related to a ton of ice once used for refrigeration. At the conversion rate of 12,000 BTUs/ton of cooling, a full rack 42U IT equipment consumes more than 10 kilowatts of electricity, produces over

34,000 BTUs of heat, and requires almost 3 tons of cooling. If 400 racks were placed in a 1,200 square foot room, 1,200 tons of cooling and 4 megawatts of electricity would be required.

Power and cooling limit the number of IT equipment in a data center rather than space. The relationship among power, cooling, and IT equipment can explain by PUE. The useful power for data center is the power delivered to the IT loads, where does the rest of power convert to? Fig. 4 demonstrates where power flow in each equipment or system of data center [14]. Around 47 percent power flows to IT equipment and the rest 57 percent converts to heat that interprets by $PUE = 2.13$.

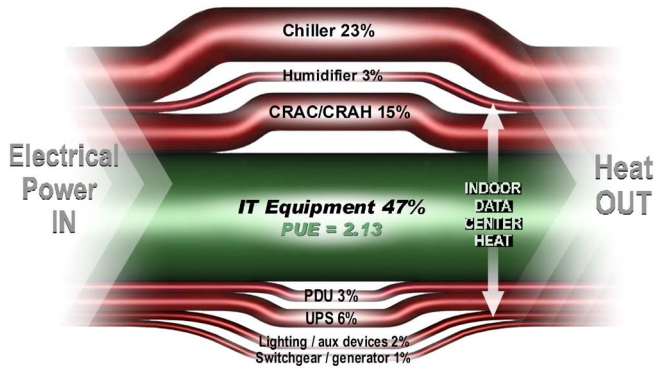


Figure 4. Energy conversion from power to heat.

III. DATA ANALYSIS

In 2020, data from the National Broadcasting and Telecommunications Commission (NBTC) Thailand reports it has 30 Internet data center providers (IDCs). They are only 10 IDCs that operate more than 400 racks, which provide for local and international customers. In order to collect and analyze data, 6 data centers were conducted as research samples, in Bangkok Area. All samples have verified subject to design capacity, occupancy rate, and Tier level, as shown in Table I.

TABLE II. DATA CENTER SAMPLES (IN BANGKOK AREA)

Data Centers	White Space (Sq.m.)	No. of Racks	Occupancy (%)	Design Capacity (Watt/Racks)	Tier Classification
A	4,800	1400	55	10,000	4
B	3,700	1200	30	15,000	3
C	1,700	650	50	5,000	3
D	1,500	500	60	6,000	3
E	1,500	450	35	5,000	3
F	1,200	400	40	5,000	3

Researcher has conducted in depth interview with 6 operation managers of each data center subject to below topics:

- Design capacity and installed capacity
- Expected load and actual load
- Ramp-up time to reach expected load
- Each company growth model

The survey data from Table II shows all data centers have the same issue on oversizing design. The three factors have

been taken for consideration; nameplate, multiplier safety factor (MSF), and redundancy topology.

A. Data Center Life Cycle (DCLC)

A data center project is a unique set of processes composing of coordinated and controlled activities with start and finish dates, undertaken to achieve and objectives [13]. Each data center has designed in different objectives such as location, sizing, risk acceptance or Tier, density operations, energy efficiency, or cost saving. Therefore, design requirements are very important while understood DCLC reduces project distortion subject to more accuracy in requirements transformation to build and operate and maintenance. DCLC composes of 6 fundamental phases [8]: plan, design, build, commission, operate and maintenance, and assess, as illustrated in Fig. 5. Energy conservation in data center starts from planning design requirements need to input from the outset of project plan as one of the key project parameter. To operate more energy efficiency means first design must be proposed for energy efficiency, if not it called cost reduction or saving. Consequently little upfront changes in plan phase will have major cost effect downstream when data center wades the build phase.

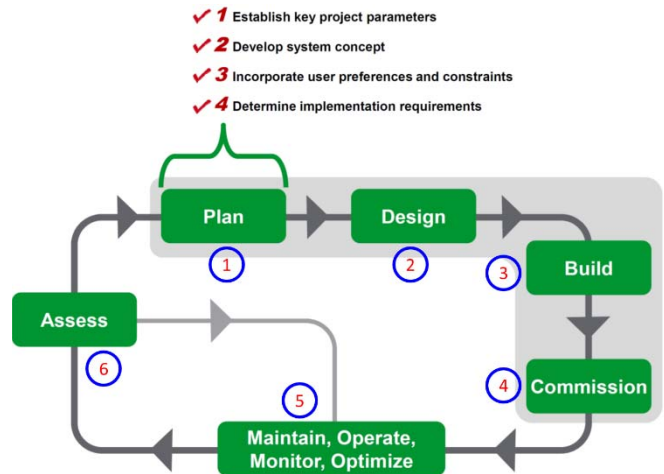


Figure 5. The plan phase of the data center life cycle.

B. Data Center Oversizing Design

The basic concept of data center planning and design is matching the power and cooling requirements of the IT equipment with capacity of infrastructure systems to support its. Sizing of power requirement for data center needs an understanding of the amount of power required by IT equipment, UPS system, and cooling system. The 3 major factors that influence data center oversizing design are nameplate, multiplier safety factor (MSF), and redundancy topology.

The nameplate uses to define the value of the server power supply unit (PSU). Manufacturers derive the nameplate value in term of the DC output rating on the PSU. The nameplate information comprises of input voltage, rated input amperage, kilovolt-ampere (kVA), frequency, and phase information. Nameplate values are accepted as superseding power consumption levels that surpass actual power consumption of IT equipment under normal usage. Regrettably, nameplate

value is a simple method used by data center consultants in the design, planning, and implement of power and cooling systems. Based on consequence of this applications result in immoderate capital expenditures (CAPEX) and operating expenses (OPEX) due to lower efficiencies that result from overestimating power requirements. The research from The Green Grid [11] shows in Fig. 6 that the 51 servers and their PSUs tested results reveal actual power draw of servers around 50 percent only or nameplate/2.

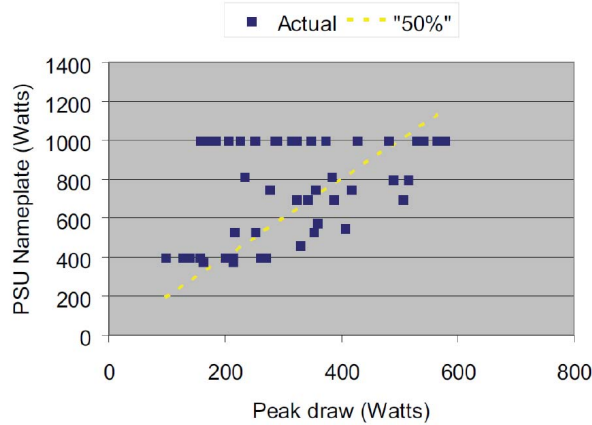


Figure 6. PSU nameplate rating versus peak power draw.

It is impractical to load any power system to its full capacity. A load factor of 90 percent is recommended with a 95 percent maximum loaded, leaving a design of 10 to 5 percent safety factor. For practical designs, also apply the code-mandate safety factor required for continuous load cloud as the design safety factor. If the total speculated data processing load has a capacity criterion of N, the multiplier safety factor (SFM) for each subsystem within the power distribution system (PDS), as demonstrated in Fig. 7, will provision sufficient capacity to encounter normal equipment layout diversity and scalability [12].

Redundancy systems have designed for increasing systems reliability while reduces load on any given equipment, with a corresponding impact on efficiency, CAPEX, and OPEX. Equipment selection should deliberate OPEX at the speculated load levels in addition to CAPEX and space utilization.

C. Total Cost of Ownership (TCO)

A simplified formula of TCO is $TCO = CAPEX + OPEX$. The results from in depth interviews revealed that all data center managers or ICT project managers were very concerned about TCO of data center. Most of them need to concentrate on these factors more carefully:

1) *Location*: or called site selection, this will be directed impact on asset costs or CAPEX on occupied land or may be some project considers as long-term leased. The natural risks and man-made must be taken to consider as the first priority. Location shall be closed to utility substations, fiber optic routes, human supplied facilities, and transportations.

2) *Tier Classification*: Top management or C-suits must make decision to gather on what is the minimal risk acceptance on data center with which Tier that they are

accepted before go to engineering design process, this is first time investment as CAPEX (Tier certified design costs).

3) *Data Center Strategy*: means number of racks and density per rack (kilo Watt/hour) because this must be discussed among marketing strategic team, ICT team, engineering team, and data center consultant before they come up with conceptual desing such as data center layout phasing, number of racks in each phase, and construction plan with buget, time, and Tier Classification.

4) *Sizing*: normally this will be consequence from number of racks and what is the designed density per rack (kilo Watt/hour), this is long-term OPEX;

No. of racks multiply by kW per rack = Total power draw from UPSs that will sizing ATSSs, generators, and transformers respectively. As the same time, the total power draw will use to calculate kW of heat rejection systems of data center, this is first time investment as CAPEX.

No. of racks will sizing number of cabling infrastructure and networking systems, this is first time investment as CAPEX

No. of racks will sizing fire protection system, surveillance system, water leak system, warning system, etc., this is first time investment as CAPEX.

5) *Compliant Standards*: ISO 9000, ISO 14000, ISO 20000, ISO 22301, ISO 27001, ISO 50000, and the Payment Card Industry Data Security Standard (PCI-DSS) are considered as OPEX.

6) *Human Capitals*: these are considered as long-term OPEX subject to daily operations and re-skill training.

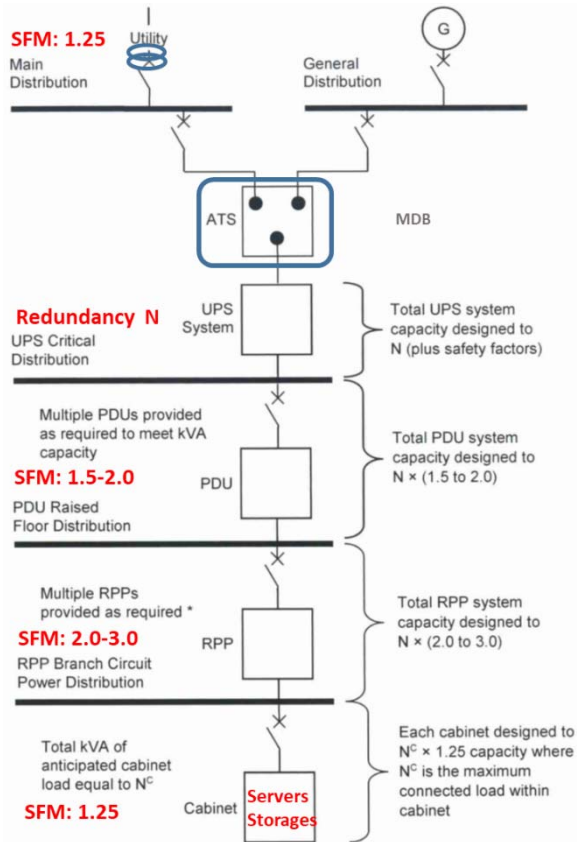


Figure 7. Safety factor multipliers of electrical distribution system components.

IV. CONSTRAINTS AND ASSUMPTIONS OF GORWTH MODEL

The growth model of data center indicates a key role in the chronological order of activity that transforms the conceptual design to detailed design and to physical infrastructure. This planning chronological order is defined in DCLC. Figure 8 shows the context of the growth model within the system planning chronological order.

Early in the planning phase, the user provides the existing IT load profiles and projection growth requirements as input to consultant. The constraints in this phase are uncertainty of IT load and unexpected future growth as projection that will affect to design, investment, and long-term operations.

Later in detailed design phase, the system capacity plan is constructed based on the Tier Classification of reference design for mapping to master plan layout and maximum capacity planning (1). By the growth model expects initial capacity (3) will fulfill maximum capacity planning or ramp-up time within (4) years [14], as shown in Fig. 8.

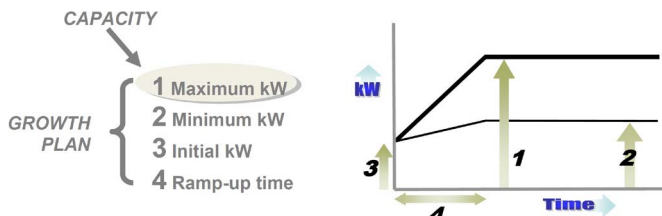


Figure 8. The growth model in the system planning.

In construction phase, contractor must build data center as detailed design from consultant that means design capacity = installed capacity [9], as depicted in Fig. 9.

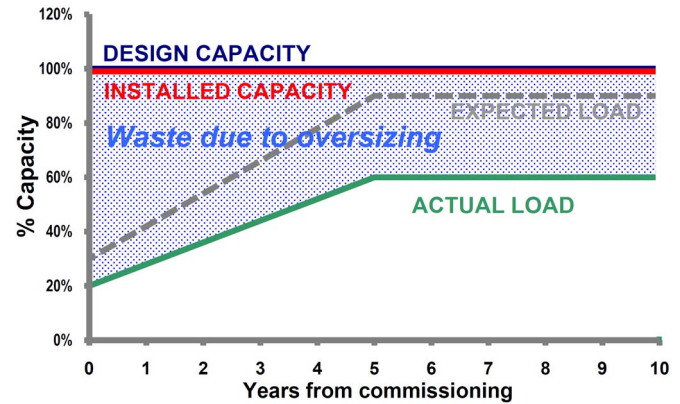


Figure 9. Design capacity equal to installed capacity.

A. Scenario Fixed CAPEX & Variable OPEX

This system design called non-scalable concept that is installed at the outset to accommodate the maximum load anticipated during the data center lifetime. The capacity of the power and cooling systems has installed as the design capacity or completely built-out from the beginning. The design planning is expected the data center load of IT equipment will start at 30 percent and ramp-up to expected load of data center within 5 years. Typically, the actual start-up load of data center operations is uncertainty normally is lower than expected load. This data has been confirmed after in depth interview with 6 operations' manager of each IDSs. The gap between expected and actual load ranges from 10-40 percent.

Case studies of non-scalable elements are the physical room size or building, electrical service entrance capacity, and pre-existing room-based air conditioning. This scenario is usually considered under situation that is not additional costs of CAPEX after operations and maintenance. This scenario affects huge of energy waste due to oversizing. The consequence of waste due to oversizing effects a long-term operating costs or OPEX as long as data center operations. On the other hand, this scenario will be feasible if actual load can reach up expected load within 5 years otherwise the return on investment, (ROI) will take longer than expectation.

B. Scenario Variable CAPEX & Variable OPEX

This system design called scalable concept that is installed, at the outset, for a lower-than-maximum load then increased over time according to the steps of the phase-in plan, as presented in Fig. 10. This new technology called modular. Modular data center systems consist of purpose-engineered modules and components to offer scalable data center capacity with multiple power and cooling options. Modular Data Center is a flexible, efficient, agile, scalable and compact modular data center can be deployed virtually at any location. This concept design reduces the gap between installed capacity and actual load that IT manager and consultant are already perceived during design phase. This helps investor prolong CAPEX as the same time save long-term operating costs or OPEX throughout the lifetime of data center operations. Examples of scalable elements are racks or

containers, rack or container-based power protection and distribution, and rack-based cooling equipment [7].

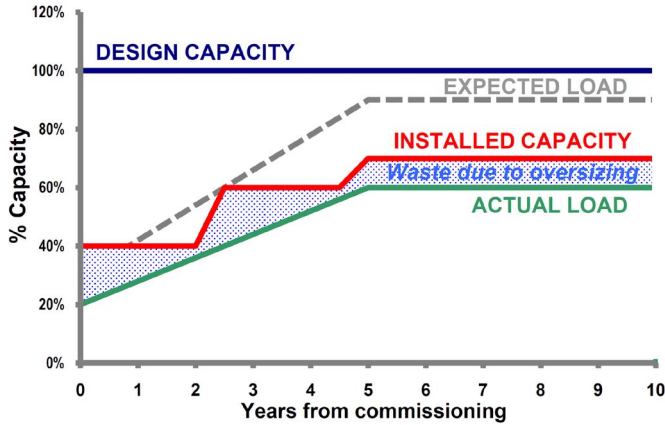


Figure 10. IT load certainty as full upfront buildout.

They two major problems of legacy data centers are investment and speed of deployment. Many data center industry took to modular approaches as a solution to get funding approved in smaller amounts and mitigate risks of long period building a data center. The other side of those two major drivers there are many rationales and advantages listed for why a modular data center approach is selected for operations and maintenance.

- **Speed of deployment:** The modular solutions have designed as prefabricated that why they are unbelievably fast timeframes from order to installation.
- **Scalability:** As the concept of repeatable, flexible, and standardized design, it is easy to match demand sizing of any data center and volume scale out infrastructure quickly.
- **Efficiency:** Modular is engineered product design that intends for matching for IT power requirements with modular power and cooling plants while save money in costly power distribution and energy loss from being so close. Within modular package, they are applied energy management platforms called data center infrastructure management (DCIM) as core operations.
- **Total cost of Ownership:** Growth as your need, installation within a few months rather than 12 months
- **Operations:** Standardized modules of components and systems with standard operations and maintenance procedures

C. Energy Waste in Operations

As normal data center operations, power losses associated with energy conversion on USPs and PDSs currently represent around 10 percent of data center's total energy consumption, while the cooling system represents 32 percent. Given the efforts made with regard to cooling, particularly through designs enabling the use of free cooling, these power losses will account for a large part of the data center's energy bill.

In order to increase the energy efficiency of the electrical infrastructure, it is therefore necessary to look into solutions that will reduce this percentage, in particular via the power supply and distribution systems. As PUE is a measure of power efficiency in data center. There are two major portions of energy consumptions; IT equipment 59 percent and cooling system 30 percent. To improve PUE from 1.7 to 1.3 needs more things to do on IT strategy, as depicted in Fig. 11. If PUE equal 2.0 that means half the energy goes for IT equipment and half for other tasks goes for waste.

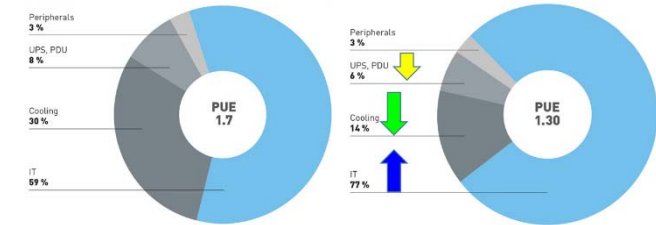


Figure 11. Reduce power losses in cooling system.

According to IT equipment, virtualization enables server to run multiple operating systems on the hardware of a single physical server (utilization IT load), while consolidation enables you to deploy multiple applications using the same operating system on a single virtual machine or server (reduce energy loss and increase CPU utilization). Moreover, Google and Facebook both use customized their's servers. They were designed to take out unnecessary components like graphic cards to minimize power consumption at the power supply unit (PSU) and minimize power loss of voltage inverter.

Ordinarily data centers are designed to only provide a few amount of outside air the data center. The new design of free cooling such as Google and Facebook begin to operate more on outside cooling air in their data center. The idea is built their data centers in countries which naturally has cold and dry climates to save their long-term OPEX. Facebook built a data center in north Sweden which physically has cold and dry weathers. Google built data center in Belgium eliminates chillers altogether, deployment evaporative cooling system that run on free cooling 100 percent.

V. CONCLUSION

As demand of power consumption in data center industry is around 2 percent of world's power consumption. Therefore, energy management must be deployed to cut loss of power consumption. Measurement the energy efficiency of data center by PUE is a guideline for improving waste due to design oversizing. This research paper reveals cause of data center designed oversizing and solutions to resolve it. To understand from the outset of data center design for energy efficiency is the right thing must do at the first time. The consequence of oversizing design data center affects investment or CAPEX and long-term operating costs or OPEX. Researcher has advised that modular data center will save TCO according to the reasons and benefits provides by modular. Google and Facebook were represented as the best practices for energy conservation and energy efficiency.

REFERENCES

- [1] E. R. Masanet, et al., *Global Data Center Energy Use: Distribution, Composition, and Near-Term Outlook*, Evanston, IL, 2018.
- [2] International Energy Agency, "Data centres and data transmission networks: Tracking clean energy progress," [online]. Available: <https://www.iea.org/tcep/buildings/datacentres/>, May 2019.
- [3] T. H. Payne, *Practical Guide to Clinical Computing Systems: Design, Operations, and Infrastructure*, Elsevier, 2008, p. 39-42.
- [4] A. S. G. Andrae and T. Edler, "On Global Electricity Usage of Communication Technology: Trends to 2030," *Challenges*, 6(1), pp. 117-157, April 2015.
- [5] Uptime, *Data Center Site Infrastructure Tier Standard: Topology*, Uptime Institute, LLC., 2018.
- [6] M. Wiboonrat, "Human Factors Psychology of Data Center Operations and Maintenance," in *The 6th International Conference on Information Management (ICIM 2020)*, Imperial College London, UK, March 27-29, 2020.
- [7] N. Rasmussen, "Avoiding costs from oversizing data center and network room infrastructure," White Paper 37, Rev 7, Schneider Electric., 2012.
- [8] N. Rasmussen, "Data Center Projects: System Planning," White Paper 142, Rev 2, Schneider Electric., 2013.
- [9] N. Rasmussen and S. Niles, "Data Center Projects: Growth Model," White Paper 143, Rev 1, Schneider Electric., 2011.
- [10] The Green Grid, *Usage and Public Reporting Guidelines for the Green Grid's Infrastructure Metrics (PUE/DCiE)*, The Green Grid, V.2.1, October 29, 2009.
- [11] The Green Grid, *Proper sizing of IT power and cooling loads white paper, White Paper #23*, The Green Grid, V. 1.0, 2009.
- [12] ANSI/BICSI 002-2019, *Data center design and implementation best practices*, BICSI, May 2019.
- [13] ISO 21500, *Guidance on project management, 1st Ed.*, ISO, Sep 2012.
- [14] E. Kotzbauer and D. Bouley, "Guide for reducing data center physical infrastructure energy consumption in federal data centers," White paper 250, Schneider Electric., 2011.