

# Joint Capacity Planning and Operational Management for Sustainable Data Centers and Demand Response

Tan N. Le, SUNY Korea & Stony Brook University  
 Zhenhua Liu, Stony Brook University  
 Yuan Chen, Hewlett Packard Labs  
 Cullen Bash, Hewlett Packard Labs

Reducing costs plays a crucial role in building and operating data centers. Internet service providers such as Facebook and Google spend billions of dollars on capacity expansion and operations of their global data centers. Traditionally, capacity planning for data centers is done separately from operational management, which incurs inefficiency. In fact, operational management has significant impacts on capacity planning. Motivated by this gap, we propose a framework that jointly optimizes both capacity planning and operational management for sustainable data centers and data centers participating in demand response programs. Numerical results based on real-world cases highlight that the proposed framework remarkably reduces up to 50% of total expenditures and 75% of greenhouse gas emissions compared to conventional methods. Additionally, our results show that participations in various demand response programs result in vastly different capacity planning decisions and lead to emission reductions of up to 60%.

CCS Concepts: •**Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; •**Networks** → Network reliability;

General Terms: Design, Algorithms, Performance

Additional Key Words and Phrases: Sustainable data centers, demand response

## ACM Reference Format:

Tan N. Le, Zhenhua Liu, Yuan Chen, Cullen Bash, 2016. Joint Capacity Planning and Operational Management for Sustainable Data Centers and Demand Response. *ACM Trans. Embedd. Comput. Syst.* 9, 4, Article 39 (June 2016), 21 pages.  
 DOI: 0000001.0000001

## 1. INTRODUCTION

The total cost of ownership (TCO) and greenhouse gas (GHG) emissions of data centers are exponentially increasing [Rana 2009] due to the explosive demand for using Internet services. Giant cloud providers like Google and Facebook spend billions of dollars every quarter on their data centers [Sverdlik 2015]. On the other hand, data centers are under pressure to reduce their emissions. Conventional data centers are mainly powered by the electricity grid that heavily depends on fossil fuel. In fact, a data center can release emissions equivalent of hundred thousands of cars [EPA 2005; Matt 2009].

The TCO of a data center is mainly the capital expense (CapEx) and the operational expense (OpEx) [Barroso et al. 2013]. CapEx of a data center is the costs that must be invested up front and then depreciated over a certain time frame, e.g., the construction cost of a data center and the purchase of servers. OpEx refers to the costs of operating

---

This work is supported by the National Science Foundation, under grant CNS-0435060, grant CCR-0325197 and grant EN-CS-0329609.

Author's addresses: Tan N. Le and Zhenhua Liu, Computer Science Department, Stony Brook University; Yuan Chen and Cullen Bash, Hewlett Packard Labs.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2016 Copyright held by the owner/author(s). 1539-9087/2016/06-ART39 \$15.00

DOI: 0000001.0000001

a data center, including electricity costs, software and hardware maintenance, and repairs, salaries for human resources, etc.

The CapEx of data centers is actually an interesting topic. For instance, renewable energy is normally considered to be free. However, the cost of deploying a renewable power plant is far more expensive than building a traditional power plant [Abdmouleh et al. 2015]. Fortunately, the price of renewable energy equipment is going down in the long-run with improved technologies. On the other hand, renewable energy sources are not dispatchable because their generations heavily depend on weather and geographical conditions. Such issues are critical in efficient use of renewable energy sources.

The OpEx of a data center is becoming more dynamic than ever as a result of increasing and abundant work focusing on power demand management. For instance, the workload in data centers can be shaped to achieve a given objective, such as minimizing electricity cost [Urgaonkar et al. 2011; Pakbaznia and Pedram 2009]. Furthermore, the workload demands can even be balanced among the geographically distributed data centers [Liu et al. 2011; Qureshi et al. 2009].

The relationship between CapEx and OpEx in a data center is bi-directional. CapEx and OpEx are the costs associated with the capacity planning and operational management, respectively. Capacity planning is the process of determining the infrastructure for a data center. In order to have efficient capacity planning, it is necessary to consider how the data center would operate in the long-run. Meanwhile, capacity planning may have significant impacts on OpEx because the cost of operational management varies under different settings. Traditionally, a data center is built to serve the peak workload demands [Ren et al. 2012]. However, it may lead to over-provisioning and cause a huge waste of capital and maintenance costs since the workload can actually be shaped to reduce the peak. On the other hand, under-provisioning would not meet the quality of service (QoS) requirements, e.g., latency, which severely debilitates the business.

In this paper, we propose an optimization framework for joint capacity planning and operational management in Section 3. The optimization framework is to minimize both CapEx and OpEx. Meaning, the framework can deal with the inter-dependency of capacity planning and operational management to outperform the traditional methods. To evaluate the proposed framework, we carry out numerical simulations for two scenarios: sustainable data centers and data center demand response.

Sustainable data centers (SDC) [Weihl et al. 2011] are designed to reduce costs and GHG emissions. For example, HP designed “Net-zero Energy Data Centers”, which utilize the various renewable energy sources to reduce energy cost as well as emissions [Arlitt et al. 2012]. Furthermore, sustainable data centers may integrate some of advanced power management techniques, such as server consolidation [Lin et al. 2013; Zhang et al. 2012; Lin et al. 2011], network consolidation [Zhang et al. 2010; Andrews et al. 2012; Sharma et al. 2015], colocation of workloads [Aksanli et al. 2012], cooling power optimization [Liu et al. 2012; Pakbaznia and Pedram 2009], batch job scheduling [Mukherjee et al. 2009; Garg et al. 2011], and using energy storages [Urgaonkar et al. 2011; Liu et al. 2012; Liu et al. 2013]. How to optimize the design of sustainable data centers is still challenging as they are far more complicated than that of traditional data centers.

The participations of data centers in demand response (DR) programs can potentially contribute to the electricity grid [Wierman et al. 2014; Liu et al. 2014]. In fact, data centers are large loads and can be considered as giant virtual batteries to help improve the reliability of electricity grid [Wierman et al. 2014]. Despite such great potential, lots of important questions still remain open. How do the participations in DR programs affect a data center in terms of cost, capacity planning, and power management? How well does a data center respond to the DR signals? What are the

subsequences of the changes, e.g., in GHG emissions? We address these questions in Section 5.

**Our contributions are three-fold.**

*First, we develop an optimization framework for joint capacity planning and operational management* in Section 3. The optimization framework is based on the model of sustainable data centers and general enough to be applied to traditional data centers. The model includes multiple power demand and supply components, i.e., IT workload demand, cooling power, renewable energy sources, non-renewable energy sources, and electricity grid. The joint optimization framework provides an optimal capacity planning decision to construct, expand and operate the data center annually. In addition, the model can estimate the emissions of a data center. As the framework requires predictions for capacity planning in the long-run, prediction errors are incorporated. Moreover, we extend the optimization framework to include Net-Zero Energy Data Centers and data center demand response (DCDR) in Section 4.4 and Section 5, respectively. Unlike conventional data centers, Net-zero Energy Data Centers (NEDC) can be run by stand-alone micro-grids mainly powered by renewable resources [Arlitt et al. 2012].

*Second, we evaluate the proposed framework on sustainable data centers* in Section 4. The evaluation is based on the real design of a data center, EcoPOD designed by HP [hpE 2014]. The data center can provision power from photovoltaic (PV) generation, gas engine (GE) generation, and electricity grid.

- *We highlight the benefits of using our proposed framework* in Section 4.2. We compare the proposed framework with three baseline methods. The comparisons demonstrate that the proposed framework achieves up to 50% of cost savings and 75% of emission reductions. Additionally, the simulation results in annual capacity planning show that the proposed framework tends to increase the use of renewable energy and decrease emissions over time.
- *We study the impacts of prediction errors on our proposed framework* in Section 4.3. Under large prediction errors, the proposed framework still achieves significant cost savings and emission reductions.
- *We provide sensitivity analysis on the proposed framework for a NEDC* in Section 4.4. As NEDC are mainly powered by the local energy resources, the framework is extended to include a net-zero energy constraint. We study various factors, i.e., electricity price, gas price, shape of interactive workload, and ratio of flexible workload. This analysis provides lots of interesting insights. For instance, there are trade-offs between PV and GE. Additionally, while the high ratio of flexible workload has very positive impacts on using more PV, the shapes of interactive workload affect little on the capacity planning and operational management of NEDC.

*Last but not least, we evaluate data center's capacity planning and operational management when participating in demand response programs.* Extensive numerical simulations in Section 5 show that this results in different capacity planning decisions, and some of them reduce emissions up to 60%. Moreover, we demonstrate that the proposed framework allows data centers to adapt to each DR program very well.

## 2. BACKGROUND AND PRIOR-WORK

### 2.1. Sustainable data centers

The burdens of financial costs, energy resources, and emissions have been heavily put on data centers [Koomey 2008]. Thus, the concept of sustainable data centers has been defined to cut the electricity usage, utilize renewable energy sources, and reduce emissions [Weihl et al. 2011]. A sustainable data center can be powered by multiple energy

resources, while renewable energy sources, such as photovoltaic (PV) and wind, are preferred.

From social perspectives, reducing GHG emissions becomes critical. Various countries are developing the policies and regulation on emissions. For example, cap and trade is the government-mandated and market-based approach that controls pollution by providing economic incentives for achieving reductions in emissions [Stavins 2003]. The violation of regulated emission caps may lead to penalties [Revelle 2009]. Therefore, such heavy power loads like data centers are under pressure to reduce their emissions.

**What we model and study:** We model sustainable data centers in Section 3. In our evaluation, we study our proposed optimization framework on a sustainable data center in terms of costs and emissions. Furthermore, we extend the evaluation to a Net-zero Energy Data Center, a special case of sustainable data centers.

## 2.2. Data center power management

The major operational cost is dependent on data center power management. A data center power management scheme, run by human or computer program, strives to reduce the costs and emissions. Data center power management can be divided into multiple topics, such as server consolidation [Lin et al. 2013; Zhang et al. 2012; Lin et al. 2011], network consolidation [Andrews et al. 2012; Zhang et al. 2010; Sharma et al. 2015], colocation of workloads [Aksanli et al. 2012], cooling power optimization [Liu et al. 2012; Pakbaznia and Pedram 2009], batch job scheduling [Mukherjee et al. 2009; Garg et al. 2011], geographical load balancing [Qureshi et al. 2009; Liu et al. 2011], and using energy storages [Urgaonkar et al. 2011; Liu et al. 2012; Liu et al. 2013]. Now, we discuss colocation of workloads in more details.

*Server consolidation:* Huge amounts of power are wasted due to the large idle power consumption of servers and low utilization. Therefore, significant power can be saved by consolidating workloads on the right amount of servers and switching the remaining servers to low power modes or even turning them off [Zhang et al. 2012; Lin et al. 2011; 2013].

*Network consolidation:* Thousands of network devices can also be switched off to save the power consumption. The key idea is to turn off the unused network devices, such as the switches connected to the servers that have been turned-off [Andrews et al. 2012; Zhang et al. 2010; Sharma et al. 2015].

*Colocation of workloads:* Currently, the two types of workload, i.e. interactive workloads and batch jobs, are usually served on different servers, making it difficult to save power through consolidation. On the other hand, the power saving potential is great because there are different resource and performance requirements for interactive workloads and batch jobs. A promising way is to run interactive workloads with high priority that keeps its performance, e.g., mean/percentile response time, (almost) unaffected while running batch jobs whenever there is spare capacity to use. This can significantly increase the server utilization and therefore save power [Aksanli et al. 2012].

*Cooling optimization:* While a large amount of power is used for keeping data centers under certain thermal constraints through cooling systems, there is great potential to reduce cooling power by optimizing the data center to use the most effective cooling in the right amount at the right time [Liu et al. 2012; Pakbaznia and Pedram 2009].

*Batch job scheduling:* The flexibility in batch jobs provides great temporal flexibility for scheduling to shape the demand. A smart scheduler can run batch jobs in the right amount at the right time to make demand more supply following [Mukherjee et al. 2009; Garg et al. 2011].

*Geographical load balancing:* When an interactive workload is served by an Internet-scale system having data centers at different locations, a spatial flexibility emerges. A global load balancer can route interactive workload request to the right data center to better align demand with supply [Qureshi et al. 2009; Liu et al. 2011].

*Energy storage:* Energy storages can be used to save energy cost by charging when supply exceeds demand or supply is cheap and discharge in the future to power the data center, which can better align power supply with demand. However, the current high cost of energy storage usually prevents large deployment, e.g., using energy storage to power the whole data center for several hours. Instead, the current practice just uses energy storages in UPS (Uninterrupted Power Supply) as a transit from the electricity grid to backup generators, which can last several minutes to tens of minutes depending on the ramp up the speed of the backup generators [Urgaonkar et al. 2011; Liu et al. 2012; Liu et al. 2013].

**What we model and study:** We incorporate one of the aforementioned techniques, i.e., colocation of workloads, into our framework, which is general enough to include other techniques. In our evaluation, we study that important role of power demand management in reducing costs and emissions.

### 2.3. Data center demand response (DCDR)

Demand response (DR) programs are defined to improve the traditional electricity markets and grids. There are generally two types of participation in DR programs, which are passive and active participations [Wierman et al. 2014].

- Passive participation: Passive participation is typical in the smart pricing services. The pricing services issue price signals to encourage electricity users to adjust their power consumption profiles. In electricity markets, there are multiple pricing services, i.e. Time-of-Use (ToU), Inclining Block Rates (IBR), Peak Pricing (PP), Coincident Peak Pricing (CPP), Day-ahead Pricing (DaP), and Real-Time Pricing (RTP). For example, peak-pricing charges a high price at peak demand to prevent power outages.
- Active participation: Active participation is diverse. Customers can use the wholesale markets, ancillary services, or voluntary reduction programs. A wholesale market allows data centers to purchase electricity directly from power suppliers instead of regional retailers. Ancillary services are defined to maintain reliable operation and security of the electricity transmission system. The basic idea is to encourage customers to adjust their loads due to the condition of the electricity grid. In voluntary reduction programs, customers can have flexible contracts with grid operators for offering services.

There is a high potential that large power loads like data centers participate in demand response programs. In addition to the flexibility of power demand, data centers can utilize their energy storages as well as UPS to increase the flexibility of their loads during the demand response events [Liu et al. 2014; Wierman et al. 2014].

**What we model and study:** In Section 5, we model and study the participation of data centers in DR programs. We incorporate and evaluate several popular DR programs, i.e., ToU, CPP, IBR, an ancillary service, and wholesale markets. We conduct simulations to answer the following questions: How do DR programs change the power profile of data centers? How do the DR programs impact on the capacity planning and operational decisions?

### 3. OPTIMIZATION FRAMEWORK

#### 3.1. Modeling sustainable data centers

We consider the problem of capacity planning and operational management for sustainable data centers in  $Y$  years, where each year is discretized into  $T$  time slots. Since the data center can be expanded annually, we model the data center for each year  $y \in \{1, 2, \dots, Y\}$  as follows.

**Power demand.** Power demand is mainly from two subsystems: IT subsystem and cooling subsystem [Barroso et al. 2013]. The IT subsystem serves the IT workloads, i.e., interactive workload and flexible workload (batch jobs). The cooling subsystem reduces the heat generated by the IT equipment to keep the inner temperature in an acceptable range. We use the model of interactive workloads and batch jobs similar to [Liu et al. 2012].

*Interactive workload demand:* There are  $N$  interactive workloads. For interactive workload  $i$ , we assume that the IT power is allocated to interactive workload  $i$  at time  $t \in \{1, 2, \dots, T\}$  of year  $y$ , denoted by  $a_i(y, t)$ . Here  $a_i(y, t)$  can be derived from either analytic performance models [Urgaonkar et al. 2005] or real-world data traces.

*Batch job demand:* Batch jobs are the sequence of computer commands to be processed. We assume there are  $J$  classes of batch jobs. Class  $j$  has total power demand  $B_j(y)$ , starting time  $S_j$ , and deadline  $E_j$ . Let  $b_j(y, t)$  denote the amount of capacity allocated to class  $j$  jobs at time  $t$  of year  $y$ . Hence,  $b_j(y, t)$  can be allocated such that

$$\sum_{t=S_j}^{E_j} b_j(y, t) = B_j(y) \quad \forall y, j. \quad (1)$$

The total IT power demand  $P_{IT}(y, t)$  at time  $t$  of year  $y$  is then computed as  $P_{IT}(y, t) = P_{idle}(y, t) + \sum_{i=1}^N a_i(y, t) + \sum_{j=1}^J b_j(y, t)$ , where  $P_{idle}(y, t)$  is the idle power consumption of the data center, which can be computed based on the number of active servers [Lin et al. 2013].

The amortized infrastructure cost of IT subsystem per Watt per year is  $I_{IT}(y)$  (\$/W). The operational and maintenance cost at time  $t$  is  $p_r(y, t)$ . Let  $C_{IT}(y)$  be the capacity of IT subsystem at time  $t$  of year  $y$ . The total IT power demand  $P_{IT}(y, t)$  is capped by

$$P_{IT}(y, t) \leq C_{IT}(y), \quad \forall y, t. \quad (2)$$

Using power usage efficiency (PUE) [Barroso et al. 2013], the total power demand  $P(y, t)$  at time  $t$  of year  $y$  is

$$P(y, t) = PUE(y, t) * P_{IT}(y, t),$$

where  $PUE(y, t)$  is the PUE at time  $t$  of year  $y$ . Here, the power demand of the cooling subsystem is  $(PUE(y, t) - 1) * P_{IT}(y, t)$ .

**Power supply.** At supply side, we model renewable generation, non-renewable generation, the electricity grid, and energy storages.

*Renewable generation (RG).* A data center may have  $R$  renewable energy sources, e.g., on-site PV panels, on-site/off-site wind farms, etc. The amortized infrastructure cost of source  $r$  per Watt per year is  $I_r(y)$  (\$/W). The operational and maintenance cost at time  $t$  is  $p_r(y, t)$ . Let  $C_r(y)$  denote the capacity of RG  $r$  in year  $y$ . So, let  $c_r(y, t)$  be the power generation of RG  $r$  at time  $t$  of year  $y$ . The renewable generation  $c_r(y, t)$  is often uncontrollable and formulated as  $c_r(y, t) = CF_r(y, t) * C_r(y)$ , where  $CF_r(y, t)$  is the capacity factor at time  $t$  of year  $y$ , which is the ratio of actual output to the potential output.

*Non-renewable generation (NG).* A data center may have  $S$  non-renewable sources, e.g., gas engines. The amortized infrastructure cost of source  $s$  per Watt per year is  $I_s(y)$  (\$/W). The operational and maintenance cost at time  $t$  is  $p_s(y, t)$ . Let  $C_s(y)$  denote

the power capacity of NG  $s$  in year  $y$ . So, the power generation  $c_s(y, t)$  of NG  $s$  at time  $t$  of year  $y$  satisfies

$$c_s(y, t) \leq C_s(y), \quad \forall y, t. \quad (3)$$

**Electricity grid.** At time  $t$  of year  $y$ ,  $p_g(y, t)$  and  $p_b(y, t)$  respectively denote the electricity usage based charging price, (\$/kWh) and the sell back price (\$/kWh). The sell back price is applied when the data center sells their unused local generation back to the electricity grid. At time  $t$  of year  $y$ , the grid power consumption is  $c_g^+(y, t)$  and the sell-back power is  $c_g^-(y, t)$ . Let  $C_g(y)$  be the power capacity of the electricity grid in year  $y$ . In fact, this is usually set to the maximum grid power of the data center since the infrastructure cost is relative small compared with the other utility charges [Barroso et al. 2013]. In particular,

$$C_g(y) = \max\{c_g^+(y, t)\}_{t \in \{1, 2, \dots, T\}}, \quad \forall y, \quad (4)$$

where  $c_g^+(y, t)$  and  $c_g^-(y, t)$  are both non-negative.

**Energy storages.** The total capacity of the energy storage at year  $y$  is  $C_e(y)$ .

**Emissions.** Emissions are from RG sources, NG sources and the electricity grid. In year  $y$ , the emissions rates of RG source  $r$ , NG source  $s$ , and the electricity grid are  $q_r(y)$ ,  $q_s(y)$ , and  $q_g(y)$ , respectively. We do not impose emission cap as it is still regional. However, and our model is general enough include a constraint of emission cap.

Table I. The description of important notations in year  $y$ .

	Symbol	Description
IT	$a_i(y, t)$	Interactive workload power demand at time $t$ of year $y$
	$B_j(y)$	Total batch job workload power demand in year $y$
Prices	$p_g(y, t)$	Electricity price at time $t$
	$p_b(y, t)$	Sell-back price at time $t$
Infra.	$I_{IT}(y)$	Amortized cost of IT and cooling subsystems
	$I_r(y)$	Amortized cost of RG $r$
	$I_s(y)$	Amortized cost of NG $s$
O&M	$p_r(y, t)$	O&M cost of RG $r$ at time $t$
	$p_s(y, t)$	O&M cost of NG $s$ at time $t$
Emissions	$e_r(y)$	Emissions rate of RG $r$
	$e_s(y)$	Emissions rate of NG $s$
	$e_g(y)$	Emissions rate of electricity grid

**Prediction errors:** Since our proposed framework does capacity planning for long-term data center operation, it requires predictions of workload demand, renewable generation, and electricity prices. In practice, prediction errors are inevitable. At time  $t$  of year  $y$ , the prediction errors of interactive workload, batch jobs, capacity factor, electricity prices, sell-back prices, and the O&M cost of NG sources are  $\delta_a(y, t)$ ,  $\delta_b(y)$ ,  $\epsilon_r(y, t)$ ,  $\rho_g(y, t)$ ,  $\rho_b(y, t)$ , and  $\rho_s(y, t)$ , respectively, such that

$$\begin{aligned} \delta_a(y, t) &= a_i(y, t) - \hat{a}_i(y, t), \\ \delta_b(y) &= B_j(y) - \hat{B}_j(y), \\ \epsilon_r(y, t) &= CF(y, t) - \hat{C}F(y, t), \\ \rho_g(y, t) &= p_g(y, t) - \hat{p}_g(y, t), \\ \rho_b(y, t) &= p_b(y, t) - \hat{p}_b(y, t), \\ \rho_s(y, t) &= p_s(y, t) - \hat{p}_s(y, t), \end{aligned}$$

where  $\hat{a}_i(y, t)$ ,  $\hat{B}_j(y)$ ,  $\hat{C}F(y, t)$ ,  $\hat{p}_g(y, t)$ ,  $\hat{p}_b(y, t)$ , and  $\hat{p}_s(y, t)$  are respectively the predicted values of interactive workload, batch job, capacity factor, electricity price, sell-back price, and O&M cost of NG source  $s$  at time  $t$  in year  $y$ . We do not consider the prediction errors of O&M cost for RG sources since they are very small and usually stable for a long-time.

### 3.2. Optimization problem formulation

Table II. Summary of objective components.

	Expression
<i>UtilBill</i>	$\sum_{y=1}^Y \sum_{t=1}^T (\hat{p}_g(y, t) c_g^+(y, t) - \hat{p}_b(y, t) c_g^-(y, t))$
<i>RGEx</i>	$\sum_{y=1}^Y (\sum_{r=1}^R (I_r(y) C_r(y) + \sum_{t=1}^T p_r(y, t) c_r(y, t)))$
<i>NGEx</i>	$\sum_{y=1}^Y (\sum_{s=1}^S (I_s(y) C_s(y) + \sum_{t=1}^T \hat{p}_s(y, t) c_s(y, t)))$
<i>ITEx</i>	$\sum_{y=1}^Y I_{IT}(y) C_{IT}(y)$

**Objective function.** The objective function includes costs from both supply and demand sides. The power supply cost is the predicted CapEx and OpEx of purchasing the electricity (*UtilBill*) and using distributed generations (*RGEx* and *NGEx*). At the demand side, there are predicted CapEx the IT subsystem (*ITEx*). Hence, the objective function of operational cost is defined as follows.

$$OPT : UtilBill + RGEx + NGEx + ITEx.$$

Table III. Summary of decision variables in year  $y$ .

	Symbol	Description
Capacity planning	$C_g(y)$	Grid power capacity
	$C_r(y)$	Capacity of RG $r$
	$C_s(y)$	Capacity of NG $s$
	$C_{IT}(y)$	IT capacity
Operational management	$c_g^+(y, t)$	Grid power usage at time $t$
	$c_g^-(y, t)$	Grid sell-back power at time $t$
	$c_s(y, t)$	Output of DG $s$ at time $t$
	$b_j(y, t)$	Power for batch job $j$ at time $t$

**Decision variables.** There are two types of decision variables which are capacity planning and operational management.

- *Capacity planning variables* are the capacities of IT subsystem,  $C_{IT}(y)$ , RG source  $r$ ,  $C_r(y)$ , NG source  $s$ ,  $C_s(y)$ , and the electricity grid,  $C_g(y)$ .
- *Operational management variables* decide (i) how much electricity would be imported from the electricity grid,  $c_g^+(y, t)$ ? (ii) how much electricity would be sold to the electricity grid,  $c_g^-(y, t)$ ? (iii) How much energy would be generated by NG source  $s$ ,  $c_s(y, t)$ ? (iv) How much power is allocated to serve batch job,  $b_j(y, t)$ ?

The summary of decision variables is in Table III.

#### Constraints.

*Supply-demand balance.* To prevent the data center from power outages, the total supply generation is always greater than or equal to the total power demand as

$$\sum_{r=1}^R c_r(y, t) + \sum_{s=1}^S c_s(y, t) + c_g^+(y, t) - c_g^-(y, t) \geq P(y, t), \quad \forall y. \quad (5)$$



**Capacity caps.** Under capacity planning, the capacities of IT  $C_{IT}(y)$ , electricity grid  $C_g(y)$ , and non-renewable distributed generation  $C_s(y)$  cannot be violated as in (2), (3), and (4), respectively.

**Batch job deadlines.** As the batch jobs  $b_j$  have to be completed during the starting time  $S_j$  and the ending time  $E_j$ , the constraint (1) is included.

**Computational complexity.** The objective function of the framework is actually linear on the decision variables. In addition, the aforementioned constraints are also linear except the constraint for the maximum grid power (4), which can be easily converted into a set of linear constraints. Thus, the framework can be efficiently solved by using a linear programming tool. In particular, we use CVX [Grant et al. 2008] to solve the optimization for our simulation.

## 4. EMPIRICAL EVALUATION FOR SUSTAINABLE DATA CENTERS

### 4.1. Experimental setup

We carry out evaluation based on the settings of a HP EcoPOD data center [hpE 2014].

**Demand side.** The power demand capacity of the data center is 1MW. The power demand is from both IT equipment serving both interactive and batch workloads, and from cooling facilities removing the heat. The peak power usage is 720kW. The Peak-to-Mean Ratio (PMR) of interactive workload is set at 3. Flexible workloads (batch jobs) are 50% of the total IT workload demand with flexibility 24 hours. The utilization of interactive workloads in a server is 40%, and the maximum utilization for all workloads is 90%. PUE is set at 1.2. We study the PMR and flexible workload ratios in terms of capacities and costs in Figure 8 and 9.

**Supply side.** We consider a power micro-grid to supply the data center. The power micro-grid consists of on-site photovoltaic (PV) array, general electric (GE) natural gas engines, and the electricity grid (grid).

**PV array:** The amortized infrastructure cost of PV array after rebate is \$2.15/W [CleanTechnica 2015]. The maximum size of PV array is 1MW. The operational cost and maintenance cost of PV array are respectively \$0/kWh and \$0.005/kWh. The PV capacity factor average is 18% which is from the trace of Houston PV generation.

**GE engines:** The amortized infrastructure cost of GE engines is \$1/W. Maximum size of GE engines can be installed is 1.4MW. The operational cost is \$0.06/kWh (natural gas, \$5/Mcf [gas 2015], 30% efficiency), while the maintenance cost is \$0.005/kWh. We vary the natural gas price in Figure 6(a) to study the impacts of natural gas prices.

**Electrical grid:** The base electricity price is fixed at \$0.056/kWh (Texas) [ele 2015]. There is no sell back price.

**Emissions:** The emissions rate from PV array, GE generation, and electricity grid are 0.034g/kWh, 0.443g/kWh, and 0.5g/kWh, respectively.

**Prediction errors.** Prediction errors are assumed to be mutually independent and follow the normal distribution with zero mean. To study the impacts of prediction errors, we vary the normalized RMSEs (root mean squared errors) in Figure 4.

### 4.2. Optimizing traditional data centers with renewable energy

In this subsection, we evaluate the joint framework on planning and operating a sustainable data center. We answer three questions: How does the optimization framework plan annually? How much benefits can the optimization framework achieve? How do prediction errors impact on the proposed framework?

**Annual capacity planning.** In practice, the electricity prices, gas prices, and workload demand tend to increase in the long-term. The average annual-increasing rates of electricity prices, gas prices, and workload demand are 1.05, 1.01, and 1.09 [EIA

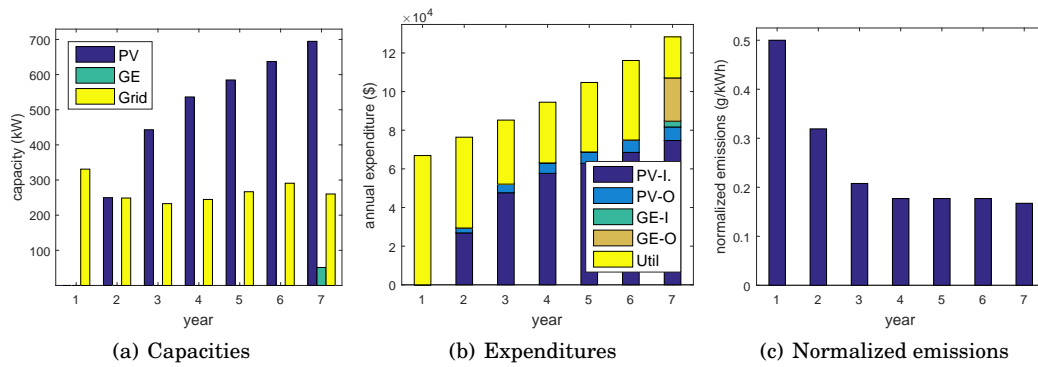


Fig. 1. Annual capacity planning. The data center is going to use more PV generation and GE generation while reducing the imported energy from the electricity grid. The major expenditure concentrates on the infrastructure of PV. The normalized emissions go down due to the high penetration of PV generation.

2016; gas 2015; Vaughan-Nichols 2014]. Meanwhile, the amortized cost of PV array decreases 12% annually [CleanTechnica 2015].

Figure 1(a) presents the capacities of power sources in 7 years. In general, the data center increases the capacity of PV annually but not the peak grid power consumption. In the first year, the data center prefers the electricity grid to other power sources because of the low electricity price. However, the data center significantly expands PV generation capacity from year 2 as the electricity prices increase and the infrastructure cost of PV decreases. In year 7, the data center provisions GE generation since the slow natural gas becomes relatively cheaper than the imported electricity. Although the data center prefers to use PV, the peak grid power consumption is still noticeably large. The intuition behind this is that PV generation is not available during night time which requires the data center to provision grid power.

Figure 1(b) shows the annual breakdown expenditures of the data center. In the first year, the utility bill (Util) of the imported energy from the electricity grid is dominant. Meanwhile, the cost of PV infrastructure (PV-I) quickly increases because PV amortized cost is added by installing more PV every year. The PV O&M (PV-O) expenditure linearly increases as the PV capacity goes up. In year 7, the GE O&M cost (GE-O) is 15% the total annual expenditure. GE-O mainly comes from the amount of natural gas supplied to the GE.

The annually normalized emissions of the data center, defined as the ratios of total emissions to the total power demand, are plotted in Figure 1(c). Since the normalized emissions of the electricity grid are high, the normalized emissions of the data center are highest in year 1. The increase of the PV generation can reduce normalized emissions for the data center. The normalized emissions sharply go down to 36% in year 4 as compared to the first year. However, there is little change from year 4 to year 7 as the ratio of PV capacity to other power sources is not considerably decreased.

**How much cost savings and emission reductions does the proposed framework achieve?** To highlight the benefits of the proposed framework, we compare the proposed framework (PROP) with three baseline methods, namely grid-only, supply-only, and demand-only.

- Grid-only (GRID): The grid-only method only uses the grid power from the electricity grid to provision the power demand. It does not use any power demand management techniques.

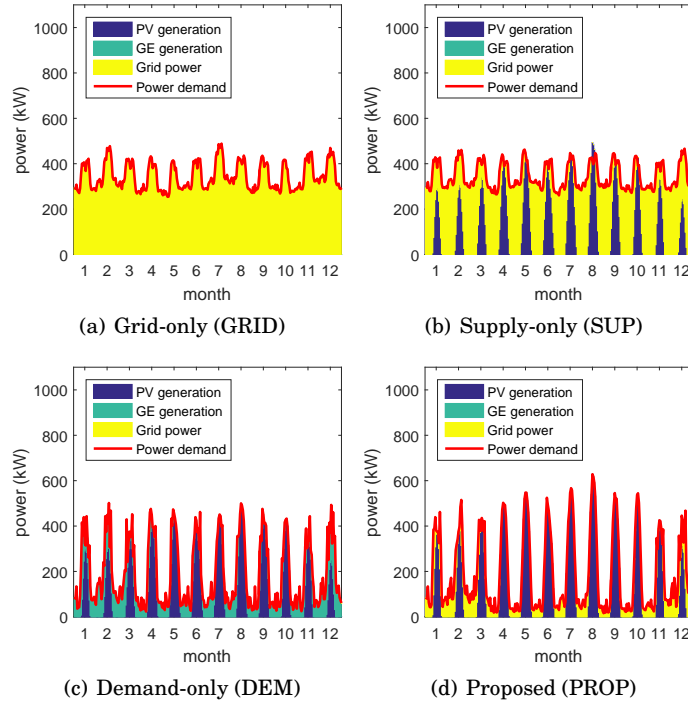


Fig. 2. The power profiles of the baseline methods and the proposed framework in year 5. GRID provisions the only grid power. SUP optimizes the power sources only at the supply side. DEM optimizes the power demand, i.e., scheduling the batch jobs.

- Supply-only (SUP): Given the power demand, the supply-only method optimizes capacity planning at the supply side. This method can optimize the use of energy sources among PV, GE, and the public electricity grid.
- Demand-only (DEM): At the supply side, the capacities of PV and GE generation are set at 50% and 70% of the power demand capacity, respectively. The demand-only method optimizes the power demand, i.e., scheduling the batch jobs, to reduce the operational cost.

In fact, PROP combines the SUP and DEM, and therefore can provide the best cost reductions.

The power profiles of these four methods in twelve typical days representing for the twelve months in year 5 are shown in Figure 2. GRID provisions power only from the electricity grid. SUP prefers the PV sources to the electricity grid and GE. Meanwhile, DEM utilizes installed GE generators because the electricity price is relatively more expensive than the O&M cost of GE sources. However, PROP uses only PV generation and grid power. In Figure 2(c) and 2(d), DEM and PROP shape the power demand to follow the PV generation while GRID and SUP are dependent on imported electricity.

We evaluate the four methods in terms of costs and emissions in Figure 3. It shows that PROP remarkably reduces the total expenditure by 50% while it achieves very close emissions to the lowest one, i.e., DEM. In Figure 3(a), SUP slightly reduces the total cost as it still depends much on the electricity grid. However, DEM shows that power management at demand side is very effective because it makes the power demand follow the PV generation.

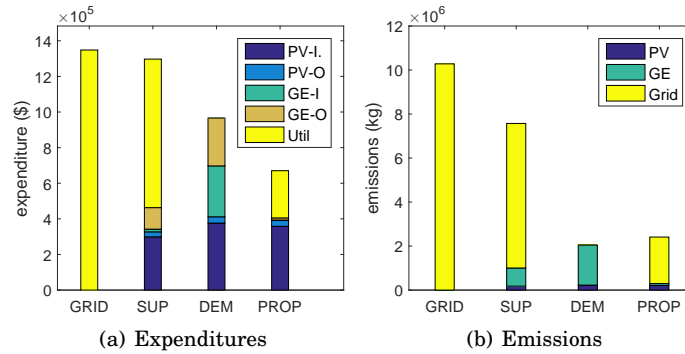


Fig. 3. Comparisons with baseline methods. The proposed framework reduces up to 50% of the total expenditures, and significantly cuts down 75% greenhouse gas emissions.

**Key insights:** (i) The proposed framework not only remarkably reduce the total cost, but also utilizes renewable energy very well. (ii) As the renewable installation becomes more cost-effective, the proposed framework prefers to use renewable energy and reduce the dependence of sustainable data centers on the electricity grid.

#### 4.3. Impacts of prediction errors

Prediction errors are generally negligible during operational management, which happens in real-time. Hence, GRID and DEM are not affected by the prediction errors because they do not need capacity planning. On the other hand, SUP and PROP suffer from prediction errors because they provide the capacity planning decisions, and then operate the data center in real-time based on the planned capacities. We normalize the root mean squared errors (RMSE) compared to the means of interactive workloads, batch jobs, capacity factors of PV, electricity prices, and gas prices, respectively. For instance, when normalized RMSEs are 0.2, the RMSEs of all the aforementioned predictions are 20% of their means.

Figure 4(a) shows the impacts of prediction errors on the total expenditures of the four methods. As the prediction errors become large (more than 10%), the total costs of SUP and PRO go up while the costs of GRID and DEM stay unchanged. Interestingly, total cost of the proposed framework is still the best and achieves the significant cost savings, i.e., 58% of GRID. The intuition behind this is that the operational management is cost-efficient in using the various power sources and scheduling the batch jobs to compensate for the prediction errors.

The impacts of prediction errors on emissions are presented in Figure 4(b). As the prediction errors increase, the emissions of PROP go up. However, PROP still reduces 49% of emissions reduction compared to GRID. Specially, the emissions stay unchanged when RMSE is greater than 0.2.

**Key insights:** Under large prediction errors, our proposed framework still achieves significant cost and emission reductions compared to the baseline methods.

#### 4.4. Sensitivity analysis

We carry out the experiments based on a real data center, called Net-zero Energy Data Centers (NEDC) invented by HP [Arlitt et al. 2012]. In NEDC, the total local generation (i.e. PV and GE generations) is greater than the total power consumption.

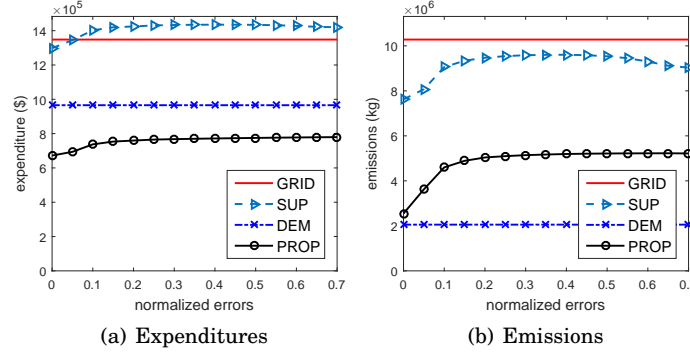


Fig. 4. Impacts of prediction errors. Under large prediction errors, the proposed framework still achieves the significant cost and emission reductions.

NEDC have an additional constraint

$$\sum_{t=1}^T \sum_{r=1}^R c_r(y, t) + \sum_{t=1}^T \sum_{s=1}^S c_s(y, t) \geq \sum_{t=1}^T P(y, t), \quad \forall y,$$

where the left hand side and the right hand side are the total power generation and the total power consumption in year  $y$ , respectively.

We focus on studying the impacts of supply and demand factors on the data centers during the first year. The supply factors include electricity price, and gas price. The demand factors include shape of interactive workload and ratio of flexible workload. Besides the capacities and expenditures of the data center, we study the payback period, which is the number of years for the data center to recoup the investment in the infrastructure costs of PV and GE instead of using only the electricity grid. The shorter payback period is, the more financial benefit the proposed framework provides.

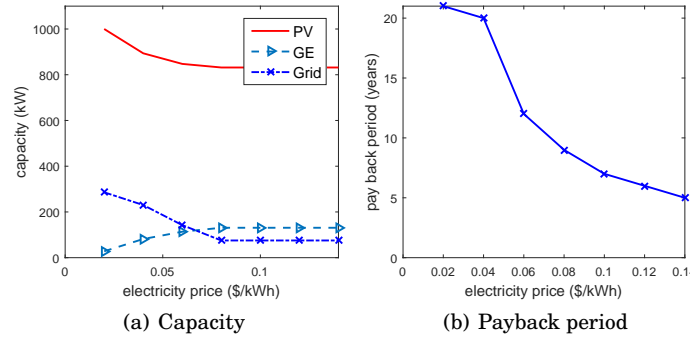


Fig. 5. Impacts of electricity prices. As the electricity price increases, the capacity of GE is increased to compensate for the PV generation during nighttime. Interestingly, it results in reducing the capacity of PV.

**4.4.1. Impacts of supply factors. Electricity price.** Figure 5 presents the impacts of electricity prices on the proposed framework. Figure 5(a) shows that the data center uses more GE generation and less grid power when the electricity price increases. However, the data center surprisingly keeps reducing the capacity of PV. It is because the data center starts to use more GE to provide power during nighttime and replace PV during daytime. It is because the costs of GE are relatively lower than the infrastructure of PV as PV is not fully utilized around its peak generation. In addition, when the grid

power becomes more expensive, the payback period sharply decreases as in Figure 5(b). Hence, the NEDC can significantly gain financial benefits when the electricity prices are high. *Gas price.* Figure 6 shows the impacts of gas prices on capacities and

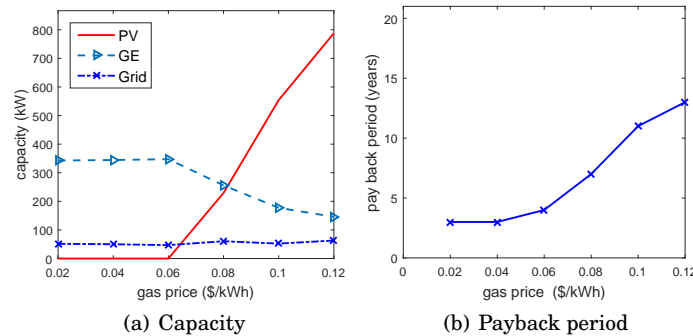


Fig. 6. Impacts of gas prices. As the gas prices increase, the capacity of PV increases quickly because the data center cannot import too much electricity grid.

payback periods. In Figure 6(a), the data center should use the GE generation at low gas prices, but it switches using PV when the gas price is more expensive. Especially, there is the sharp increase of PV capacity when the gas price is greater than 0.06. Due to the non-dispatchability of solar energy, the data center needs the large capacity of PV generation to compensate for the reduction of GE generation. As the gas price increases, the payback period goes up as in Figure 6(b) because the data center needs more PV.

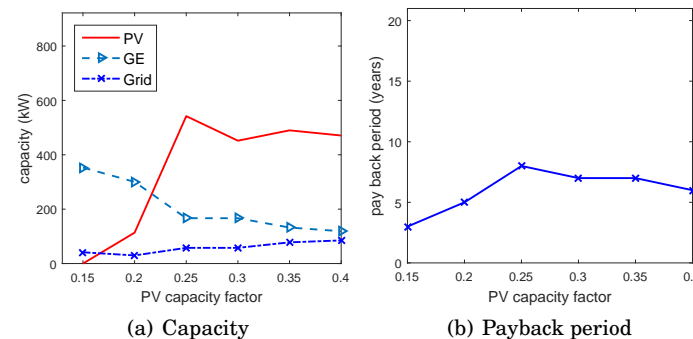


Fig. 7. Impact of different PV capacity factors. The curves of PV are very interesting which goes up and down because the high capacity factors have strong impact on the capital cost and operational cost of PV.

*PV capacity factor.* To understand how PV capacity factors affect on the proposed framework, we run the experiments with various capacity factors of PV. Recall the capacity factor is the average power generated, divided by the rated peak power. Figure 7 shows the dependence of the capacities and payback periods of the data center on PV capacity factors. The capacity of PV increases and then varies at high capacity factors. The capacity factor of PV array varies due to diverse reasons, such as geographical conditions. When PV has enough efficiency, i.e. capacity factor varies from 0.15 to 0.25, the data center starts to use PV. At a very high capacity factor (0.25-0.4), it is not necessary to increase the PV capacity because even the lower capacity with high capacity factors can still provide enough generation.

**Key insight:** The capacities of PV, GE, and peak power of grid power consumption adapt to the variety of supply factors accordingly. (i) The impacts of electricity prices and gas prices show the trade-offs between PV and GE. (ii) Under the increase of electricity price, the capacity of PV unexpectedly goes down together with the peak grid power. (iii) Under a certain gas price, the data center does not provision PV generation but only use GE and grid power instead.

4.4.2. *Impacts of demand factors.* Besides the supply factors, it must be interesting to study how the demand factors impact on capacity planning and costs of the data center.

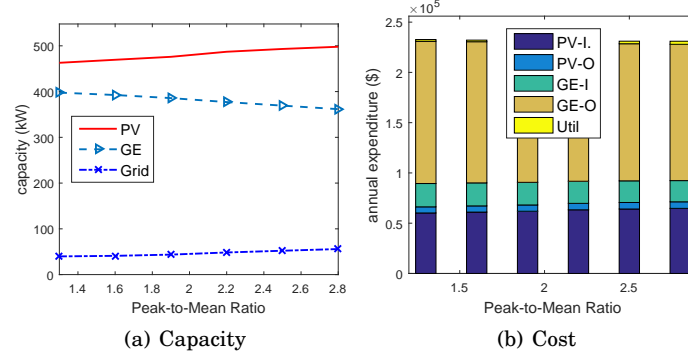


Fig. 8. Impacts of interactive workload shapes. The shapes of interactive workload have limited impacts on the capacities and expenditures of GE and PV.

*Shape of interactive workload:* Interactive workload is non-flexible and required to be processed with high responsive speed in data centers. We use the peak-to-mean ratio (PMR) to study the impacts of interactive workload shape. Figure 8 shows that shape of interactive workload has limited influence on both the capacities and expenditures of data centers. In particular, the capacity of PV and the peak grid power consumption slightly increase as the PMR increases. On the other hand, the data center slightly reduces GE when it has more PV power generation. As the capacities of PV and GE do not vary much, the breakdown expenditures are almost the same when varying the PMRs in Figure 8(b).

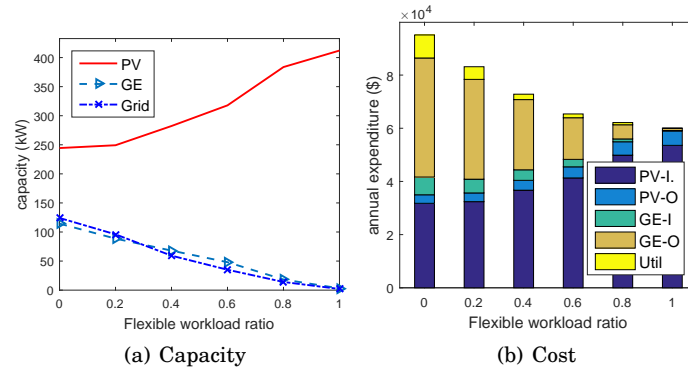


Fig. 9. Impacts of flexible workload ratios. The ratio can significantly reduces the total expenditure and change the capacity structure of the data center.

*Ratio of flexible workload* is the ratio of batch jobs to the IT workload demand. The higher ratio of flexible workload means the more flexibility data centers have in power demand management. In particular, Figure 9(a) highlight that the data center is more aggressive in using renewable energy as the ratio of flexible workload increases. It shows that the flexible workloads can promote the use of renewable sources. Especially, when workloads are totally flexible, there is no need to use GE and grid power as the workloads can be scheduled to totally follow PV generation. Figure 9(b) shows that the flexible workloads can significantly reduce the total expenditure, e.g., 28% reduction of the total cost.

**Key insights:** (i) The shape of interactive workload does not affect much on the capacity planning and operational management of data centers. (ii) The flexible workloads promote the use of renewable energy and significantly reduces the total cost.

## 5. DATA CENTER DEMAND RESPONSE

In this section, we extend the framework to study data center demand response (DCDR).

### 5.1. Modeling data center demand response

Table IV. Demand response rates.

Symbol	Value (\$/kWh)
$p_{tou}(t)$	0.05 (night), 0.219 (peak), 0.06 (off)
$p_{ibr}^l$	0.2 ( $l = 1$ : 50kW), 0.5 ( $l = 2$ : 100kW),
$p_{cpp}$	11.2
$p_{sr}$	0.02
$p_{ws}$	0.05

To model DCDR, the proposed framework is modified for including the costs of participating in DR programs into the objective function a part of *UtilBill*. In particular, we consider the following five DR programs:

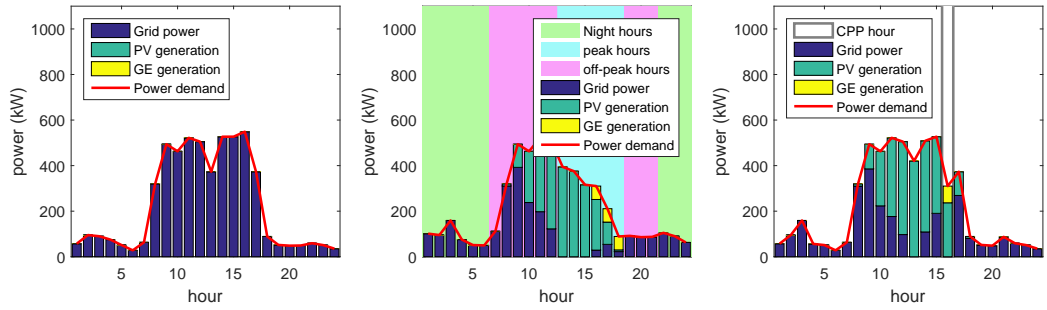
- Time-of-Use (ToU) rates,  $p_{tou}(t)$ , are defined based on the different times during a day, such as night time, peak time, and off-peak time [Smith 2011].
- Inclining block rates (IBR) encourage customers to consume electricity under some level,  $l$ , by charging a higher price, denoted by  $p_{ibr}^l$ , for the exceeding electricity usage.
- Under coincidence peak pricing (CPP) programs, industrial consumers like data centers are charged at a very high price,  $p_{cpp}$ , (e.g. than 200 times) for the usage during coincident peaks [Liu et al. 2013]. For example, the CPP time is an hour per month selected by the utility company.
- In spinning reserve (SR) service, electricity customers are rewarded based on predefined SR rates if they reduce their load after receiving an SR signal command.
- The rates  $p_{ws}$  in wholesale markets are typically cheaper than the regular electricity prices. The participations of data centers in wholesale markets allow electricity suppliers to efficiently plan their generation.

The DR rates are summarized in Table IV.

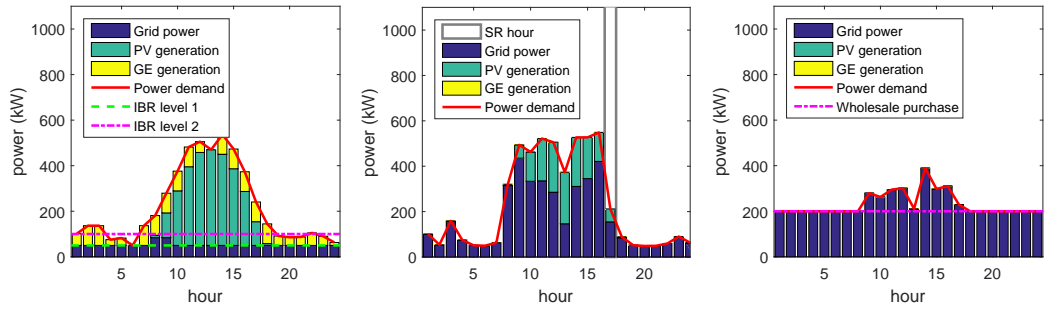
### 5.2. Numerical results

The numerical results are to study the impacts of DR programs on a data center using our proposed framework within a year. We answer the following two questions: How does the power profile of data center look like when participating in DR programs? How do DR programs impact on costs and emissions?





(a) Without DR programs. Data centers only use the grid power as provision 3 types of power sources GE generation and PV generation in CPP hours. The base prices of electricity are during peak-hours. cheap.



(d) IBR. The data center almost provisions its grid power under the center reduces the power demand during SR hour to earn the SR reward. (e) Spinning reserve (SR). The data center flattens the power demand to follow the pre-purchased power. (f) Wholesale (WS). The data center follows the pre-purchased power.

Fig. 10. The power profiles of the data centers participating in different DR programs.

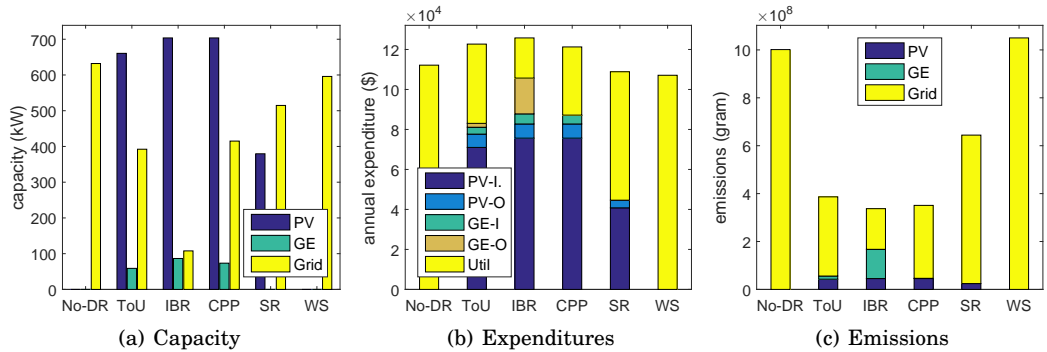


Fig. 11. Impacts of DR programs on capacities, expenditures, and emissions of the data center.

Figure 10 presents the typical daily power profiles of data centers participating in the DR programs. There are six cases: (1) Data centers do not participate in any DR programs, and data centers participate in (2) the ToU pricing program, (3) the CPP program, (4) the IBR program, (5) the SR ancillary service program, and (6) the wholesale market.

In the case of without any DR programs in Figure 10(a), the data center prefers to provision only the grid power because the electricity price is relatively cheaper than

using GE and PV. In addition, the proposed framework does not need to schedule the batch jobs as the electricity price is fixed.

Figure 10(b) illustrates the power profiles of the data center in the ToU program. We have 3 ToU price levels for off-peak hours  $\$0.06/kWh$ , night hours  $\$0.05/kWh$ , and peak hours  $\$0.219/kWh$  [Shahan 2011]. The rates of night hours are the cheapest while the electricity rates of peak hours are the most expensive. The data center provisions PV and GE generation to reduce the electricity bill during peak hours. As the peak of PV generation is in the peak hours, PV generation contributes the most power during peak hours. Compared to Figure 10(a), the power demand is shifted to match the PV generation and avoid the high peaks during the peak hours. When the PV generation goes down within the ToU peak hours, the data center operates GE to serve the power demand. During off-peak hours, the data center prefers to use PV generation if available, then imports the energy from the electricity grid if necessary. Furthermore, the peak of power demand is not at the peak hours compared to the case of without any DR programs. However, why does the data center not shift the power demand from the off-peak hours (7 am - 12 pm) to the night hours (9 pm - 6 am)? Because the interactive workload and batch jobs are colocated in the same servers, scheduling the batch jobs may increase the idle power.

Figure 10(c) shows the case of the data center participating in the CPP program. During the CPP hour, the data center needs to avoid using the grid power because the CPP rate is too high. Therefore, the data center provisions GE generation and utilizes PV generation during the CPP hour. A larger amount of PV generation is used compared to the case without DR programs. The power demand management schedules the batch job workload to react to the shape of PV generation, which has the peak in the afternoon.

We study the power profile of the data center with respect to two IBR levels which are level 1 (50kW) and level 2 (100kW) as in Figure 10(d). The electricity prices of exceeding the level 2 grid power is  $\$0.5/kWh$ , more expensive than the level 1, i.e.,  $\$0.2/kWh$ . The idea of IBR is to regulate the power demand under the two load levels. As we expected, the data center adapts to the IBR program very well. In particular, the grid power is mostly under the level 1 and never exceeds the level 2. In order to provision power under the IBR levels, the data center actually requires a lot of PV generation and GE generation. The batch job workload is shifted to the high peak of PV generation during day time. On the other hand, GE is used in the whole day.

Figure 10(e) presents the operation of the data center in the SR program. In the SR program, the data center can earn financial benefits if it reduces the power demand compared to the baseline consumption. In this simulation, the baseline consumption is the power consumption of the data center without participating in any DR program. We run the optimization framework to minimize the total cost for the data center in the DR program. During the SR hour, the data center reduces 30% their grid power consumption as compared to the case without DR programs.

In the wholesale market, it is assumed that the data center provisions 200kWh every hour at a cheaper price ( $0.05 \$/kWh$ ) than the base price ( $0.056 \$/kWh$ ). The power profile of data center is in Figure 10(f). It is seen that the power demand is flattened to follow the pre-purchased electricity in the wholesale market. The data center significantly reduces its peak compared to the case without DR programs, which can be very beneficial to reduce the peak demand of the electricity grid.

The comparison of the six cases of data center demand response in terms of capacities, costs, and emissions are shown in Figures 11. In general, the proposed framework enables the data center to adapt to each DR program very well. The data center increases the capacities of GE and/or PV under the ToU, CPP, IBR, and SR programs as in Figure 11(a). Hence, the DR programs can indirectly change the capacity planning

of data centers. In Figure 11(b), the total costs of the data center increase in the ToU, IBR, and CPP programs but they decrease in the SR program and the wholesale market. The lowest expenditure is in the wholesale market because of the cheap electricity price but it releases the most emissions as the grid power is generated by mainly using the fossil fuel. Participating in ToU, IBR, and CPP programs causes the data center spends slightly more expenditures, but data centers can remarkably reduce emissions by using other environmentally friendly power sources rather than the grid power.

**Key insights:** The proposed framework enables data centers to adapt to DR programs very well. (i) The data center uses a lot of PV and/or GE generation under ToU, IBR, CPP, and SR. (ii) While the total cost of the data center slightly increases in the ToU, CPP, and IBR programs, it can be reduced in the SR program and wholesale market. (iii) Data centers participating in ToU, IBR, CPP, and SR programs can cut down their emissions by provisioning other power sources rather than grid power.

## 6. CONCLUSIONS AND FUTURE WORK

In this paper, we propose an optimization framework for joint capacity planning and operational management that not only plans the capacities for sustainable data centers but also takes the operational management into account. The proposed framework can actually cut down significant expenditures by integrating the optimizations on both supply and demand sides. Numerical evaluation based on real-world case studies highlights the benefits to data center operators by using the proposed framework. In particular, it can achieve up to 50% cost savings and 75% emission reductions.

There are a lot of interesting future directions. For instance, tackling the stochastic characteristics of workload and renewable energy in capacity planning and operational management is challenging and important. Another promising direction is to extend the framework from a single data center to the system of geographically distributed data centers, which has more flexibility on planning their IT capacities because the workload demand can be shifted among different data centers. These can result in further cost and emission reductions.

## Acknowledgments

Part of this work was done when Zhenhua Liu visited Hewlett Packard Labs. This work is partially supported by NSF through CNS-1464388 and the MSIP “ICT Consilience Creative Program” (IITP-2015-R0346-15-1007) of Korea.

## REFERENCES

- 2014. *IT efficiency races forward in eBay Inc.'s data centers with HP EcoPODs*. Technical Report. HP.
- 2015. Texas Electricity Rates & Consumption. <http://www.electricitylocal.com/states/texas/>. (2015). <http://www.electricitylocal.com/states/texas/>
- 2015. United States Natural Gas Industrial Price. <https://www.eia.gov/dnav/ng/hist/n3035us3m.htm>. (2015). <https://www.eia.gov/dnav/ng/hist/n3035us3m.htm>
- Zeineb Abdmouleh, Rashid AM Alammari, and Adel Gastli. 2015. Review of policies encouraging renewable energy integration & best practices. *Renewable and Sustainable Energy Reviews* 45 (2015), 249–262.
- Baris Aksanli, Jagannathan Venkatesh, Liuyi Zhang, and Tajana Rosing. 2012. Utilizing green energy prediction to schedule mixed batch and service jobs in data centers. *ACM SIGOPS Operating Systems Review* 45, 3 (2012), 53–57.
- Matthew Andrews, Antonio Fernández Anta, Lisa Zhang, and Wenbo Zhao. 2012. Routing for power minimization in the speed scaling model. *IEEE/ACM Transactions on Networking (TON)* 20, 1 (2012), 285–294.
- Martin Arlitt, Cullen Bash, Sergey Blagodurov, Yuan Chen, Tom Christian, Daniel Gmach, Chris Hyser, Niru Kumari, Zhenhua Liu, Manish Marwah, and others. 2012. Towards the design and operation of net-zero energy data centers. In *Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, 2012 13th IEEE Intersociety Conference on. IEEE, 552–561.

- Luiz André Barroso, Jimmy Clidaras, and Urs Hölzle. 2013. The datacenter as a computer: An introduction to the design of warehouse-scale machines. *Synthesis lectures on computer architecture* 8, 3 (2013), 1–154.
- CleanTechnica. 2015. Solar Costs Will Fall Another 40% In 2 Years. Here's Why. <http://cleantechnica.com/2015/01/29/solar-costs-will-fall-40-next-2-years-heres/>. (2015).
- EIA. 2016. Electricity. <https://www.eia.gov/forecasts/steo/report/electricity.cfm>. (2016). <https://www.eia.gov/forecasts/steo/report/electricity.cfm>
- EPA. 2005. Greenhouse gas emissions from a typical passenger vehicle. <https://www.epa.gov/greenvehicles/greenhouse-gas-emissions-typical-passenger-vehicle-0>. (2005).
- Saurabh Kumar Garg, Srinivasa K Gopalaiyengar, and Rajkumar Buyya. 2011. SLA-based resource provisioning for heterogeneous workloads in a virtualized cloud datacenter. In *Algorithms and Architectures for Parallel Processing*. Springer, 371–384.
- Michael Grant, Stephen Boyd, and Yinyu Ye. 2008. CVX: Matlab software for disciplined convex programming. (2008).
- Jonathan G Koomey. 2008. Worldwide electricity used in data centers. *Environmental Research Letters* 3, 3 (2008), 034008.
- Minghong Lin, Adam Wierman, Lachlan LH Andrew, and Eno Thereska. 2011. Online dynamic capacity provisioning in data centers. In *Communication, Control, and Computing (Allerton), 2011 49th Annual Allerton Conference on*. IEEE, 1159–1163.
- Minghong Lin, Adam Wierman, Lachlan LH Andrew, and Eno Thereska. 2013. Dynamic right-sizing for power-proportional data centers. *IEEE/ACM Transactions on Networking (TON)* 21, 5 (2013), 1378–1391.
- Zhenhua Liu, Yuan Chen, Cullen Bash, Adam Wierman, Daniel Gmach, Zhikui Wang, Manish Marwah, and Chris Hysler. 2012. Renewable and cooling aware workload management for sustainable data centers. In *ACM SIGMETRICS Performance Evaluation Review*, Vol. 40. ACM, 175–186.
- Zhenhua Liu, Minghong Lin, Adam Wierman, Steven H Low, and Lachlan LH Andrew. 2011. Greening geographical load balancing. In *Proceedings of the ACM SIGMETRICS joint international conference on Measurement and modeling of computer systems*. ACM, 233–244.
- Zhenhua Liu, Iris Liu, Steven Low, and Adam Wierman. 2014. Pricing data center demand response. In *The 2014 ACM international conference on Measurement and modeling of computer systems*. ACM, 111–123.
- Zhenhua Liu, Adam Wierman, Yuan Chen, Benjamin Razon, and Niangjun Chen. 2013. Data center demand response: Avoiding the coincident peak via workload shifting and local generation. *Performance Evaluation* 70, 10 (2013), 770–791.
- Matt. 2009. Carbon Footprints of Servers Can Vary By 10X. <http://www.vertatique.com/carbon-footprints-servers-can-vary-10x>. (2009). <http://www.vertatique.com/carbon-footprints-servers-can-vary-10x>
- Tridib Mukherjee, Ayan Banerjee, Georgios Varsamopoulos, Sandeep KS Gupta, and Sanjay Rungta. 2009. Spatio-temporal thermal-aware job scheduling to minimize energy consumption in virtualized heterogeneous data centers. *Computer Networks* 53, 17 (2009), 2888–2904.
- Ehsan Pakbaznia and Massoud Pedram. 2009. Minimizing data center cooling and server power costs. In *Proceedings of the 2009 ACM/IEEE international symposium on Low power electronics and design*. ACM, 145–150.
- A. Qureshi, R. Weber, H. Balakrishnan, J. Gutttag, and B. Maggs. 2009. Cutting the electric bill for internet-scale systems. In *Proc. of ACM Sigcomm*.
- Vidhan Rana. 2009. The World of Data Centers. <http://whittakerassociates.com/the-world-of-data-centers/>. (2009). <http://whittakerassociates.com/the-world-of-data-centers/>
- C. Ren, D. Wang, B. Urgaonkar, and A. Sivasubramaniam. 2012. Carbon-Aware Energy Capacity Planning for Datacenters. In *MASCOTS*. IEEE, 391–400.
- Eleanor Revelle. 2009. CAP-AND-TRADE VERSUS CARBON TAX TWO APPROACHES TO CURBING GREENHOUSE GAS EMISSIONS. *League of women voters of the United States: Climate change task force* (2009).
- Zachary Shahan. 2011. Time of Day Pricing in Texas. <http://cleantechnica.com/2011/12/27/time-of-day-pricing-in-texas/>. (2011).
- Abhigyan Sharma, Xiaozheng Tie, Hardeep Uppal, Arun Venkataramani, David Westbrook, Aditya Yadav, Antonio AA Rocha, Ramesh Sitaraman, Jim Kurose, Dipankar Raychaudhuri, and others. 2015. *Shrink: A Cluster Manager for Greening Content Datacenters*. Technical Report.

- Dan Smith. 2011. TXU Energy Offers Deep Nighttime Discounts for Electricity. <http://www.businesswire.com/news/home/20111117005294/en/TXU-Energy-Offers-Deep-Nighttime-Discounts-Electricity>. (2011).
- Robert N Stavins. 2003. Experience with market-based environmental policy instruments. *Handbook of environmental economics* 1 (2003), 355–435.
- Yevgeniy Sverdlik. 2015. The Billions in Data Center Spending behind Cloud Revenue Growth. <http://www.datacenterknowledge.com/>. (2015). <http://www.datacenterknowledge.com/>
- B. Urgaonkar, G. Pacifici, P. Shenoy, M. Spreitzer, and A. Tantawi. 2005. An analytical model for multi-tier internet services and its applications. In *Proc. of ACM Sigmetrics*.
- R. Urgaonkar, B. Urgaonkar, M.J. Neely, and A. Sivasubramaniam. 2011. Optimal power cost management using stored energy in data centers. In *Proc. of the ACM Sigmetrics*.
- Steven J. Vaughan-Nichols. 2014. Cisco projects data center-cloud traffic to triple by 2017. (2014).
- Bill Weihl, Erik Teetzel, Jimmy Clidaras, Chris Malone, Joe Kava, and Michael Ryan. 2011. Sustainable data centers. *XRDS: Crossroads, The ACM Magazine for Students* 17, 4 (2011), 8–12.
- Adam Wierman, Zhenhua Liu, Iris Liu, and Hamed Mohsenian-Rad. 2014. Opportunities and challenges for data center demand response. In *Green Computing Conference (IGCC), 2014 International*. IEEE, 1–10.
- Mingui Zhang, Cheng Yi, Bin Liu, and Beichuan Zhang. 2010. GreenTE: Power-aware traffic engineering. In *Network Protocols (ICNP), 2010 18th IEEE International Conference on*. IEEE, 21–30.
- Q. Zhang, M.F. Zhani, Q. Zhu, S. Zhang, R. Boutaba, and J. Hellerstein. 2012. Dynamic Energy-Aware Capacity Provisioning for Cloud Computing Environments. In *ICAC*.

Received February 2007; revised March 2009; accepted June 2009