



A distributed approach to emergency demand response in geo-distributed mixed-use buildings

Chuan Pham^{a,*}, Nguyen H. Tran^{b,c}, Shaolei Ren^d, Choong Seon Hong^c, Kim Khoa Nguyen^a, Mohamed Cheriet^a

^a Synchromedia – École de technologie supérieure, Université du Québec, Canada

^b School of Information Technologies, The University of Sydney, NSW 2006, Australia

^c Department of Computer Science and Engineering, Kyung Hee University, Republic of Korea

^d Department of Electrical and Computer Engineering, University of California at Riverside, USA

ARTICLE INFO

Keywords:

Emergency demand response

Mixed-use building

Geographically distributed datacenters

ABSTRACT

Emergency Demand Response (EDR) has attracted research attention in recent years with its critical role in smart grids. Even though there are numerous potential participants for EDR, we especially focus on EDR, especially within datacenters and buildings, due to their huge power consumption yet flexible control knobs for power shedding. To reduce the deployment cost, many edge datacenters now are co-located inside buildings, which are responsible for power and IT infrastructure (called mixed-use buildings). In this paper, we consider a scenario that has not been addressed in the literature, in which multiple loads in geographically Distributed Mixed-use Buildings (geo-MUBs) can team up to participate EDR. We then design a mechanism that can coordinate tenants and geo-distributed buildings to minimize the system cost for EDR based on a robustly distributed framework, Alternating Direction Method of Multipliers (ADMM). In this mechanism, we also design a privacy-preserving scheme to conceal all tenants' transactions by using a lightweight algorithm. Simulation results show that our proposed method can reduce the total cost by 48.8% compared to existing approaches while satisfying all tenants constraints.

1. Introduction

Recent years, the electric power industry considers the Demand Response (DR) program as an increasingly valuable resource option, which has high potential and capacity to expand via grid modernization efforts [1]. The most important benefit of DR is to improve resource-efficiency of electricity production through various pricing schemes, such as Real Time Pricing (RTP), Time of Use (TOU) and critical peak pricing (CPP) [2]. According to a report by the Federal Energy Regulatory Commission (FERC) [3], the DR resource contribution from all U.S. DR programs is estimated to be about 5.8% of the peak demand in summer 2008 and will increase up to 14% of the peak demand by 2019. Among three types of DR, in this work, we consider the Emergency Demand Response (EDR), which has been widely adopted throughout the world and can even be executed multiple times each day in some developing countries where power infrastructure is increasingly fragile. In some cases of emergency situations, such as extreme weather conditions, EDR is performed by a coordination of many large energy consumers for power demand reductions in order to defend for power

grids without cascading blackouts [4]. Furthermore, as more renewables are incorporated into the grid and result in a higher volatility in power supply, we anticipate that EDR will be playing an even more crucial role [5].

As a promising candidate in EDR, datacenters receive much attention in both business and literature. In particular, [6] showed that datacenters have been identified by U.S.EPA as valuable assets for EDR. The amount of electricity usage of datacenters is large and increases rapidly; for example, in 2011, datacenters consumed 1.5% of all electricity worldwide [7]. Moreover, datacenters are now equipped with highly automated and monitoring functions, including power loads and IT equipment statuses (e.g., servers, storages) [8], and cooling facility. Such equipment can perform load curtailments to reduce the power consumption [9]. Also, there exist many well-known scheduling techniques used to control the workload demands [10–12] in datacenters. These techniques can help datacenters perform EDR flexibly, such workload load balancing, workload consolidation, shifting workloads to low-cost periods.

Another important candidate for EDR is the building. As reported in

* Corresponding author.

E-mail address: chuan.pham.1@ens.etsmtl.ca (C. Pham).

[7], buildings use about 40% of the global energy. Nonetheless, these EDR resources (datacenters and buildings) are not independent or isolated from each other. According to a report by Green Grid [13], the majority of datacenters are located within mixed-use facilities or mixed-use buildings (MUB). A further study [14] showed that owner-operated datacenters only account for 4% of the total datacenter energy consumption, whereas the remaining 96% is used by other types of datacenters (e.g., scientific computing clusters, colocation datacenters, and server rooms) that are mostly located in MUBS. Therefore, a *coordination between datacenters and non-datacenters colocated in MUB* is an urgent need for efficient EDR but under-explored in the literature. This issue becomes more significant since colocations of MUBs often reside in dense metropolitan areas, where DR programs are critical for cutting load at peak demands [9].

Coordinated EDR for a single MUB. Even though there are some studies focusing on the coordinated workload scheduling [15–17], most of them take into account the colocation datacenters or geo-distributed datacenters in terms of neglecting non-datacenter loads. These existing studies are used price-based or incentive-based mechanisms, in which the network operator needs to follow reactively the electricity price and load reduction signal in DR. The new trend of solutions in this area is shown in [18,19], where these MUB tenants (datacenters and offices) can be coordinated and aligned to participate in the EDR. With those approaches, they demonstrate the lower total incurred cost during compared to the individual control.

Although [18,19] can reduce significantly the energy consumption compared to other un-coordinated approaches, they leave two un-addressed aspects of MUBS that we consider in this work. The first issue is *the workload distribution between MUBs*. Recently, there are many buildings currently located in different locations, which are managed by the same entity to participate in EDR. Loads in those buildings, including offices and datacenters, share infrastructure (e.g., electric lines, Heating, Ventilation, and Air Conditioning (HVAC) systems), called geo-distributed MUBs. MUB energy management needs efficient mechanisms that can use migration techniques to distribute workloads among distributed datacenters for reducing energy consumption. Second, none of research in EDR for MUB considers *the privacy-preserving communication* to secure the sensitive information of MUB tenants in the network even in iterative distributed methods that require thousands of message passing in execution. Such a secure mechanism in implementation to ensure a fairly coordinated manner, in which all loads are fairly controlled to distribute precisely an amount of energy as they need without any cheating as well as exposing sensitive data in networks.

EDR of geo-distributed MUBs.

“As aforementioned discussion, MUB energy management needs dynamic manners through coordinated management across geo-distributed buildings in the smart grid. A key idea is to leverage the energy demand in multiple buildings flexibly based on the demand workload migration and temperature awareness of each building. However, existing methods based on those aspects [20–22], mainly focus on individual control systems (non-datacenters and dedicated datacenters), which cannot be applied directly to MUB, especially in terms of geo-distributed MUBs. Furthermore, when participating in EDR individually, such tenants may not satisfy the amount of energy reduction as in the EDR contract agreement due to high delay-sensitive workloads [18,19,23] (e.g., datacenters cannot turn off more servers due to ensuring the quality of services) or the limitation of discomfort setting from users during EDR time slot [24]. A major additional reason to investigate a coordinated manner in geo-MUBs is that in recent years, many cloud providers have built large datacenters in geographically distributed locations to achieve reliability while minimizing the operational cost [25]. In that geo-distributed scenario, these datacenter tenants should coordinately control the demand workloads and efficiently manage the energy usage across MUBS to exploit the potential of DR programs. For example, when receiving high workloads at the peak

price of electricity, an MUB datacenter cannot handle load curtailment strategies to reduce the energy consumption due to violate the Quality of Service (QoS) constraints. That datacenter can migrate dynamically workloads to other tenants in the coordinated system.

In this paper, we study a coordination in EDR of office tenants and edge datacenters in geo-distributed MUBs, which lease spaces from building managers. Even though these tenants are located in geo-distributed buildings, they are managed by the same entity to team up for an EDR. In particular, we consider one datacenter operator that has multiple geo-distributed datacenters. Each datacenter collocates with offices in a building to share building infrastructures. Such participants (e.g., datacenter operator, offices) have agreed to limit their total energy usage below the maximum energy consumption by signing a shaving load contract via the building operator (as a broker). An example of this model can be referred to Equinix. In the case of a multi-tenant MUB, each load/tenant may be small to participate directly in EDR and hence, they can sign an EDR contract via the building operator. During EDR period, the MUB energy manager (the main controller), which manages the whole geo-distributed MUB system, observes the electricity price, temperature conditions and energy load of each MUB to determine the amount of workloads distributed to MUB's datacenters in terms of minimizing the system cost (defined later) incurred by reducing energy. Within each MUB, there is a local operator that can handle the thermal load of offices as well as the backup generator system to satisfy the energy usage limitation according to the EDR contracts.

Privacy-preserving energy management for EDR. Last but not least, we study the privacy-preserving issue to secure the control messages in the coordinated manner of geo-distributed MUBs, where MUB tenants need to transmit their information to the network controller. A need of a privacy-preserving mechanism in a distributed manner is to guarantee a fair control between tenants, where all control messages are encrypted to ensure an uncheatable communication. We ingeniously integrate a privacy-preserving mechanism into the iteratively distributed algorithm to encrypt all control messages in a safety way. Since control messages contain tenant private and sensitive information, such as the discomfort cost [24], the discomfort temperatures, and/or the service latency, a practical solution for a coordinated manner should be carefully considered to secure such information.

The main contribution in this paper is the development of a *coordinated and privacy-preserving mechanism for all tenants in geo-distributed MUBs to participate an EDR program*, in which key contributions are summarized as follows:

- We study the well-known problem of edge computing providers, who house edge datacenters located in geo-distribution locations (including buildings). All the tenants in buildings are managed by the same entity to participate an EDR program. We investigate an optimization problem to formulate for multiple MUB tenants (including non-datacenter tenants and datacenter tenants) under their specific constraints (called P_{GMUB}). This optimization problem is transformed into a solvable Alternating Direction Method of Multipliers (ADMM) form. Using the ADMM framework, we decompose P_{GMUB} into sub-problems corresponding to each MUB, where workloads are managed under an energy aware mechanism.
- We also solve the subproblem at each MUB by applying ADMM, where all tenants and the backup generator system within the MUB are controlled to satisfying the delay-sensitive workloads distributed from the main controller as well as the limitation settings of offices and the energy usage contract. Furthermore, we ingeniously *integrate a privacy-preserving mechanism into ADMM mechanism to secure the control messages* carrying sensitive information of MUB tenants.
- Finally, using real trace workloads and dynamic electricity prices, we *performed many case studies to validate our proposed method*. The results show that our algorithm outperforms the existing non-coordinated approaches in terms of reducing the total cost.

The rest of the paper is organized as follows. Section 2 discusses the current literature, related to our proposed method. Section 3 presents the system model and problem statement. To solve the problem, we discuss a solution applying the ADMM method and privacy-preserving techniques in Section 4 and 5. We then simulate and evaluate our work in Section 6. Finally, we conclude in Section 7.

2. Related work

EDR of MUB tenants. At present, there exist several EDR programs that focused on ideal participants, such as buildings and datacenters. With buildings (office tenants), these tenants can reduce the power consumption through lighting power reduction, global thermostat set-point setback control, supply air temperature adjustment, pre-cooling, and use of a discharging energy storage device (e.g., battery) [26–28]. For MUB datacenter tenants, many resource management approaches are proposed and implemented in real-systems by consolidation [29], scaling down CPU frequencies [30], dynamically turning-on/off servers in owner-operated datacenters [8,31] and performing load balancing of the workloads [32,33]. Generally, unless considering non-datacenter loads, our proposed geo-distributed MUBs model can be simplified to the existing green load balancing model in geo-distributed datacenters [20–22], where workloads can be driven or migrated between datacenters to participate EDR efficiently. Having said that, existing mechanisms cannot be applied in geo-MUBs due to following reasons. An individual small datacenter in MUBs cannot shed enough amount of energy to participate DR programs. Similarly, by controlling separately non-datacenter loads, the controller can meet the limitation of the discomfort settings. To shed enough reducing energy following the signed DR contract, the energy manager needs to handle the backup generator system, which leads to increase carbon footprint [34].

To mitigate aforementioned disadvantages, an alternative approach for cost-effective EDR by MUBs is proposed in [18] via the coordinated energy management. These works illustrated the effectiveness of coordinated energy management compared to individual control mechanisms. They continued and developed their mechanisms on this approach in [35] with an incentive mechanism that can distribute the shedding energy reduction to each MUB tenant in order to minimize the total incurred cost. Despite achieving the minimum total cost in their mechanisms compared to the state-of-the-art methods, these works did not focus on joint load balancing workloads and energy management in geo-distributed MUBs, where the workload distribution in datacenters needs to be energy-aware not only by datacenters but also offices.

Last but not least, the current distributed mechanisms used in MUB have another drawback in that it discloses tenant information when communicating with the MUB operator in the network. In this work, we discuss some privacy-preserving methods that can be used in MUB.

Privacy-preserving communication mechanisms. There is a number of papers on private computations [36–39] that focus on encrypted data in the network. Among these privacy-preserving approaches, not all of them can be applied to encrypt control messages in an iterative distributed algorithm. Due to thousands of iterations, a light-weight mechanism with a less computation is suitable to be applied. In this work, we consider the method proposed in [36] that is applied successfully in many areas, such as data mining, cloud services and medical data. In this work, our designed method is designed to embed an encrypted protocol into an iteratively distributed algorithm. Here, we advocate an encrypted homomorphic model that is applicable for large scale system, such as geo-distributed MUBs. We utilize the simple secure summation protocol in [36] with low computational cost to protect private information.

3. System modeling

We consider a typical geo-distributed MUBs model in a particular region: a cloud provider is running cloud services on a set $\mathcal{S} = \{1, \dots, I\}$

Table 1
Notation.

\mathcal{S}	Geo-distributed MUB datacenters
S_i	A set of MUBS, indexed by $i = 1, \dots, I$.
s_i^{dl}	The number of servers at MUB datacenter i .
s_i^b	The number of active servers for serving delay-sensitive workloads at datacenter i .
s_i^b	The number of active servers for serving delay-tolerant workloads at datacenter i .
\bar{s}_i^b	The threshold for serving delay-tolerant workloads at MUB i .
Λ	The total incoming workloads.
λ_i	The incoming workloads at MUB datacenter i .
e_i^b	The amount of energy consumption for serving delay-tolerant workloads at MUB i .
e_i^{dl}	The amount of energy consumption for serving delay-sensitive workloads at MUB i .
\underline{e}_i^{dl} and \bar{e}_i^{dl}	The upper and lower bound settings of energy consumption for delay-tolerant workloads at MUB i , respectively.
s_i	The amount of servers turned-off at MUB datacenter i .
ω^{sla} , ω^{wat} and ω^b	The weight representing the unit cost of SLA, wear-and-tear and delay-tolerant workloads, respectively.
$d_i(\cdot)$	The delay function of delay-sensitive workloads at MUB datacenter i .
D_i	The SLA threshold of MUB datacenter i .
μ_{u_i}	The server's service rate delay at MUB datacenter i .
$p_{i,a}$ and $p_{i,s}$	The static and active powers of each server at MUB i , respectively.
PUE_i	The power usage effectiveness of MUB datacenter i .
\mathcal{N}_i	MUB office tenants
T_{ij}^c	A set of MUB office tenants at MUB i , indexed by $j = 1, \dots, N_i$.
R_{max} and R_{min}	The cooling TCL temperature of MUB office tenant j at MUB i .
$\zeta_i(t)$	The on/off state of the cooling system TCL.
Δt_{hvac}	The discretization time step.
C_a	The thermal capacitance.
R_e	The thermal resistance.
C_p	The Coefficient Of Performance (COP).
P_n	The related power.
T_a	The ambient temperature.
$w(t)$	The process noise.
e_{ij}^{cl}	Energy consumption of MUB office j at MUB i .
e_{ij}^{cf}	Comfort energy setting of MUB office j at MUB i .
E_{ij}^m	The active power of the cooling system of office j at MUB i .
e_i^{bg}	The backup generator
	Energy supplied by the backup generator at MUB i .
β^{cl} , β^{dc} and β^{bg}	The optimization problem.
C^{cl} , C^{dc} and C^{bg}	The weighted factors.
Q_i	The cost of MUB office tenants, MUB datacenter tenants and the backup generator, respectively.
	Energy consumption signed contracts with EDR providers at MUB i .

of MUBS. We take into account the locations of MUBS within the same region, such that they receive the EDR signals at the same time (e.g., the power system in California has some energy shortfalls, and hence, requires all the cities within California to cut loads. The emerging edge computing providers that house their servers in edge datacenters located in geo-distributed locations can be seen as examples for this model). At MUB i , datacenter i is colocated with a set \mathcal{N}_i of offices and the backup generator. Similar to [40], we consider one time-slot Δt demand response (e.g., 15 min or 1 h) that matches an interval for which the operator's decisions can be updated. We summarize all notations of the system model in Table 1.

3.1. Geo-distributed MUBs

We consider that each MUB datacenter i consists of S_i servers, which

are assumed to be homogeneous in this work. A datacenter with heterogeneous servers can be viewed as multiple virtual datacenters, each having homogeneous servers. In general, there are two types of workloads in MUB datacenters: *delay-sensitive workloads*, such as Internet services, and *delay-tolerant workloads* such as indexing data, batch jobs, etc. For delay-sensitive workloads, a response time is strictly imposed (usually in milliseconds), while delay-tolerant workloads can be flexibly scheduled to run any time as long as they can be completed within the deadlines. In this paper, we define the variables s_i^{dl} and s_i^b as the number of servers needed to serve the delay-sensitive and tolerant workloads at MUB datacenter i , respectively. Further, we denote $\bar{s}_i^b \leq S_i$ as the maximum number of servers that need to be active for executing the delay-tolerant workloads. Thus, we have the constraint $0 \leq s_i^b \leq \bar{s}_i^b, \forall i \in \mathcal{J}$.

In terms of delay-sensitive workloads, the total amount of incoming workloads Λ is driven to the main controller in a specific scheduling period Δt of EDR. The front-end controller will distribute the delay-sensitive workloads to appropriate MUB datacenters with the following objectives: *balancing the workloads and minimizing the total incurred cost*. We define the parameter λ_i as an amount of interactive workloads distributed to datacenter tenant i , satisfying the following load balance constraint:

$$\sum_{i \in \mathcal{J}} \lambda_i = \Lambda. \quad (1)$$

Energy consumption of MUB datacenters. There are various control knobs for reducing the power consumption in MUB datacenters as mentioned in Section 2. One widely used approach is turning-off idle servers, which has attracted much attention in recent years since switching servers can be easily handled without manual effort [8,17,41].

Inspired by the calculation in [17,42], the power consumption of a server comprises the static/idle and dynamic power. For the static/idle power, this term is measured based on the physical capacity of servers. For the dynamic power, which depends on server utilization, we assume that all delay-sensitive workloads driven to datacenter i will be evenly distributed to all active servers. Thus, the energy needed for MUB datacenter i to serve that workloads is as follows: $s_i^{dl} \left(p_{i,s} + p_{i,a} \frac{\lambda_i}{s_i^{dl} \mu_i} \right) PUE_i$, where $p_{i,s}$ and $p_{i,a}$ are the static and active powers of each server at MUB datacenter i , respectively; μ_i is a server's service rate measured in terms of the amount of workloads processed per unit time at datacenter i ; $\frac{\lambda_i}{s_i^{dl} \mu_i}$ is the server utilization with s_i^{dl} active servers, and PUE_i is the Power Usage Effectiveness of datacenter tenant i , which is measured by IT plus non-IT power consumption divided by the IT power consumption.

During EDR, the total energy consumption for serving the delay-sensitive workloads is calculated as follows:

$$e_i^{dl}(s_i^{dl}) = s_i^{dl} p_{i,s} PUE_i \Delta t + p_{i,a} \frac{\lambda_i}{\mu_i} PUE_i \Delta t, \forall i \in \mathcal{J}. \quad (2)$$

Based on (2), the maximum power consumption of MUB datacenter i corresponds to all servers S_i used to serve the delay-sensitive workload. Hence, we define $\bar{e}_i^{dl} := e_i^{dl}(S_i)$ as the threshold of consumed energy for serving λ_i . Given the amount of energy consumption of an MUB datacenter i needed to serve given workloads λ_i , we can derive the number of active servers s_i^{dl} .

Similar to the delay-sensitive workloads, the energy usage for executing the delay-tolerant workloads is as follows:

$$e_i^b(s_i^b) = p_{i,a} s_i^b PUE_i \Delta t, \forall i \in \mathcal{J}. \quad (3)$$

For ease of notations, we define the threshold of energy consumption for serving delay-tolerant workloads as $\bar{e}_i^b := e_i^b(\bar{s}_i^b)$.

Based on the EDR contract, we assume that MUB datacenter i manages its energy consumption by turning-off s_i servers in the EDR time slot. The amount of servers that need to be turned-off can be derived as follows: $s_i = S_i - s_i^{dl} - s_i^b, \forall i \in \mathcal{J}$.

In contrast to the energy reduction, frequently turning-off servers can negatively affect tenant performance, resulting in increased incurred costs in datacenters, such as the wear-and-tear cost and Service Level Agreement (SLA) cost that are discussed in the next part.

MUB datacenters cost. We consider three types of cost induced by MUB energy reduction. The first one is the *SLA cost*, which is especially important in geographically distributed datacenters since it directly affects the QoS. Using the M/M/1 queue model, the constraint of datacenter i 's workloads delay is as follows:

$$d_i(s_i^{dl}) := \frac{1}{\mu_i - \frac{\lambda_i}{s_i^{dl}}} \leq D_i, \forall i \in \mathcal{J}, \quad (4)$$

where D_i is the SLA threshold at MUB datacenter i . This constraint corresponds to the maximum tolerant workloads delay at MUB datacenter i , which is widely used as a measurement of QoS [43,44]. Specifically, turning-off servers at MUB datacenter i increases the delay d_i . Therefore, the controller has to guarantee the SLA constrain (4) during participating EDR. From this constraint, we can derive a threshold for the number of servers that can be turned-off:

$$s_i^{dl} \geq \frac{\lambda_i}{\mu_i - 1/D_i}, \forall i \in \mathcal{J}. \quad (5)$$

From (2) and (5), we derive a constraint for the energy consumption that guarantees the mean service delay of the delay-sensitive workloads as follows:

$$e_i^{dl} := e_i^{dl} \left(\frac{\lambda_i}{\mu_i - 1/D_i} \right) \leq e_i^{dl} \leq \bar{e}_i^{dl}, \forall i \in \mathcal{J}. \quad (6)$$

We assume that s_i is a continuous variable with thousands of servers in the datacenter tenants following [45]. According to (4), when turning-off servers, the workloads delay increases, which can be modeled as the SLA cost in datacenters [45] as follows: $\omega^{sla} \lambda_i d_i(s_i^{dl}), \forall i \in \mathcal{J}$, where ω^{sla} is the weight representing the unit cost of SLA.

The second cost is the *wear-and-tear cost*. Even though turning-off servers can reduce the energy consumption, the controller should not frequently turn-on/off arbitrary servers due to inducing the wear-and-tear cost, as shown in [8,46], which is detrimental to the lifetime of physical servers. In details, the wear-and-tear cost can be calculated as a linear function depending on the number of servers that are turned-off, as follows: $\omega^{wat} s_i, \forall i \in \mathcal{J}$, where ω^{wat} is the weight representing the wear-and-tear unit cost.

Finally, to characterize the delay performance of executing the delay-tolerant workload, we model the cost as a *penalty* function, which is defined based on the number of allocated servers as follows: $\omega^b (1 - s_i^b / \bar{s}_i^b), \forall i \in \mathcal{J}$, where ω^b is the weight representing the unit cost of the delay-tolerant workload.

Therefore, during the EDR, the total cost of MUB datacenter i is as follows:

$$C_i^{dc} = \omega^{wat} s_i + \omega^{sla} \lambda_i d_i(s_i^{dl}) + \omega^b \left(1 - \frac{s_i^b}{\bar{s}_i^b} \right), \forall i \in \mathcal{J}. \quad (7)$$

As the aforementioned discussion, given an amount of workloads λ_i and energy for serving delay-sensitive workloads and for executing delay-tolerant workload, we can derive the variables s_i^{dl}, s_i^b and s_i . Hence, for ease formulation, we define $C_i^{dc}(e_i^{dl}, e_i^b)$ as the cost function of the MUB's datacenter i , based on the energy consumption variables $\{e_i^{dl}, e_i^b\}$.

3.2. HVAC energy management for MUB office tenants

We consider a set \mathcal{N}_i of offices at MUB $i \in \mathcal{J}$, where each office is equipped with a smart thermostat that can measure the room temperature and communicate with the coordinator. Thermostatically Controlled Loads (TCLs) are good candidates for EDR that significantly

affect the performance of the control strategies.

Following the TCL model in [47], we denote the cooling TCL temperature at time t by $T^c(t)$ and the on/off state at time step t by $\zeta_i(t) \in \{0, 1\}$. For both “on” and “off” states, the thermal dynamics of a TCL system can be typically modeled as a linear system [47] as follows:

$$T^c(t+1) = aT^c(t) + (1-a)(T_\alpha - \zeta_i(t)R_e P_n) + w(t), \quad (8)$$

$$\zeta_i(t+1) = \begin{cases} 1 & \text{if } T^c(t+1) \geq T_{\max}, \\ 0 & \text{if } T^c(t+1) \leq T_{\min}, \\ \zeta_i(t) & \text{otherwise,} \end{cases} \quad (9)$$

where $a = e^{-\Delta t_{\text{hvac}}/(C_a R_e)}$, Δt_{hvac} is the discretization time step, C_a is the thermal capacitance, R_e is the thermal resistance, C_p is the Coefficient of Performance (COP), P_n is the rated power, T_α is the ambient temperature, and $w(t)$ is the process noise. Further, T_{\max} and T_{\min} are the upper and lower temperature dead-band limits, respectively.

In this paper, we assume each office has its own temperature setting that is separately controlled, rather than being centrally handled. Given a setting temperature, the energy consumption of each load during a time period can be derived by calculating the portion of time that the cooling system is on [48].

Energy consumption of MUB offices during EDR.

Adopted from the energy-temperature correlation model in [48,49], the HVAC power of an office tenant $j \in \mathcal{N}_i$ depends on the TCL temperature T_{ij}^c , corresponding to the setpoint $[T_{\min}, T_{\max}]$ and the state $\zeta_{i,j}(t)$ at office j . During Δt , the energy consumption of each office (the cooling system) can be calculated using the portion of time that the system is on. Therefore, the energy consumption of each office can be represented as

$$e_{ij}^{cl}(T_{ij}^c) = E_{ij}^m \int_0^{\Delta t} \zeta_{i,j}(t) dt, \quad \forall i \in \mathcal{I}, \forall j \in \mathcal{N}_i, \quad (10)$$

where E_{ij}^m is the active power of the cooling system of office j at MUB i Fig. 1.

Inspired by the Pacific Northwest GridWise National Laboratory demonstration project [50], each office tenant can easily control the temperature in the room and automatically adjust the setpoint and active load corresponding to the amount of load power allocation by equipping smart thermostats. Given an amount of energy e^{cl} , the office tenant adjusts the setting temperatures T_{\max} and T_{\min} as well as the active duration Δt_{hvac} for curtailing the energy expense. For example, in the summer, the cooling system is turned on when T^c is greater than 25°C and turned off when T^c is less than 20°C (these values are named as the comfort thresholds). With this setting, the mean temperature T^c in the office is kept around 23°C (called the comfort temperature T^{cf}) and the total energy use is 100 KWh (named the comfort energy e^{cf}). To participate in EDR, the controller aims to reduce power usage to 50 KWh by controlling the dynamic load (e.g., increasing the setpoint temperature as depicted in Fig. 2).

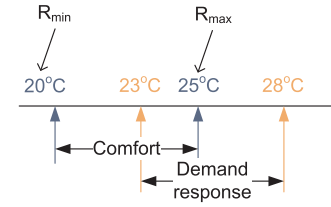


Fig. 2. The setting temperatures in the EDR program.

Having said that, changing a thermostat for shedding reduction energy always results in uncomfortable situations in the occupants. To improve the comfort of the occupant in offices, smart thermostats or programmable thermostats can be used in buildings [51]. They are integrated “anticipator” or hysteresis sensors to prevent excessively rapid cycling of devices when the temperature reaches close to the setpoints, and also reduce the magnitude of temperature variations caused by a switching status.

HVAC cost model. Although it is difficult to exactly measure the user comfort cost; fortunately, this is an abstract concept that has recently attracted many researchers [48,52,19]. Based on the dynamic load model above, we consider that e^{cf} is the most comfortable energy consumption for which the HVAC can satisfy the user needs based on the thermostat settings. To participate in the EDR program, each office tenant attempts to reduce its power consumption to satisfy the energy limitation. Therefore, we define the total user comfort cost of all office tenants at MUB i as follows:

$$C_i^{cl}(e_i^{cl}) = \sum_{j \in \mathcal{N}_i} w_i^{cl}(e_{ij}^{cf} - e_{ij}^{cl})^2, \quad \forall i \in \mathcal{I}, \quad (11)$$

where w_i^{cl} is a monetary weight (i.e., \$/ kWh) at MUB i . This cost model reflects that the user discomfort increases quadratically with respect to the deviation of the controlled energy e^{cl} from the comfort e^{cf} ($e^{cf} \geq e^{cl}$).

We choose a quadratic function for user discomfort because i) there exists many fields (e.g., control, signal processing, network communication) that widely use quadratic functions to model cost functions, and ii) it also provides analysis tractability. Furthermore, our study focuses on the *coordination* of multiple tenants in EDR; hence, a simple and intuitive user discomfort is chosen.

3.3. MUB backup generator cost model

As the aforementioned discussion, HVAC tenants and datacenter tenants together may not sufficiently satisfy the amount of required energy reduction according to the EDR requirement. In that case, the MUB operator has to turn on the backup generator system (as a traditional method [35]). Even though this option is costly, it is necessary to keep the system without any interruption.

Let q_i^{bg} and η_i denote the power capacity and efficiency of the generator at MUB i , respectively; then its energy production during time

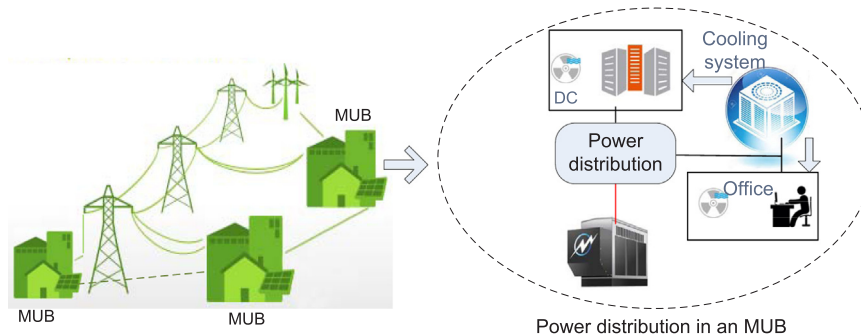


Fig. 1. The system model of geo-distributed MUBs.

slot Δt is

$$e_i^{bg} = \eta_i q_i^{bg}, \forall i \in \mathcal{I}, \quad (12)$$

where $\eta_i := \eta'_i \Delta t$. Since a generator can control its power capacity to produce a certain amount of energy during a time slot, following [53], the backup generator cost is expressed as follows:

$$C_i^{bg}(e_i^{bg}) = w_{i,0}^{bg} + w_{i,1}^{bg} e_i^{bg} + w_{i,2}^{bg} (e_i^{bg})^2, \quad \forall i \in \mathcal{I}, \quad (13)$$

where $w_{i,0}^{bg}$, $w_{i,1}^{bg}$, $w_{i,2}^{bg}$ are the fuel cost coefficient at site i

3.4. Problem formulation and challenges

It is critical for the geo-distributed MUBs to satisfy the EDR signals without causing many negative impacts on the SLA performance of datacenter jobs or the user comfort in buildings. Consequently, we consider the geo-distributed MUBs' cost minimization problem for EDR, which uses the weighted factor $\beta = \{\beta^{cl}, \beta^{dc}, \beta^{bg}\}$ to combine cost functions of MUB tenants as follows:

\mathbf{P}_{GMUB} :

$$\min. \sum_{i \in \mathcal{I}} \beta^{cl} C_i^{cl}(e_i^{cl}) + \beta^{dc} C_i^{dc}(e_i^{dl}, e_i^b) + \beta^{bg} C_i^{bg}(e_i^{bg}) \quad (14)$$

$$\text{s. t. } \sum_{i \in \mathcal{I}} \lambda_i = \Lambda, \quad (15)$$

$$e_i^{dl} + e_i^b + \sum_{j \in \mathcal{N}_i} e_{ij}^{cl} \leq Q_i + e_i^{bg}, \forall i \in \mathcal{I}, \quad (16)$$

$$0 \leq e_{ij}^{cl} \leq e_{ij}^{cf}, \quad \forall i \in \mathcal{I}, \forall j \in \mathcal{N}_i, \quad (17)$$

$$\underline{e}_i^{dl} \leq e_i^{dl} \leq \bar{e}_i^{dl}, \quad \forall i \in \mathcal{I}, \quad (18)$$

$$0 \leq e_i^b \leq \bar{e}_i^b, \quad \forall i \in \mathcal{I}, \quad (19)$$

$$0 \leq e_i^{bg} \leq Q_i, \quad \forall i \in \mathcal{I}, \quad (20)$$

$$\text{var. } \{e^{cl}, e^{dl}, e^b, e^{bg}, \lambda\}. \quad (21)$$

In \mathbf{P}_{GMUB} , the objective is to minimize the MUB's total cost incurred in the system for satisfying the limitation of the energy consumption signed in EDR contracts as well as the delay and the comfort constraints for the MUB tenant services. The load balancing delay-sensitive workload constraint is given in (15). When participating in EDR, the MUB controller at MUB i has to limit the total power usage of tenants so that it is not greater than the threshold Q_i (which is calculated based on the amount of energy consumption from contracts signed with the EDR providers and the production capacity of generation), as reflected in (16).

Based on the aforementioned problem formulation, we briefly describe our model as follows. In the colocation system of MUBS, datacenter workloads are managed by an IT workload component. It is responsible for distributing workloads Λ among distributed datacenters and an amount of tolerant workloads. This component needs to handle a set of active servers (related to $\{s_i\}$) for serving workloads. For receiving and routing workloads, the workload management module can be implemented as a gateway that operates similarly to the green load balancing datacenters [54]. In order to control temperature and humidity in an MUB, the HVAC component is responsible for controlling the thermal building based on the temperature setting $\{T_i^c\}$. These mentioned components are connected to an energy component that is responsible for controlling energy demand in the building, including energy for datacenters, HVAC and the backup generator system $e_i = \{e_i^{cl}, e_i^{dl}, e_i^b, e_i^{bg}\}$. All components are connected to the main controller that operates the system by solving \mathbf{P}_{GMUB} .

To make the solution for \mathbf{P}_{GMUB} feasible in practice, the challenge is to find a distributed manner for this model since MUBS are located in different locations and physically self-managed at each MUB. All

tenants teaming up in \mathbf{P}_{GMUB} have their own controllers that are independent of each other MUB. Thus, in order to utilize these independent controllers, a distributed implementation is essential. One approach is using the standard dual decomposition with (sub) gradient methods [55] for distributed optimization. However, since discomfort cost and backup generator cost are linear functions, the dual decomposition approach, which requires the cost function to be strictly convex, is not applicable [56].

Last but not least, a “privacy-preserving” communication mechanism for EDR in geo-distributed MUBs should be considered to secure sensitive information of MUB tenants via control message passing. A distributed and privacy-preserving mechanism in a coordinated manner of geo-distributed MUBs will make it more significant and visible in practice.

4. Energy reduction mechanism via distributed computation

The alternating direction method of multipliers (ADMM) [56] is a robust framework that can solve convex optimization problems in a distributed manner. In this section, we describe the basics of the ADMM framework and design a distributed mechanism to solve \mathbf{P}_{GMUB} .

4.1. Basic of the ADMM framework

We consider the following convex optimization problem with N ($N \geq 2$) blocks of variables:

$$\min_{\mathbf{x}} \sum_{i=1}^N f_i(x_i) \quad (22)$$

$$\text{s. t. } \sum_{i=1}^N a_i x_i = c, \quad (23)$$

where $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$.

The augmented Lagrangian function for (22) is as follows:

$$L(\mathbf{x}, \lambda) = \sum_{i=1}^N f_i(x_i) - \lambda \left(\sum_{i=1}^N a_i x_i - c \right) + \left(\rho/2 \right) \left\| \sum_{i=1}^N a_i x_i - c \right\|_2^2, \quad (24)$$

where λ is the Lagrangian multiplier and $\rho > 0$ is a penalty parameter. The iterative scheme of ADMM is outlined below:

$$\begin{cases} x_1^{k+1} := \arg\min_{x_1} L_p(x_1, x_2^k, \dots, x_N^k, \lambda^k), \\ x_2^{k+1} := \arg\min_{x_2} L_p(x_1^{k+1}, x_2, \dots, x_N^k, \lambda^k), \\ \dots \\ x_N^{k+1} := \arg\min_{x_N} L_p(x_1^{k+1}, x_2^{k+1}, \dots, x_{N-1}^{k+1}, x_N, \lambda^k), \\ \lambda^{k+1} := \lambda^k - \rho \left(\sum_{i=1}^N a_i x_i^{k+1} - c \right). \end{cases} \quad (25)$$

Even though each variable x_i can be carried out in a distributed manner, they cannot be solved in parallel, as seen in (25). Moreover, this also leads to exhibit the message passing overhead during operation. To make the solution of \mathbf{P}_{GMUB} distributed and parallel, a mechanism should enable an MUB to control and execute its computation separately. Thus, we introduce a new variable \mathbf{z} to fully decouple variable \mathbf{x} [57], where $x_i = z_i$, $i = 1, \dots, N$. The problem now is equivalent to

$$\min_{\mathbf{x}} \sum_{i=1}^N f_i(x_i) \quad (26)$$

$$\text{s. t. } a_i x_i - z_i = c/N, i = 1, \dots, N, \quad (27)$$

$$\sum_{i=1}^N z_i = 0. \quad (28)$$

The partial augmented Lagrangian function now is given by

$$L(\mathbf{x}, \mathbf{z}, \boldsymbol{\lambda}) = \sum_{i=1}^N f_i(x_i) - \sum_{i=1}^N \lambda_i \left(a_i x_i - z_i - c/N \right) + \left(\rho/2 \right) \sum_{i=1}^N \|a_i x_i - z_i - c/N\|_2^2. \quad (29)$$

Since all x_i and z_i are now fully decoupled, x_i can be solved parallel while \mathbf{z} is updated later by the main controller.

1) **x_i -sub-problem.** The x_i -sub-problem is represented as follow [57]:

$$x_i^{k+1} = \underset{x_i}{\operatorname{argmin}} f_i(x_i) + \left(\rho/2 \right) \|a_i x_i - z_i - c/N - (\lambda_i^k/\rho)\|_2^2. \quad (30)$$

2) **\mathbf{z} -sub-problem.** Especially, the \mathbf{z} -sub-problem now is a quadratic problem that can be obtained a closed-form solution [57], as follows:

$$z^{k+1} = \underset{\sum_{i=1}^N z_i=0}{\operatorname{argmin}} \left(\rho/2 \right) \sum_{i=1}^N \|a_i x_i - z_i - c/N - (\lambda_i^k/\rho)\|_2^2. \quad (31)$$

3) **Dual-update.** Finally, the dual variable $\boldsymbol{\lambda}$ is updated as follows:

$$\lambda_i^{k+1} = \lambda_i^k + \rho(x_i^{k+1} - z_i^{k+1}), i = 1, \dots, N. \quad (32)$$

Note that in ADMM, \mathbf{x} , \mathbf{z} and $\boldsymbol{\lambda}$ are updated sequentially instead of jointly as in the dual ascent algorithm [58]. Further, this order can be reversed, leading to a variation on ADMM. The optimality and convergence of the ADMM algorithm are given and proven in [59].

4.2. A distributed mechanism for \mathbf{P}_{GMUB}

\mathbf{P}_{GMUB} can be seen as a multi-block optimization problem with a set \mathcal{J} of functions $f_i(\cdot)$ that corresponds to the cost functions $C_i(\cdot)$ of MUB i . Hence, based on the ADMM decoupling technique, we introduce a set of auxiliary variables $z_i = \lambda_i, \forall i \in \mathcal{J}$. The partial augmented Lagrangian function of \mathbf{P}_{GMUB} is presented as follows

$$L = \sum_{i \in \mathcal{J}} L_i \left(e_i^{dl}, e_i^b, e_i^{cl}, e_i^{bg}, \lambda_i, v_i^k, z_i^k \right) = \sum_{i \in \mathcal{J}} C_i \left(e_i^{dl}, e_i^b, e_i^{cl}, e_i^{bg}, \lambda_i \right) - v_i \left(\lambda_i - z_i - \Lambda/I \right) + (\rho/2) \|\lambda_i - z_i - \Lambda/I\|_2^2, \quad (33)$$

where $C_i(e_i^{dl}, e_i^b, e_i^{cl}, e_i^{bg}, \lambda_i)$ is the total cost at MUB i and v is the Lagrangian multiplier variable. The distributed mechanism for solving \mathbf{P}_{GMUB} is detailed as follows:

1) **MUB sub-problem.** Based on this augmented Lagrangian (33) and the ADMM framework, the main operations of \mathbf{P}_{GMUB} can be decomposed to calculate at each MUB as follows. At iteration $k+1$, the MUB controller at MUB i solves the following sub-problem:

$$\mathbf{P}_{\text{MUB}}: \min L_i(e_i^{dl}, e_i^b, e_i^{cl}, e_i^{bg}, \lambda_i, v_i^k, z_i^k) \quad (34)$$

s. t. constraints (16) – (20). (35)

2) **\mathbf{z} -update.** Consider \mathbf{z} as a control variable in the system that is used to adjust the value of $\boldsymbol{\lambda}$ approaching the optimal value. This procedure can be invoked at the MUB controller to obtain the optimal \mathbf{z} based on the following sub-problem:

$$z^{k+1} = \underset{\sum_{i=1}^N z_i=0}{\operatorname{argmin}} \left(\rho/2 \right) \sum_{i=1}^N \|\lambda_i^{k+1} - z_i - \Lambda/I - (\lambda_i^{k+1}/\rho)\|_2^2. \quad (36)$$

3) **Dual-update.** Finally, \mathbf{P}_{GMUB} updates the dual variable v as follows:

$$v_i^{k+1} = v_i^k + \rho(x_i^{k+1} - z_i^k), \forall i \in \mathcal{J}. \quad (37)$$

By using the ADMM method, we decompose \mathbf{P}_{GMUB} horizontally into I sub-problems corresponding to MUBS. We next discuss about the solution of the sub-problem \mathbf{P}_{MUB} at each MUB.

4.3. Distributed solution to \mathbf{P}_{MUB}

Within each MUB, office tenants, datacenter tenants and backup generator systems team up to find an optimal energy consumption of each tenant underlying a signed EDR contract. Each MUB receives amount of delay-sensitive workloads from the main MUB controller, then determines a power distribution scheme to satisfy the EDR contract by solving the sub-problem \mathbf{P}_{MUB} . As mentioned in Section 3.4, we relax the inequality constraint in (16) to the equality constraint as follows:

$$e_i^{dl} + e_i^b + \sum_{j \in \mathcal{N}_i} e_{ij}^{cl} - e_i^{bg} = Q_i, \forall i \in \mathcal{J}. \quad (38)$$

The sub-problem \mathbf{P}_{MUB} is now in the ADMM form mentioned in Section 4.1 under the coupling equality constraint (38). We introduce a new variable \mathbf{q} , where $e_j = q_j, j = 1, \dots, J, N_i = |\mathcal{N}_i|$ and $J = 2 + N_i + 1$ is the number of loads at MUB i that are participant in EDR. For convenience, we define again the index of office tenants and the cost function $f_j(e_j), j = 1, \dots, J$ of each load. \mathbf{P}_{GMUB} now can be rewritten as follows:

$$\underset{e_j \in [0, e_j^{\max}]}{\min} \sum_{j=1}^J f_j(e_j) \quad (39)$$

$$\text{s. t. } a_j e_j - q_j = Q_i/J, j = 1, \dots, J, \quad (40)$$

$$\sum_{j=1}^J q_j = 0, \quad (41)$$

where vector $\mathbf{a} = \{a_1, \dots, a_J\}$ includes that the first two elements represent for delay-sensitive workloads and delay-tolerant workloads, the next N_i elements represent for the set n_i of office tenants, and the final element represents for the backup generator at MUB i ($\mathbf{a} = [\underbrace{11}_{\text{DC } N_i \text{ offices}}, \underbrace{-1}_{\text{Generator}}]^T$), and $f_j(\cdot)$ is the cost function corresponding to each tenant with the energy variable, written as

$$f_j(e_j) = \begin{cases} C_i(e_i^{dl}) & \text{if } j = 1, \\ C_i(e_i^b) & \text{if } j = 2, \\ C_{ij}(e_{ij}^{cl}) & \text{if } j = 3, \dots, N_i + 2, \\ C_i(e_i^{bg}) & \text{if } j = J. \end{cases} \quad (42)$$

and

$$e_j^{\max} = \begin{cases} \bar{e}_i^{dl} & \text{if } j = 1, \\ \bar{e}_i^b & \text{if } j = 2, \\ e_{ij}^{cf} & \text{if } j = 3, \dots, N_i + 2, \\ Q_i & \text{if } j = J. \end{cases} \quad (43)$$

Following the update steps of the ADMM framework in Section 4.1, the algorithm to solve the sub-problem \mathbf{P}_{MUB} is represented as follows:

1) **e -update.** Given the amount of sensitive workloads received from the main MUB controller, tenant j at each MUB iteratively adjusts his energy e_j as follows:

$$e_j^{l+1} = \underset{e_j \in [0, e_j^{\max}]}{\operatorname{argmin}} f_j(e_j) + \varphi/2 \left\| a_j e_j - q_j^l - \frac{Q_i}{J} - \frac{\sigma_j^l}{\varphi} \right\|_2^2, \quad (44)$$

where $\sigma = \{\sigma_1, \dots, \sigma_J\}$ is the Lagrangian multiplier variable.

2) **q-update.** The MUB controller collects the energy information of all tenants and updates the communication variable $q_j, j = 1, \dots, J$, as follows:

$$q_j^{l+1} = \underset{\sum_{j=1}^J q_j = 0}{\operatorname{argmin}} \left(\varphi/2 \sum_{j=1}^J \left\| a_j e_j^{l+1} - q_j - \frac{Q_i}{J} - \frac{\sigma_j^l}{\varphi} \right\|_2^2 \right). \quad (45)$$

3) **Dual-update.** Finally, each MUB updates the dual variable σ as follows:

$$\sigma_j^{l+1} = \sigma_j^l + \varphi \left(e_j^{l+1} - q_j^{l+1} \right), j = 1, \dots, J. \quad (46)$$

Following the steps of Algorithm 2, at iteration l , each tenant has to broadcast its energy consumption e_j to the controller, as in (45). This information is used to execute the load curtailment strategies (e.g., controlling the numbers of active servers and temperature setpoints) to satisfy the EDR contracts. However, this sensitive information should not be disclosed to others. In order to secure the tenants' information, in the next section, we integrate a privacy-preserving technique into Algorithm 2, where the control messages can be exchanged between tenants and the MUB controller without disclosing real information in the network.

Algorithm 1. Distributed algorithm for \mathbf{P}_{GMUB}

- 1: **Initialization:** $k = 0$, set a random value λ_i to satisfy (15);
- 2: **repeat**
- 3: $k \rightarrow k + 1$;
- 4: Adjust the energy consumption of MUB tenants by solving \mathbf{P}_{MUB} (Algorithm 2);
- 5: Update λ_i by solving \mathbf{P}_{GMUB} ;
- 6: Operator updates z_i as in (36) and v_j as in (37).
- 7: **until** $\|v^{(k)} - v^{(k-1)}\|_2 < \epsilon$.

Algorithm 2. Adjusting energy consumption for MUB tenants.

- 1: **Initialization:** $l = 0$, set a random value e_j to satisfy (38) and (43);
- 2: **repeat**
- 3: $l \rightarrow l + 1$;
- 4: Updates energy e_j of tenant j by solving (44);
- 5: Operator updates q_j as in (45) and σ_j as in (46);
- 6: **until** $\|q_j^{(l)} - q_j^{(l-1)}\|_2 < \epsilon$.

5. Privacy-preserving Communication of MUB Tenants for Participating EDR

In this section, we propose a *privacy-preserving* mechanism that avoids side information leakage associated with reporting the power usage information while participating in EDR. This mechanism encrypts sensitive information of the tenants that is not exposed in the network when executing Algorithm 2.

5.1. Observations and assumptions

Regarding the sensitive and private information of tenants, in Algorithm 2, the controller of each MUB has to gather the energy value each tenant e_j for the q -update step (Line 5). The closed-form of (45) can be derived as follows:

$$q_j^{l+1} = \left(a_j e_j^{l+1} - \frac{Q_i}{J} - \frac{\sigma_j^l}{\varphi} \right) - \frac{1}{J} \left(\sum_{j=1}^J a_j e_j^{l+1} - \frac{Q_i}{J} - \frac{\sigma_j^l}{\varphi} \right). \quad (47)$$

From (47), we observe that the controller only needs to know the average energy consumption $\frac{1}{J} \left(\sum_{j=1}^J a_j e_j^{l+1} \right)$. Therefore, we can use the privacy-preserving update, which gathers the summation of energy information of all tenants without knowing the individual values.

Without loss of generality, we assume that e_j^{l+1} can be represented as an integer based on [39]. Furthermore, we assume that every iteration of Algorithm 2 invokes the privacy-preserving update; hence, we omit the iterative superscript l .

5.2. A non-directional ring for the privacy-preserving update procedure

Recently, there are a few lightweight techniques used for securing summation protocols in data mining due to the development of geographically distributed datacenters. One of the popular approaches is the mechanism proposed in [36].

To apply the mechanism in [36], we assume that the tenants are arranged in a non-directional ring that includes a master (in our model, a master can be seen as the MUB controller), who creates the random key during the procedure. The tenants also do not have any collusion for cheating the data. Further, we assume that the total value $e_{\text{sum}} = \sum_{j=1}^J e_j$ belongs to $[0, \eta]$. η can be set by the limitation of energy Q_i of MUB i . This range can be derived from constraint (38). The following procedure is also depicted in an example in Fig. 3. First, the master selects a random generator number R and adds this value to its value (the value of master can be set 0). The next tenant i in the ring adds its value with the received value as follows:

$$e_j^{\text{send}} = e_j^{\text{rcv}} + (e_j \bmod \eta), \quad (48)$$

and sends this to the next tenant. The final tenant in the ring performs the same procedure and sends its value to the master. When receiving the return value, the master subtracts R to obtain the actual result. Since each tenant receives a value that is encrypted by the key R , tenant j learns nothing about the actual value of the others Fig. 4.

Note that the backup generator tenant in our system can be represented by a negative value, as shown in the value of tenant 3 in the example of Fig. 3. This lightweight method can be applied directly in Algorithm 2 with assumption of no collusion [36]. Furthermore, the algorithm \mathbf{P}_{MUB} divides the computation into tenants so that they can execute independently their procedures.

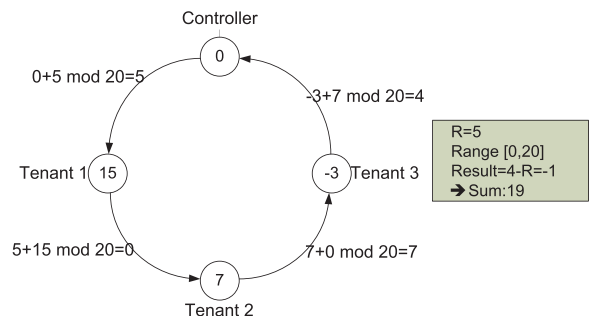


Fig. 3. Non-directional ring for the privacy-preserving update procedure.

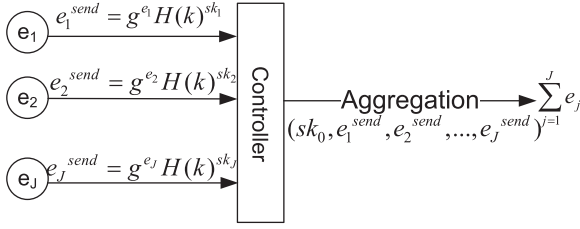


Fig. 4. A distributed aggregation for the privacy-preserving update procedure.

Table 2

Parameters for simulation of the MUB office tenants (cooling systems).

Parameters	Interpretation	Setting values
T_a	Ambient temperature	27–32 °C
T_{min}	Minimum temperature	20–22 °C
T_{max}	Maximum temperature	23–25 °C
C_a	Thermal capacitance	10 kWh/°C
R_e	Thermal resistance	2 °C/kW
P_n	Rated power	14 kW
C_p	Coefficient of Performance (COP)	2.5

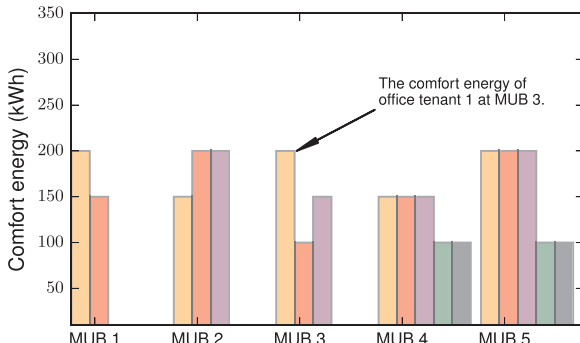
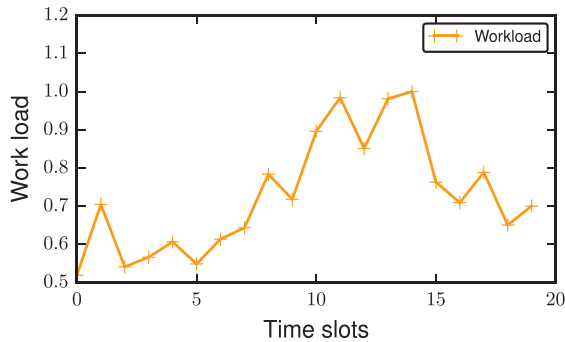


Fig. 5. The comfort energy setting for MUB office tenants.

5.3. A distributed aggregation for the privacy-preserving update procedure

To address the above-mentioned challenges of the non-directional ring method, we apply the proposed method in [60], which can be performed in a distributed aggregation process. The privacy-preserving procedure is represented as follows. Each tenant j has a secret key s_j^k and the controller has a secret key $s_0^k = -\sum_{j=1}^J s_j^k$. Using their keys, each tenant encrypts its data as $e_j^{send} = s_j^k + e_j$. The controller sums all the received values and decrypts the actual sum as follow: $e_{sum} = \sum_{j=1}^J e_j^{send} - s_0^k = \sum_{j=1}^J e_j$.



(a) The trace workloads.

Furthermore, to secure the sensitive information of the tenants, secret keys should be created for each iteration of Algorithm 2. We rely on a hash function H that maps an integer to a cyclic group of prime order p [60]. Let \mathcal{G} denote a cyclic group of prime order p for which the Decisional Diffie-Hellman problem is hard [60]. A hash function is denoted as follows: $H: \mathcal{Z} \rightarrow \mathcal{G}$, where \mathcal{Z} is the set of random secret keys $\mathcal{Z} = \{s_1^k, s_2^k, \dots, s_J^k\}$. The distributed aggregation for the privacy-preserving update procedure is represented as follows:

- Setup: The controller chooses a random generator $g \in \mathcal{G}$ and $J + 1$ random secret keys $s_0^k, s_1^k, \dots, s_J^k \in \mathcal{Z}_p$ such that $s_0^k = -\sum_{j=1}^J s_j^k$.
- Encrypting and sending data: At iteration l , each tenant j calculates their energy based on (44) and encrypts the data as follows:

$$e_j^{send} = g^{e_j} H(k)^{s_j^k}. \quad (49)$$

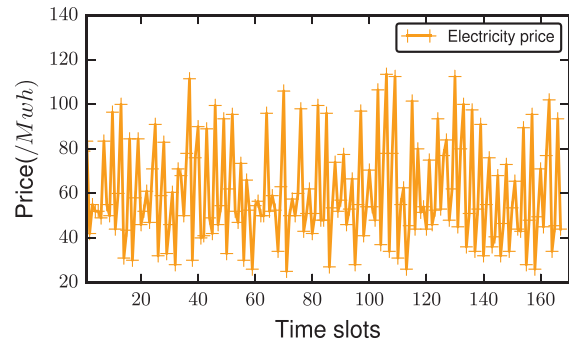
- Receiving and decrypting data: At the MUB controller, all the data are aggregated and calculated as follows:

$$\begin{aligned} e_{sum} &= \log_g \left(H(k)^{s_0^k} \prod_{j=1}^J e_j^{send} \right) \\ &= \log_g \left(H(k)^{s_0^k} \prod_{j=1}^J g^{e_j} H(k)^{s_j^k} \right) \\ &= \log_g \left(g^{\sum_{j=1}^J e_j} \right) \\ &= \sum_{j=1}^J e_j. \end{aligned} \quad (50)$$

Since $\sum_{j=0}^J s_j^k = 0$, we have $\prod_{j=1}^J H(k)^{s_j^k} = 1$, which we can leverage to construct a scheme for aggregating MUB tenant's information without communication between tenants in each iteration. Consequently, this privacy-preserving aggregation can be executed separately by each MUB tenant in Algorithm 2. The computational overhead of this approach, which includes hashing, multiplication in a Diffie-Hellman group, and modular exponentiation, is more complicated than the non-directional ring. The robustness for defending against collusion attacks and the efficiency of the computation over multiple iterations were discussed and proven in [60].

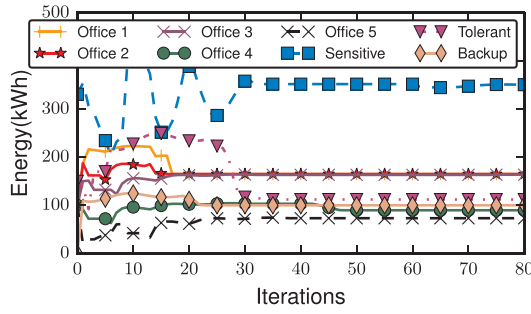
6. Numerical results

In this section, we first describe the simulation settings and then present the results of the proposed mechanism to validate their efficacy.

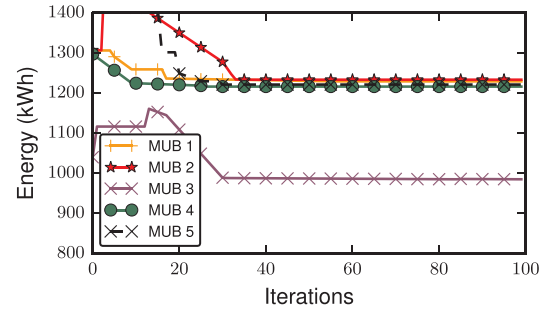


(b) The price of electricity.

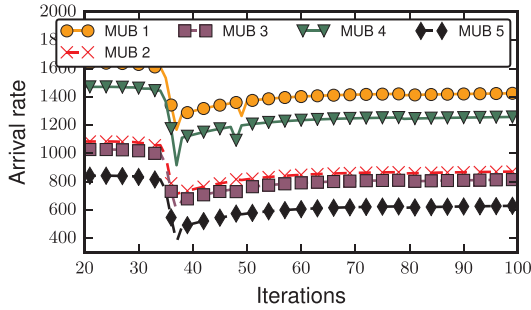
Fig. 6. The trace data for simulation.



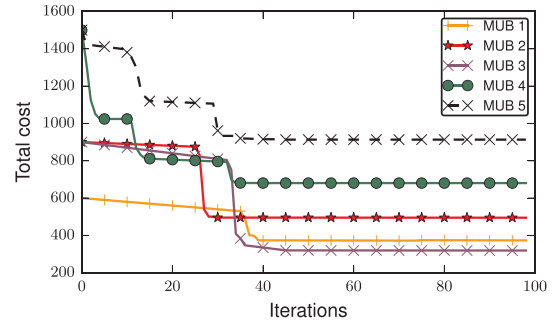
(a) Convergence of energy consumption of agents at MUB 5 in Alg. 2.



(b) Convergence of energy consumption in Alg. 1.



(c) Convergence of delay-sensitive workloads in Alg. 1.



(d) Convergence of total cost of geo-distributed MUBs in Alg. 1.

Fig. 7. Evaluation of convergence.

6.1. Settings

Geo-distributed MUBS setting. We consider a set \mathcal{J} of geo-distributed MUBs, which consists of five MUBs, $|\mathcal{J}| = 5$, and the limitation of energy consumption in the EDR contracts of geo-distributed MUBs is set between 900 and 1000 kWh. Each datacenter tenant has $S_i = 2000$ servers, and each server has an idle power $p_{i,s} = 200$ W and a peak power $p_{i,a} = 400$ W. This setting is referred to [61] with $\Delta t = 1$ hour. The average PUE in our system is set to 1.5, which is a typical setting for datacenters [22], i.e., whenever an MUB tenant consumes 1 kWh energy, the corresponding energy consumption at the geo-distributed MUBs level is 1.5 kWh. Using the assumption in Section 3, we set the service rate of servers μ_i to be increased from 2 to 4. The delay threshold D_i is set to 1 s, i.e., the maximum delay-sensitivity service is 1 s.

For MUB office tenants, we assume that in the summer, the comfort temperature is lower than the indoor temperature. Based on the TCL model, all of the specific settings of the office tenants are shown in Table 2. To maintain the comfort temperature of each tenant according to that setting, we set the energy comfort e^{cf} of office tenants from 100 kWh to 200 kWh. A snapshot of the comfort energy setting for all office tenants at a specific time slot is shown in Fig. 5, where each bar represents an amount of the comfort energy setting for each office tenant in MUB.

MUB's datacenter workloads. In this work, we collect the trace workloads from Facebook [62] during 20 time slots, as depicted in Fig. 6a. The data of workloads is normalized with respect to the datacenter tenants' capacities in our settings.

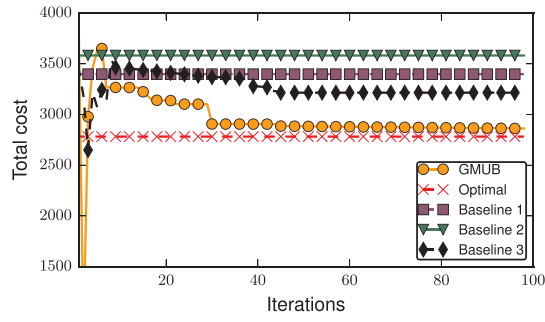
Electricity price. To reflect the impact of the electricity price on the energy management of the geo-distributed MUBs, we use various traces of electricity price in [3] for our simulation. The representative of the electricity price at all MUBs is depicted in Fig. 6b.

6.2. Simulation results

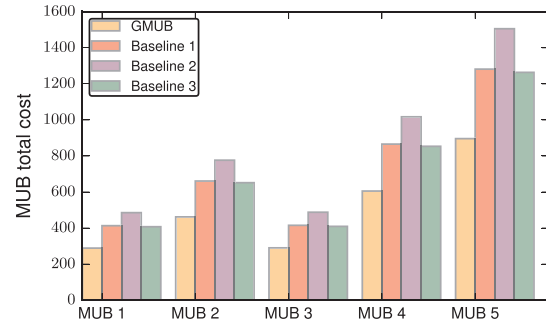
In order to evaluate the convergence, the optimality of our proposed method, we compare the result of our algorithm with three following baselines:

- **Baseline 1:** This baseline only controls the energy of each tenant without relying on the incurred cost at each MUB to adjust the delay-sensitive workload allocation. The demand workloads is driven equally to each MUB by the controller (i.e., $\lambda_i = \Lambda/I$). This baseline is used to evaluate the efficiency of workload management in terms of energy awareness.
- **Baseline 2:** Similar to the existing works [63,17], this baseline only relies on datacenters and the backup generator in each MUB without considering office tenants, i.e., $e_{ij}^{cl} = 0, \forall i \in \mathcal{J}, \forall j \in \mathcal{N}_i$. Thus, constraint (38) is reduced to $e_i^{dl} + e_i^b - e_i^{bg} = Q_i, \forall i \in \mathcal{J}$. For workloads, it is also equally distributed among all datacenter tenants (i.e., $\lambda_i = \Lambda/I$).
- **Baseline 3:** As aforementioned discussions, without considering office tenants, our model can be simplified into a green load balancing datacenters [12]. Hence, we concern this case as baseline to evaluate our performance. In particular, it only relies on datacenters and the backup generator for the EDR (i.e., $e_{ij}^{cl} = 0, \forall i \in \mathcal{J}, \forall j \in \mathcal{N}_i$).

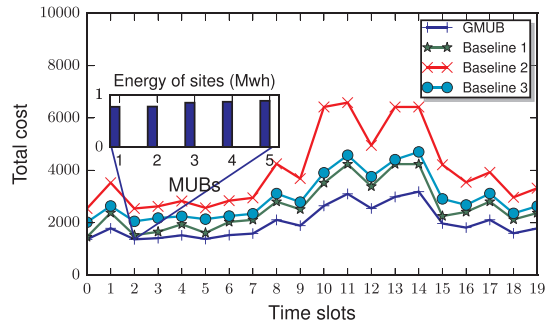
Algorithm convergence. Fig. 7 plots the convergence property of Algorithm 1 and 2. We first illustrate the outcome of energy consumption of an MUB in Algorithm 2 (for sub-problem \mathbf{P}_{MUB}). We select MUB 5 (the largest MUB in our simulation) to conduct the convergence in Algorithm 2. The result of energy consumption of all loads at MUB 5 is depicted in Fig. 7a, where the energy used for serving delay-sensitive workloads dominates others. We next evaluate the convergence of Algorithm 1 in Fig. 7b, Fig. 7c and Fig. 7d. The energy distribution in



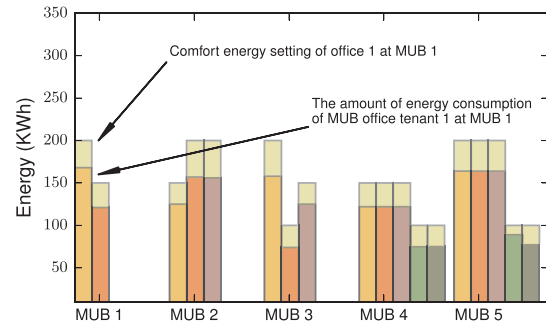
(a) Comparison of total cost with baselines and optimal solution.



(b) Evaluation of total cost incurred in the system within 20 time slots.

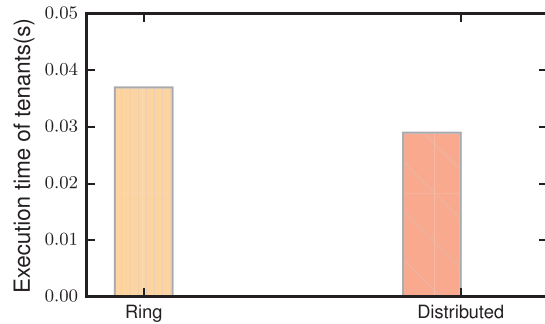


(c) Comparison of total cost with baselines within 20 time slots.

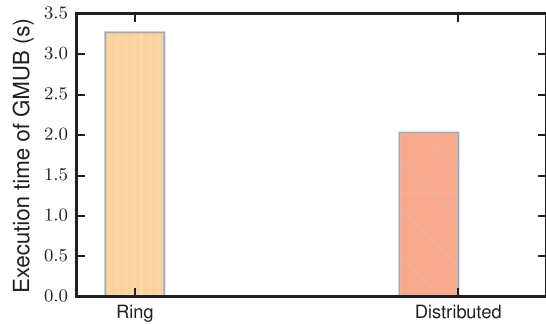


(d) Energy consumption of the MUBs' office tenants at time slot 10.

Fig. 8. Evaluation of the MUB total cost in short-term and over long-term.



(a) Evaluation of the average execution time.



(b) Evaluation of the average convergence time.

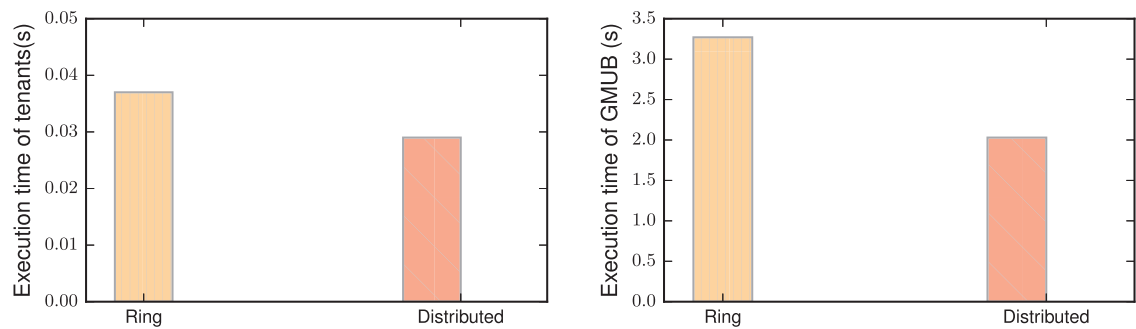
Fig. 9. Evaluation of the impact of reducing energy on the office temperatures during 20 time slots.

our system is plotted in 7b corresponding to the distribution of workload distribution in Fig. 7c and our settings. Due to lower values of service rates in settings, datacenters in MUB 1 and MUB 4 receive more workloads compared to others. Although the rate of workloads at MUB 5 is low, the number of office tenants in MUB 5 is high (mentioned in the prior settings). Such a setting results in the high energy consumption of MUB 5 as well as the highest cost incurred in MUB 5 in Fig. 7d. From all the convergence figures, we observe that Algorithm 1 converges quickly to the optimal solution, usually within 40 iterations. By using ADMM, Algorithm 1 can achieve not only the fast convergence but also reducing the calculation due to a complete distributed approach.

Evaluation of the MUB total cost. We evaluate the performance of our algorithm compared with the optimal solution (solved by the

centralized solver JuMP) and baselines in one-time slot (short-term). Fig. 8a shows the convergence of GMUB to the optimal solution with a small gap compared to Optimal. By distributing equally the workloads to all MUB datacenters, Baseline 1 and Baseline 2 induce the highest cost incurred in the system, as demonstrated in this figure. The incurred costs are not changed in these baselines because the arrival rates to all datacenter tenants are fixed for all iterations (i.e., $\lambda_i = \Lambda/I$). By teaming up three kinds of tenants, such as office tenants, datacenter tenants, and backup generator system, GMUB can reduce the total cost up to 11.2% compared to Baseline 3.

For the long-term consideration, we use electricity price settings for all MUBS to conduct the performance of our proposed mechanisms compared with baselines. Fig. 8b shows that GMUB outperforms others at all MUBS overall consideration time slots. During 20 time slots,



(a) Evaluation of the average execution time. (b) Evaluation of the average convergence time.

Fig. 10. Evaluation of the average execution time.

GMUB can reduce by up to 44.7%, 69.8% and 47.3% compared with Baseline 1, Baseline 2 and Baseline 3, respectively. In particular, without considering office tenants, Baseline 2 sheds the most reducing energy from the backup generator compared to others. Thus, it causes the highest cost in all MUBs. In order to show the impact of workload management during 20 time slots, we illustrate the variation of the total cost in Fig. 8c, compared to baselines. By teaming up with multiple tenants in EDR, GMUB induces significantly the total cost. Comparing to Baseline 2, GMUB can reduce the total cost averaged over all time slots by 47.3%, while comparing to Baseline 1 and Baseline 3, GMUB can reduce this by 29.1% and 35.2%, respectively. A snapshot of the energy consumption corresponding to a light workload distribution (at time slot 2) is also depicted in Fig. 8c, where all MUBs do not spend over the limitation of energy at that time. Furthermore, we also illustrate the efficiency of our method at high workload duration (time slot 10). Fig. 8d shows the amount of energy consumption of office tenants compared to the comfort energy settings. Using the workload migration, the MUB manager can control flexibly the energy consumption in buildings instead of handling the backup generator in a lot of times.

Evaluation of the office temperatures. In order to evaluate the impact of our mechanism on the office temperature, we illustrate the change of office temperatures during 20 time slots as shown in Fig. 9. Compared to Baseline 1 (without considering discomfort cost awareness), office tenants suffer lower temperatures during peak workload period (from time slot 10–15) when applying the GMUB mechanism. The efficiency of our mechanism is reflected well in all MUBs and significantly in MUB 5 with the largest setting.

Impact of the privacy-preserving protocol. Finally, we evaluate the privacy-preserving mechanism in our proposed Algorithm 2. We analyze the average execution time of the MUB tenants. Fig. 10 shows the comparison of two privacy-preserving protocols mentioned in Section 5. In Fig. 10a, we evaluate the average execution time of an MUB tenant in one iteration of Algorithm 2. The average computation time of the distributed aggregation protocol is 0.029 s, while the non-directional ring protocol takes 0.037 s to execute. Due to the sequential computation according to the ring as mentioned in Section 5, the average execution time of this method is higher than the distributed aggregation protocol. Therefore, the distributed protocol in our approach converges faster in Algorithm 1 than the non-directional ring method by about 1.24 s, as shown in Fig. 10b.

7. Conclusion

In this paper, we studied the well-known problem of edge computing providers who house edge datacenters located in geo-distribution locations. All the tenants in buildings are managed by the same entity to participate an EDR program. We investigate an optimization problem to formulate for multiple MUB tenants (including non-datacenter tenants and datacenter tenants) under their specific constraints.

This coordination problem is transformed into a solvable ADMM form. Based on our proposed mechanism, the delay-sensitive workloads can be appropriately distributed to the datacenter tenants, and is used to adjust the energy consumption of MUBs such that the total cost incurred in the system is minimized. Furthermore, we design an iterative distributed algorithm that can secure all tenants communication data by a privacy-preserving mechanism. Finally, we conducted many case studies to validate our proposed mechanism, for which the results show that our mechanism can reduce significantly the total cost while outperforming uncoordinated or partially coordinated approaches.

Acknowledgments

This work is partly funded by the NSERC Canada Research Chair on Sustainable Smart Eco-Cloud for Prof. Mohamed Cheriet. This work was supported in part by the U.S. NSF under grants CNS-1551661 and ECCS-1610471. This work was supported in part by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education under Grant NRF-2016R1D1A1B01015320.

References

- [1] P. Siano, Demand response and smart grids – a survey, *Renew. Sustain. Energy Rev.* 30 (2014) 461–478, <http://dx.doi.org/10.1016/j.rser.2013.10.022> (URL <<http://www.sciencedirect.com/science/article/pii/S1364032113007211>>).
- [2] H.-p. Chao, Price-responsive demand management for a smart grid world, *Electr. J.* 23 (1) (2010) 7–20.
- [3] US Federal Energy Regulatory Commission. URL <<http://www.ferc.gov/>>.
- [4] P. Interconnection, Emergency demand response (load management) performance report 2012/2013, PJM. Retrieved March 13 (2012) 2014.
- [5] Q. Sun, S. Ren, C. Wu, Z. Li, An online incentive mechanism for emergency demand response in geo-distributed colocation data centers, in: *ACM e-Energy, e-Energy '16*, 2016.
- [6] EnerNOC, Ensuring U.S. grid security and reliability: U.S. EPA's proposed emergency backup generator rule (2013).
- [7] F. Jazizadeh, A. Ghahramani, B. Becerik-Gerber, T. Kichkaylo, M. Orosz, Human-building interaction framework for personalized thermal comfort-driven systems in office buildings, *J. Comput. Civ. Eng.* 28 (1) (2014) 2–16, [http://dx.doi.org/10.1061/\(ASCE\)CP.1943-5487.0000300](http://dx.doi.org/10.1061/(ASCE)CP.1943-5487.0000300).
- [8] M. Lin, A. Wierman, L.L. Andrew, E. Thereska, Dynamic right-sizing for power-proportional data centers, *IEEE/ACM Trans. Netw. (TON)* 21 (5) (2013) 1378–1391.
- [9] L.A. Barroso, J. Clidaras, U. Hözl, The datacenter as a computer: an introduction to the design of warehouse-scale machines, *Synth. Lect. Comput. Archit.* 8 (3) (2013) 1–154.
- [10] Z. Zhou, F. Liu, Z. Li, H. Jin, When smart grid meets geo-distributed cloud: An auction approach to datacenter demand response, in: *Proceedings of IEEE Conference on Computer Communications (INFOCOM)*, 2015, pp. 2650–2658. <http://dx.doi.org/10.1109/INFOCOM.2015.7218656>.
- [11] Y. Li, D. Chiu, C. Liu, L. Phan, Towards dynamic pricing-based collaborative optimizations for green data centers, in: *IEEE ICDEW, Brisbane, Australia*, 2013, pp. 272–278.
- [12] Z. Liu, M. Lin, A. Wierman, S. Low, L.L.H. Andrew, Greening geographical load balancing, *IEEE/ACM Trans. Netw.* 23 (2) (2015) 657–671, <http://dx.doi.org/10.1109/TNET.2014.2308295>.

- [13] V. Avelar, D. Azevedo, A. French, E.N. Power, *Pue™: A comprehensive examination of the metric*.
- [14] Data center efficiency assessment. URL <https://www.nrdc.org/sites/default/files/data-center-efficiency-assessment-IP.pdf>.
- [15] S. Ren, Y. He, F. Xu, Provably-efficient job scheduling for energy and fairness in geographically distributed data centers, in: Proceedings of the 32nd IEEE International Conference on Distribution Computer System (1) (2012) 22–31. <http://dx.doi.org/10.1109/ICDCS.2012.77>.
- [16] H. Xu, C. Feng, B. Li, Temperature aware workload management in geo-distributed datacenter, ACM SIGMETRICS Perform. Eval. Rev. 41 (1) (2013) 373–374, <http://dx.doi.org/10.1145/2494232.2465539>.
- [17] S. Ren, M. A. Islam, Colocation demand response: Why do I turn off my servers?, in: USENIX ICAC, Philadelphia, PA, 2014, pp. 201–208.
- [18] N.H. Tran, C. Pham, S. Ren, C.S. Hong, Coordinated energy management for emergency demand response in mixed-use buildings, in: Ubiquitous Wireless Broadband (ICUWB), 2015 IEEE International Conference on, IEEE, 2015, pp. 1–5.
- [19] N. H. Tran, C. Pham, M. N. H. Nguyen, S. Ren, C. S. Hong, Incentive695 aligned mechanism for emergency demand response in multi-tenant mixed-use buildings, in: Proceedings of INFOCOM Workshops, 2016, pp. 112–117. doi:10.1109/INFCOMW.2016.7562056.
- [20] M. Lin, Z. Liu, A. Wierman, L.L.H. Andrew, Online algorithms for geographical load balancing, Proc. IGCC 2012 (2012) 1–10, <http://dx.doi.org/10.1109/IGCC.2012.6322266>.
- [21] Z. Zhou, F. Liu, Z. Li, H. Jin, When smart grid meets geo-distributed cloud: An auction approach to datacenter demand response, in: IEEE Conference on Proceedings of Computer Communications (INFOCOM), 2015, pp. 2650–2658.
- [22] N. Chen, X. Ren, S. Ren, A. Wierman, Greening multi-tenant data center demand response, in: IFIP WG 7.3 Performance, Sydney, Australia, 2015.
- [23] M. Islam, K. Ahmed, H. Xu, N. Tran, G. Quan, S. Ren, Exploiting spatio-temporal diversity for water saving in geo-distributed data centers, IEEE Trans. Cloud Comput. 99 (2016) 1, <http://dx.doi.org/10.1109/TCC.2016.2535201>.
- [24] S. Godo, K. Matsui, H. Nishi, Cost-effective air conditioning control considering comfort level and user location, in: Proceedings of the 40th Annual Conference of the IEEE Industrial Electronics Society, IECON 2014, 2014, pp. 5344–5349. <http://dx.doi.org/10.1109/IECON.2014.7049316>.
- [25] A. Greenberg, J. Hamilton, D.A. Maltz, P. Patel, The cost of a cloud: research problems in data center networks, ACM SIGCOMM Comput. Commun. Rev. 39 (1) (2008) 68–73.
- [26] D. Irwin, S. Barker, A. Mishra, P. Shenoy, A. Wu, J. Albrecht, Exploiting home automation protocols for load monitoring in smart buildings, BuildSys (2011).
- [27] A. Mishra, D. Irwin, P. Shenoy, J. Kurose, T. Zhu, Greencharge: managing renewable energy in smart buildings, IEEE J. Sel. Areas Commun. 31 (7) (2013) 1281–1293.
- [28] M. Maasoumy, Modeling and optimal control algorithm design for hvac systems in energy efficient buildings.
- [29] Reducing costs by consolidation strategies. URL <http://www.oracle.com/us/products/servers-storage/servers/sparc-enterprise/reducing-costs-wp-075962.pdf>.
- [30] Y. Xiang, Z. Liu, X. Qu, Cpu frequency scaling algorithm for energy-saving in cloud data centers, J. Softw. 9 (9) (2014) 2283–2290.
- [31] Q. Wu, Making facebook's software infrastructure more energy efficient with autoscale (2014).
- [32] A. Qureshi, R. Weber, H. Balakrishnan, J. Gutttag, B. Maggs, Cutting the electric bill for internet-scale systems, SIGCOMM (2009).
- [33] L. Rao, X. Liu, L. Xie, W. Liu, Reducing electricity cost: optimization of distributed internet data centers in a multi-electricity-market environment, INFOCOM (2010).
- [34] C. Ren, D. Wang, B. Urgaonkar, A. Sivasubramanian, Carbon-aware energy capacity planning for datacenters, in: Proc. of MASCOTS, 2012, 740 pp. 391–400. doi:10.1109/MASCOTS.2012.51.
- [35] N.H. Tran, C. Pham, M.N. Nguyen, S. Ren, C.S. Hong, Incentivizing energy reduction for emergency demand response in multi-tenant mixed-use buildings, GSNIC16 (2016).
- [36] C. Clifton, M. Kantarcioglu, J. Vaidya, X. Lin, M.Y. Zhu, Tools for privacy preserving distributed data mining, ACM Sigkdd Explor. Newsl. 4 (2) (2002) 28–34.
- [37] J. Vaidya, C. Clifton, Privacy preserving association rule mining in vertically partitioned data, in: Proceedings of the Eighth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2002, pp. 639–644.
- [38] A. Evfimievski, R. Srikant, R. Agrawal, J. Gehrke, Privacy preserving mining of association rules, 2002, pp. 217–228.
- [39] J.W. Bos, K. Lauter, M. Naehrig, Private predictive analysis on encrypted medical data, JBI 50 (2014) 234–243.
- [40] S. Ren, M.A. Islam, Colocation demand response: Why do i turn off my servers, Proceedings of USENIX ICAC.
- [41] L. Zhang, S. Ren, C. Wu, Z. Li, A truthful incentive mechanism for emergency demand response in colocation data centers, in: IEEE INFOCOM, Hong Kong, China, 2015.
- [42] A. Beloglazov, R. Buyya, Y.C. Lee, A. Zomaya, Chapter 3 - a taxonomy and survey of energy-efficient data centers and cloud computing systems, Vol. 82 of Advances in Computers, Elsevier, 2011, pp. 47 – 111. <http://dx.doi.org/10.1016/B978-0-12-385512-1.00003-7>. URL <http://www.sciencedirect.com/science/article/pii/B9780123855121000037>.
- [43] Z. Liu, M. Lin, A. Wierman, S.H. Low, L.L.H. Andrew, Geographical load balancing with renewables, ACM SIGMETRICS 39 (3) (2011) 62–66, <http://dx.doi.org/10.1145/2160803.2160862>.
- [44] A. Gandhi, M. Harchol-Balter, R. Das, C. Lefurgy, Optimal power allocation in server farms, ACM SIGMETRICS Perform. Eval. Rev. 37 (1) (2009) 157–168, <http://dx.doi.org/10.1145/2492101.1555368>.
- [45] M. Lin, A. Wierman, L. L. H. Andrew, E. Thereska, Dynamic Right-sizing for Power-proportional Data Centers, in: Proceedings IEEE INFOCOM, Shanghai, China, 2011. doi:10.1109/INFCOM.2011.5934885.
- [46] H. Qian, D. Medhi, Server operational cost optimization for cloud computing service providers over a time horizon, in: Proceedings of the 11th USENIX Conference, USENIX Association, 2011, pp. 4–4.
- [47] E. Vrettos, J.L. Mathieu, G. Andersson, Demand response with moving horizon estimation of individual thermostatic load states from aggregate power measurements, in: Proceedings of American Control Conference, 2014, pp. 4846–4853. <http://dx.doi.org/10.1109/ACC.2014.6859068>.
- [48] A.H.-y. Lam, Y. Yuan, D. Wang, An occupant-participatory approach for thermal comfort enhancement and energy conservation in buildings, in: Proceedings ACM e-Energy, New York, USA, 2014, pp. 133–143. <http://dx.doi.org/10.1145/2602044.2602067>.
- [49] S. Li, W. Zhang, J. Lian, K. Kalsi, Market-based coordination of thermostatically controlled loads-part I: A mechanism design formulation, in: Proceedings of IEEE Power and Energy Society General Meeting (PESGM), 2016, pp. 1–1.
- [50] J. C. Fuller, K. P. Schneider, D. Chassin, Analysis of residential demand790 response and double-auction markets, in: IEEE Power and Energy Society General Meeting, 2011, pp. 17. doi:10.1109/PES.2011.6039827.
- [51] N. Combe, D. Harrison, S. Craig, M.S. Young, An investigation into usability and exclusivity issues of digital programmable thermostats, J. Eng. Des. 23 (5) (2012) 401–417, <http://dx.doi.org/10.1080/09544828.2011.599027> (arXiv:10.1080/09544828.2011.599027). URL <https://doi.org/10.1080/09544828.2011.599027>.
- [52] F. Auffenberg, S. Stein, A. Rogers, A personalised thermal comfort model using a bayesian network.
- [53] A. Wood, B. Wollenberg, G. Shebl, Power Generation, Operation, and Control, Wiley, 2013. URL <https://books.google.ca/books?id=JDVmAgAAQBAJ>.
- [54] Z. Liu, M. Lin, A. Wierman, S. Low, L.L.H. Andrew, Greening geographical load balancing, IEEE/ACM Trans. Netw. 23 (2) (2015) 657–671, <http://dx.doi.org/10.1109/TNET.2014.2308295>.
- [55] D.P. Palomar, M. Chiang, Alternative distributed algorithms for network utility maximization: framework and applications, IEEE Trans. Autom. Control 52 (12) (2007) 2254–2269, <http://dx.doi.org/10.1109/TAC.2007.910665>.
- [56] S. Boyd, Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers, Found. Trends Mach. Learn. 3 (1) (2010) 1–122, <http://dx.doi.org/10.1561/22000000016>.
- [57] W. Deng, M.-J. Lai, Z. Peng, W. Yin, Parallel multi-block admm with $O(1/K)$ convergence, J. Sci. Comput. 71 (2) (2017) 712–736.
- [58] S. Boyd, L. Vandenberghe, Convex Optimization, Cambridge University Press, 2004.
- [59] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, Distributed optimization and statistical learning via the alternating direction method of multipliers, Found. Trends Mach. Learn. 3 (1) (2011) 1–122.
- [60] E. Shi, T.H. Chan, E. Rieffel, R. Chow, D. Song, Privacy-preserving aggregation of time-series data.
- [61] L. Minas, B. Ellison, Energy Efficiency for Information Technology: How to Reduce Power Consumption in Servers and Data Centers, Intel Press, 2009.
- [62] Facebook, Open sourcing pue, wue dashboards <https://code.facebook.com/posts/272417392924843/>.
- [63] N.H. Tran, C.T. Do, S. Ren, Z. Han, C.S. Hong, Incentive Mechanisms for Economic and Emergency Demand Responses of Colocation Datacenters, IEEE JSAC 33. <http://dx.doi.org/10.1109/JSAC.2015.2483420>.