

## Distributed Foresighted Energy Management in Smart-Grid-Powered Cellular Networks

Xinruo Zhang , Mohammad Reza Nakhai , Gan Zheng , Sangarapillai Lambotharan , and Jonathon A. Chambers

**Abstract**—This paper studies energy management in a smart grid-powered cellular network consisting of an independent system operator (ISO) and multiple geographically distributed aggregators. The aggregators have energy storage devices and can purchase energy from the electric grid via the ISO to serve their users. To account for the uncertainty of the renewable energy supply as well as the impacts of multiple aggregators on the electric grid and energy prices, a foresighted strategy combined with the adaptive  $\epsilon$ -greedy method is developed for the aggregators to distributively and adaptively minimize the long-term overall cost of the system, based on the ahead-of-time decision making of the storage pre-charging amount. Simulation results validate that the proposed strategy surpasses a recent learning-based storage management design and a myopic design.

**Index Terms**—Distributed management, energy management, online learning.

### I. INTRODUCTION

The enormous energy consumption in the next generation of dense wireless networks has been considered as one of the most challenging issues from both technical and economic perspectives. The integration of renewable energy generated from natural sources with the conventional electric grid is viewed as a promising approach towards achieving the targets of green communications and is implemented in the smart grid infrastructure [1]. The stochastic nature of renewable energy, nevertheless, causes significant uncertainty in energy generation and fluctuations of the electricity price. Such uncertainty could lead to abrupt ramping in the power plants or the adoption of peaker plants that use non-renewable sources of energy to compensate for random variations in energy generation, which is one of the most expensive operations or even may not be technically feasible for some energy generators [2]. Hence, efficient new control mechanisms are advocated to improve the flexibility and robustness of smart grid networks, so that the independent system operator (ISO) can maintain its reliable and cost-efficient operation under such variability. Furthermore, the deployment of energy storage units in the demand side has also been proposed as a viable solution to tackle these concerns [2], as it can not only compensate for the real-time energy shortage, but also be used to minimize the long-term energy consumption cost via being proactively pre-charged either from the power grid at a lower price or by the

excessive renewable energy generation. Most approaches in the literature [3]–[5] propose myopic demand side energy management designs based on an optimization method and seek for instantaneous/short-term cost minimization without taking into account either the deployment of energy storage units or the possibility of learning the system dynamics during the operation. The authors in [6] concentrate on the interaction between a single aggregator and customers and develop an online learning algorithm for stochastic storage management based on the Markov decision process. Using stochastic optimization rather than online learning, the authors in [7] propose a dynamic energy management design for the smart-grid-powered coordinated multipoint system. However, the stochastic optimization based approaches usually require a model or the statistics of the system dynamics to be known upfront. Furthermore, [6]–[8] focus merely on a single aggregator, while the impacts of the decisions of the aggregators on the power network have been neglected.

Accounting for the intermittent renewable energy generation and the impacts of aggregators on the power network, this paper focuses on designing an online energy management strategy to minimize the long-term average energy consumption cost of the ISO in a decentralized manner. This is challenging because the ISO and the distributed aggregators that are in geographically different locations have no access to the local information of one another. Hence, the conventional bandit approaches typically designed for solving centralized problems can not tackle such a distributed problem with strictly limited signalling information. Furthermore, the statistics of the system dynamics such as renewable energy generation are unknown *a priori*, and the storage pre-charging decisions have strong temporal correlations, which render the problem intractable for traditional stochastic optimization based approaches. Hence, the novelty of this paper is the introduction of a distributed online foresighted energy management strategy that requires no upfront knowledge and can optimize the long-term average energy cost while learning via distributively alternating between two decision making processes. The first process minimizes the current cost at each aggregator based on the storage pre-charging amount and local user demand using convex optimization. The second process designs the ahead-of-time storage pre-charging strategies at distributed aggregators via online learning.

### II. SYSTEM MODEL

As illustrated in Fig. 1, this paper considers a smart grid powered cellular network with one ISO and a set of geographically distributed aggregators with central processing units (A-CPU), indexed by  $\mathcal{L}_a = \{1, \dots, N\}$ . Each aggregator is associated with a cloud radio access network (C-RAN) of  $N_b$  base stations (BSs) serving  $N_i$  users, indexed as  $\mathcal{L}_i = \{1, \dots, N_i\}$ , using a shared spectrum. Featuring smart grid operations, the individual A-CPU are deployed with on-site energy storage units and can exploit renewable energy sources with harvesting facilities. Furthermore, the A-CPU have access to ancillary grid energy markets at various prices via energy trading mechanisms. The ISO operates the system through receiving energy purchase requests from the A-CPU and dispatching the power grid based on its operational status. In particular, the ISO only has access to the status information of the grid as well as the energy purchase requests from the A-CPU, whilst each A-CPU only has the status information of its associated C-RAN cluster and its storage units.

Manuscript received May 8, 2018; revised July 11, 2018 and January 21, 2019; accepted February 8, 2019. Date of publication February 18, 2019; date of current version April 16, 2019. This work was supported in part by the UK EPSRC under Grant EP/N007840/1 and in part by the Leverhulme Trust under Grant RPG-2017-129. The review of this paper was coordinated by Prof. Y. Guo. (Corresponding author: Gan Zheng.)

X. Zhang, G. Zheng, and S. Lambotharan are with the Wolfson School of Mechanical, Electrical and Manufacturing Engineering, Loughborough University, LE11 3TU Loughborough, U.K. (e-mail: x.zhang@lboro.ac.uk; g.zheng@lboro.ac.uk; s.lambotharan@lboro.ac.uk).

M. R. Nakhai is with the Centre for Telecommunications Research, King's College London, WC2B 4BG London, U.K. (e-mail: reza.nakhai@kcl.ac.uk).

J. A. Chambers is with the Department of Engineering, University of Leicester, LE1 7RH Leicester, U.K. (e-mail: jonathon.chambers@leicester.ac.uk).

Digital Object Identifier 10.1109/TVT.2019.2899464

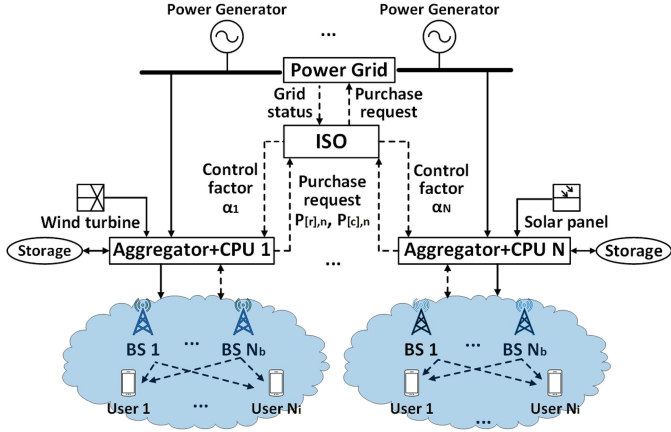


Fig. 1. Illustration of system scenario. The information flow is denoted by dashed lines and the energy flow is denoted by solid lines.

Let the time horizon  $T$  be divided into discrete time slots and indexed as  $\mathcal{T} = \{1, \dots, T\}$ . Assume that the renewable energy generation varies across time slots but remains invariant within each time slot. For convenience, the slot duration is normalized to unity, thus the terms ‘energy’ and ‘power’ are used synonymously throughout the paper. Let  $P_{[g],n}^t$  and  $P_{[r],n}^t$  denote, respectively, the amount of renewable energy generation and the amount of energy shortage to be supplied by the grid in real-time at time slot  $t$ ,  $t \in \mathcal{T}$ , to the  $n^{\text{th}}$  A-CPU. Let  $P_{[s],n}^t$  be the initial amount of energy stored in the storage device at time slot  $t$  and  $P_{[c],n}^t$  denote the amount of pre-charged energy stored in the  $n^{\text{th}}$  A-CPU prior to the actual time of energy demand at the beginning of the  $t^{\text{th}}$  time slot. The total energy consumption cost incurred by the  $n^{\text{th}}$  A-CPU at time slot  $t$ , i.e.,  $C_n^t$ , can be described as

$$C_n^t = c_{[r]}^t P_{[r],n}^t + c_{[c]}^t P_{[c],n}^t + c_{[g]}^t P_{[g],n}^t + c_{[s]}^t P_{[s],n}^t, \quad (1)$$

where  $c_{[r]}^t$  and  $c_{[c]}^t$  are, respectively, the per unit energy prices for  $P_{[r],n}^t$  and  $P_{[c],n}^t$  that will be updated by the ISO.  $c_{[g]}^t$  and  $c_{[s]}^t$  are the per unit equivalent annual cost of renewable harvesters and storage devices, respectively. It is assumed that the grid status and the real-time energy production price  $c_{[r]}^t$  are under the influence of the real-time energy production level in the previous time slot [9], i.e.,  $\{P_{[r],n}^{t-1}\}$ . Similar to [9], the real-time energy price of the grid can be modelled as

$$c_{[r]}^{t+1} \propto \underbrace{0.5 \sum_{n=1}^N |P_{[r],n}^t|}_{\text{generation price}} + \underbrace{0.1 \left( \sum_{n=1}^N \left[ |P_{[r],n}^t| - |P_{[r],n}^{t-1}| \right]^+ \right)^2 / \sum_{n=1}^N |P_{[r],n}^t|}_{\text{ramping price}}. \quad (2)$$

Let us denote by  $P_{[Tx],n}^t = \sum_{i \in \mathcal{L}_i} \|\mathbf{w}_{n,i}^t\|^2$  and  $P_{[h],n}^t$  the total transmit power consumption of the  $n^{\text{th}}$  A-CPU at the  $t^{\text{th}}$  time slot and the hardware circuit power consumption, respectively, where  $\mathbf{w}_{n,i}^t$  is the beamforming vector from all BSs in the  $n^{\text{th}}$  A-CPU to the  $i^{\text{th}}$  user. Per time slot  $t$ , the total energy consumption of the  $n^{\text{th}}$  A-CPU,  $n \in \mathcal{L}_a$ , is bounded as

$$\xi P_{[Tx],n}^t + P_{[h],n}^t \leq P_{[g],n}^t + P_{[s],n}^t + P_{[c],n}^t + P_{[r],n}^t, \quad (3)$$

where  $\xi > 0$  is the power amplifier efficiency. We simplify the constraints of the storage devices to the capacity limit only and denote by  $P_{[cap],n}^t$  the finite capacity of the  $n^{\text{th}}$  storage device. As shown in Fig. 2,

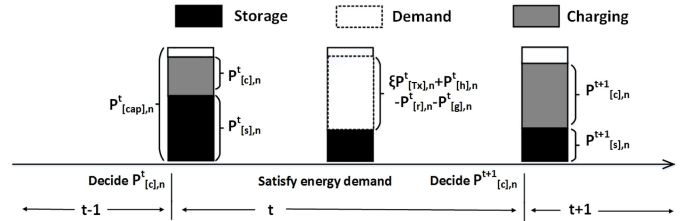


Fig. 2. Illustration of storage charging and discharging processes.

the initial amount of stored energy of the  $n^{\text{th}}$  A-CPU at time slot  $t$  is constrained as follows

$$P_{[s],n}^t = \min \left\{ [P_{[g],n}^{t-1} + P_{[s],n}^{t-1} + P_{[c],n}^{t-1} + P_{[r],n}^{t-1} - P_{[Tx],n}^{t-1} - P_{[h],n}^{t-1}]^+, P_{[cap],n}^t \right\}. \quad (4)$$

### III. FORESIGHTED ENERGY MANAGEMENT STRATEGY

#### A. Problem Formulation

The ISO aims to minimize the long-term overall energy consumption cost of the entire network, whilst the individual A-CPU need to satisfy the signal-to-interference-plus-noise ratio (SINR) requirements of their users in the presence of variability in wireless channels as well as in the power grid, e.g., uncertain renewable energy generation and real-time energy price [9]. Hence, the problem of interest covering the entire cost control in the power network and the beamforming design in the cellular network can be formulated as

$$\begin{aligned} & \min_{\{P_{[c],n}^t, P_{[r],n}^t, \mathbf{w}_{n,i}^t\}} \left\{ \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{n=1}^N C_n^t \right\} \\ \text{s.t. } & \text{C1: } \text{SINR}_i(\{\mathbf{w}_{n,i}^t\}) \geq \gamma_i, \quad \forall i \in \mathcal{L}_i, n \in \mathcal{L}_a, t \in \mathcal{T}, \\ & \text{C2: } \xi P_{[Tx],n}^t + P_{[h],n}^t - P_{[g],n}^t - P_{[s],n}^t - P_{[c],n}^t \leq P_{[r],n}^t, \\ & \quad \quad \quad \forall n \in \mathcal{L}_a, t \in \mathcal{T}, \\ & \text{C3: } P_{[s],n}^t = \min \left\{ [P_{[g],n}^{t-1} + P_{[s],n}^{t-1} + P_{[c],n}^{t-1} + P_{[r],n}^{t-1} \right. \\ & \quad \quad \quad \left. - P_{[Tx],n}^{t-1} - P_{[h],n}^{t-1}]^+, P_{[cap],n}^t \right\}, \quad \forall n \in \mathcal{L}_a, t \in \mathcal{T}, \quad (5) \end{aligned}$$

where the constraint C1 indicates the minimum SINR requirements  $\{\gamma_i\}$  for the users. Let us denote by  $C_n^0$  the total cost incurred by the  $n^{\text{th}}$  A-CPU at time slot 0, where no learning or battery pre-charging has yet been applied in the system. Then, we can consider  $C_n^0$  as a constant reference point and define the reward for the  $n^{\text{th}}$  A-CPU at time slot  $t$ , as

$$\mathcal{R}_n^t = C_n^0 - C_n^t, \quad \forall n \in \mathcal{L}_a, t \in \mathcal{T}. \quad (6)$$

Then, the problem of minimizing long-term overall energy cost in (5) is equivalent to the following optimization problem that maximizes the time-averaged accumulated reward, as

$$\begin{aligned} & \max_{\{P_{[c],n}^t, P_{[r],n}^t, \mathbf{w}_{n,i}^t\}} \left\{ \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{n=1}^N \mathcal{R}_n^t \right\} \\ \text{s.t. } & \text{C1–C3}. \quad (7) \end{aligned}$$

## B. Optimization-Assisted Sequential Decision Making Problem

The stochastic optimization problem in (7) is difficult to solve directly since we aim to minimize the long-term average system cost while the statistics of the system dynamics are unknown in advance. In particular, the time-coupling storage constraints C2 and C3 render the problem intractable for traditional solvers. Furthermore, both the ISO and the A-CPU only have access to their respective local information. For instance, the storage status is only available to its local A-CPU whilst the energy purchase decisions of the individual A-CPU are unknown to the others and will affect the overall cost of the system. Hence, the problem cannot be solved in a centralized manner in either the ISO or the A-CPU through the traditional stochastic optimization based approaches.

Thus, we employ the bandit approach that requires no upfront knowledge of the system dynamics to account for the time-coupled variables, and consider (7) as an optimization-assisted sequential decision making problem that alternates between two decision making processes and indirectly affects the decisions of the A-CPU slot-by-slot, so as to distributively and asymptotically maximize the average reward in the long run. More specifically, the online learning part of the proposed strategy aims at proactive one-slot-ahead decision making of the storage pre-charging amount for time slot  $t$ , i.e.,  $P_{[c],n}^t$ , prior to any possible random variations in the renewable generation or real-time energy-shortage at time slot  $t$ . During time slot  $t$  and based on the one-slot-ahead decision of  $P_{[c],n}^t$ , the individual A-CPU will then make decisions on the real-time energy request, i.e.,  $P_{[r],n}^t$ , for satisfying energy demand of the users via solving the following optimization problem distributively in real-time, as

$$\begin{aligned} & \min_{\{\mathbf{w}_{n,i}^t, P_{[r],n}^t \geq 0\}} P_{[Tx],n}^t + \alpha_n P_{[r],n}^t \\ \text{s.t. } & \text{SINR}_i^t = \frac{|\mathbf{h}_i^H \mathbf{w}_{n,i}^t|^2}{\sum_{j \neq i, j \in \mathcal{L}_i} |\mathbf{h}_i^H \mathbf{w}_{n,j}^t|^2 + \sigma_i^2} \geq \gamma_i, \forall i \in \mathcal{L}_i, \\ & \xi P_{[Tx],n}^t + P_{[h],n}^t \leq P_{[g],n}^t + P_{[s],n}^t + P_{[c],n}^t + P_{[r],n}^t, \end{aligned} \quad (8)$$

where  $\mathbf{h}_i$  represents the channel vector from all BSs in the  $n^{\text{th}}$  A-CPU to the  $i^{\text{th}}$  user,  $i \in \mathcal{L}_i$ , and  $P_{[s],n}^t$  can be updated in the light of (4). The objective function of problem (8) seeks an optimal schedule for beamforming vectors  $\{\mathbf{w}_{n,i}^t\}$  and the real-time purchase of energy shortage  $P_{[r],n}^t$ , in order to minimize the total cost  $C_n^t$  incurred by the  $n^{\text{th}}$  A-CPU, given the one-slot-ahead decision of  $P_{[c],n}^t$ . In particular, the penalty factor  $\alpha_n$  in the objective function in (8) is assigned and updated by the ISO according to  $\alpha_n \propto \frac{P_{[r],n}^{t-1}}{\min_{n \in \mathcal{L}_a} \{P_{[r],n}^{t-1}\}}$  at the beginning of time slot  $t$ . This is to account for the grid status at the previous time slot and to control the impacts of the A-CPU on the grid energy prices at the current time slot. For instance, a larger  $\alpha_n$  emphasizes more on the minimization of peak-time energy purchases for those A-CPU with larger real-time energy shortage in the previous time slot, in order to bring the higher real-time electricity price under control. Furthermore,  $\alpha_n$  will also be employed in the proposed online learning algorithm to penalize the reward and indirectly affect the one-slot-ahead decisions of the A-CPU, i.e.,  $\{P_{[c],n}^t\}$ .

To solve problem in (8), let us define the rank-one semidefinite matrices  $\mathbf{W}_{n,i}^t = \mathbf{w}_{n,i}^t \mathbf{w}_{n,i}^{tH}$  and  $\mathbf{H}_i = \mathbf{h}_i \mathbf{h}_i^H$ . The problem in (8) can be transformed to a tractable form after relaxing the rank-one constraint

### Algorithm 1: Foresighted Energy Management Algorithm.

- 1: **Initialize:**  $\epsilon^0 = 1$ ,  $P_{[s],n}^1 = 0$ ,  $P_{[c],n}^0 = 0$ , estimated mean reward  $\bar{\mathbf{r}}_n^t = [\bar{r}_{n,1}^t, \dots, \bar{r}_{n,E}^t] = \mathbf{0}$ , temporary reward  $\mathbf{r}_n^{[k,t]} = [r_{n,1}^{[k,t]}, \dots, r_{n,E}^{[k,t]}] = \mathbf{0}$ , total numbers of actions  $E$  with step of  $\Delta$ , time slots  $T$  and learning trials  $K$ .
- 2: **REPEAT**
- 3: With the probability of  $\epsilon^t$ :  $t$  is for **Exploration**
- 4: **While**  $k \leq K$
- 5: **if**  $k = 1$ 
  - with probability of  $\epsilon^t$ , randomly select an action  $P_{[c],n}^{t,1}$ ,
  - with probability of  $1 - \epsilon^t$ , select  $P_{[c],n}^{t,1}$  as  $P_{[c],n}^{t,1} = \operatorname{argmax}_e(\bar{\mathbf{r}}_n^t)$ ,  $e \in \mathcal{E}$ .
- 6: **end if**
- 7: Solve (9) and compute  $\mathcal{R}(P_{[c],n}^{t,k})$  as per (10).
- 8: **if**  $k = 2$  or  $\mathcal{R}(P_{[c],n}^{t,k}) > \mathcal{R}(P_{[c],n}^{t,k-1})$ 
  - $P_{[c],n}^{t,k+1} = P_{[c],n}^{t,k} + \Delta$ .
- 9: **else if**  $\mathcal{R}(P_{[c],n}^{t,k}) < \mathcal{R}(P_{[c],n}^{t,k-1})$ 
  - $P_{[c],n}^{t,k+1} = P_{[c],n}^{t,k} - \Delta$ .
- 10: **else**  $P_{[c],n}^{t,k+1} = P_{[c],n}^{t,k}$ .
- 11: **end if**
- 12: Update the temporary reward matrix  $\mathbf{r}_n^{[k,t]}$ , as  $r_{n,e}^{[k,t]} = \mathcal{R}(P_{[c],n}^{t,k})$ ,  $e = \frac{P_{[c],n}^{t,k}}{\Delta} \in \mathcal{E}$ .
- 13: Update  $k = k + 1$ .
- 14: **End While**
- 15: With the probability of  $1 - \epsilon^t$ :  $t$  is for **Exploitation**
- 16: Select best action  $P_{[c],n}^t = \operatorname{argmax}_e(\bar{\mathbf{r}}_n^t)$ ,  $e \in \mathcal{E}$ .
- 17: Solve problem (9) and compute  $\mathcal{R}(P_{[c],n}^t)$  as per (10).
- 18: Compute the estimated mean reward vector, as 
$$\bar{\mathbf{r}}_n^t = \frac{\sum_{t'=1}^t \left( \frac{\sum_{k=1}^K r_{n,1}^{[k,t']}}{K}, \frac{\sum_{k=1}^K r_{n,2}^{[k,t']}}{K}, \dots, \frac{\sum_{k=1}^K r_{n,E}^{[k,t']}}{K} \right) \beta^{(t-t')}}{t}$$
, where  $\beta$  is the discount factor.
- 19: Update  $t = t + 1$ .
- 20: **UNTIL**  $t = T$

of  $\operatorname{rank}(\mathbf{W}_{n,i}^t) = 1$ , as

$$\begin{aligned} & \min_{\{\mathbf{W}_{n,i}^t, P_{[r],n}^t \geq 0\}} \sum_{i \in \mathcal{L}_i} \operatorname{tr}(\mathbf{W}_{n,i}^t) + \alpha_n P_{[r],n}^t \\ \text{s.t. } & \gamma_i^{-1} \operatorname{tr}(\mathbf{H}_i \mathbf{W}_{n,i}^t) \geq \sum_{j \in \mathcal{L}_i, j \neq i} \operatorname{tr}(\mathbf{H}_i \mathbf{W}_{n,j}^t) + \sigma_i^2, \forall i \in \mathcal{L}_i, \\ & \xi \sum_{i \in \mathcal{L}_i} \operatorname{tr}(\mathbf{W}_{n,i}^t) + P_{[h],n}^t - P_{[s],n}^t - P_{[c],n}^t - P_{[g],n}^t \leq P_{[r],n}^t. \end{aligned} \quad (9)$$

It can be proved that the optimal solutions to problem (9) satisfy the rank-one constraint and therefore are also optimal for (8). The proof is similar to that in [4] and thus omitted.

## C. The Proposed Foresighted Algorithm

Next we turn to the online learning part of the proposed strategy to decide the one-slot-ahead storage pre-charging amount at individual A-CPU. We employ the multi-armed bandit model [10], where the agents are the A-CPU and a total number of  $E$  actions correspond to the  $E$  discrete amounts of energy with step of  $\Delta$  that can be charged as storage within an A-CPU prior to the occurrence of any possible energy

shortage. Let us denote by  $\mathcal{E} = \{1, \dots, E\}$  and  $\mathcal{A} = \{A_1, \dots, A_E\}$  the indexes and the set of  $E$  possible actions, respectively. The proposed foresighted algorithm aims to distributively and sequentially determine the one-slot-ahead storage pre-charging amount, i.e., select the best-possible action  $P_{[c],n}^{t,*} \in \mathcal{A}$ , before experiencing an energy shortage at time slot  $t$ , and adaptively react with the foresighted best response via solving the problem in (9). To accelerate the learning process, we implement  $K$  learning trials within a time slot, indexed by  $\mathcal{K} = \{1, \dots, K\}$ . Let the action selected to pre-charge the storage unit of the  $n^{\text{th}}$  A-CPU in the  $k^{\text{th}}$  trial of the  $t^{\text{th}}$  time slot,  $n \in \mathcal{L}_a, k \in \mathcal{K}, t \in \mathcal{T}$ , be denoted by  $P_{[c],n}^{t,k}$ . The instantaneous reward of  $P_{[c],n}^{t,k}$ , can be defined as

$$\mathcal{R} \left( P_{[c],n}^{t,k} \right) = \tilde{C}_n^{0,1} - \tilde{C}_n^{t,k}, \quad \forall n \in \mathcal{L}_a, t \in \mathcal{T}, \quad (10)$$

where  $\tilde{C}_n^{t,k}$  is the penalized total cost incurred by the  $n^{\text{th}}$  A-CPU and is defined as

$$\tilde{C}_n^{t,k} = \alpha_n c_{[r]}^t P_{[r],n}^{t,k} + c_{[c]}^t P_{[c],n}^{t,k} + c_{[g]}^t P_{[g],n}^t + c_{[s]}^t P_{[s],n}^{t,k}. \quad (11)$$

The penalized reward function in (10) urges the A-CPU with higher real-time energy shortage in the current round to pre-charge more energy into their storages in order to be prepared for the next round.

The proposed foresighted algorithm is governed by a trade-off between exploring other actions that may yield a better accumulated reward in the presence of uncertainties on the renewable energy supply, and exploiting current knowledge to make the empirically best decisions among a set of actions. The details of the proposed foresighted procedures to be executed at the individual A-CPU are described in Algorithm 1, where the individual time slots can either be allocated as an exploration cycle with the probability of  $\epsilon^t$ , or an exploitation cycle with the probability of  $1 - \epsilon^t$ , as explained below:

- **Exploration:** With the probability of  $\epsilon^t$ , a perturbation procedure is applied to explore actions uniformly at random, as the starting point of learning, i.e.,  $P_{[c],n}^{t,1} \in \mathcal{A}$ . With the probability of  $1 - \epsilon^t$ , the action associated with the highest mean reward so far will be selected as an initial learning point and will be gradually improved until a total number of  $K$  learning trials is attained.
- **Exploitation:** The individual A-CPU select the best-possible action that yields the highest estimated mean reward so far, based on the observed knowledge up to the  $(t - 1)^{\text{th}}$  time slot, and then feedback it to solve (9).

Instead of the traditional hand-tuning of  $\epsilon^t$  [8], we modified the value-difference based exploration method [11], such that the exploration-exploitation control parameter  $\epsilon^t, 0 \leq \epsilon^t \leq 1$ , is adaptive to the uncertainty in the learning progress, i.e.,

$$\epsilon^{t+1} = \delta \cdot \frac{1 - e^{-\frac{|\mathcal{R}_n^{t+1} - \frac{1}{t} \sum_{t'=1}^t \mathcal{R}_n^{t'}|}}{\sigma}}{1 + e^{-\frac{|\mathcal{R}_n^{t+1} - \frac{1}{t} \sum_{t'=1}^t \mathcal{R}_n^{t'}|}}{\sigma}} + (1 - \delta) \cdot \epsilon^t, \quad (12)$$

where  $\sigma$  is a positive constant indicating the inverse sensitivity and  $\delta \in (0, 1)$  [11]. More specifically, a time-decayed exploration rate can be adopted in a relative static environment, where the estimation of the mean reward process of the actions is improved with time. On the contrary, a relative high exploration rate can be employed when a sudden change in the environment or the reward is observed.

#### IV. SIMULATION RESULTS

Consider an ISO with 3 A-CPU, where each aggregator consists of a C-RAN of 3 BSs and 6 users. The renewable energy generation at each time slot varies as  $P_{[g],1}^t \in [0, 0.5]$  W,  $P_{[g],2}^t \in [0.3, 1.0]$  W

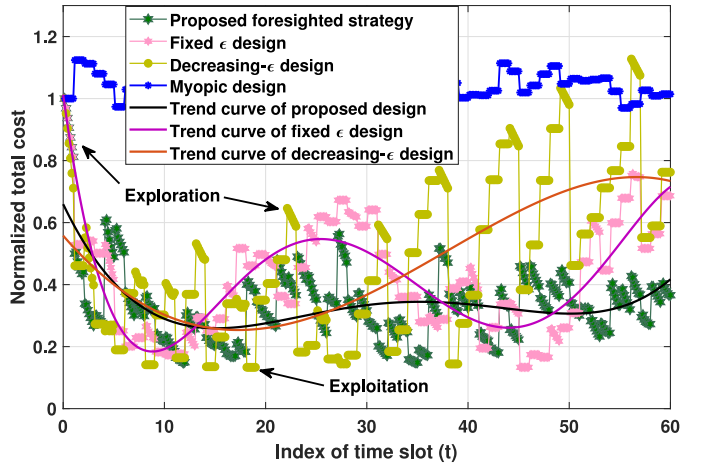


Fig. 3. Normalized total cost at individual time slots at  $\gamma = 15$  dB.

and  $P_{[g],3}^t \in [0.5, 1.5]$  W, respectively. The other simulation parameters are described as follows:  $E = 20$  with  $\Delta = 300$  mW,  $\beta = 0.95$ ,  $c_{[r]}^t = \mathcal{L}0.15/\text{W}$ ,  $c_{[c]}^t = \mathcal{L}0.07/\text{W}$ ,  $c_{[g]}^t = \mathcal{L}0.05/\text{W}$ ,  $c_{[s]}^t = \mathcal{L}0.01/\text{W}$ ,  $P_{[cap],n}^t = 40$  dBm,  $P_{[h],n}^t = 33$  dBm,  $\xi = 1$ ,  $\sigma = 5$  and  $\delta = \frac{1}{20}$ . The proposed strategy is evaluated with  $K = 7$  learning trials and  $T = 60$  time slots. Two designs are chosen as benchmark schemes: an adaptive storage management design in [8] that only considers a single A-CPU with fixed exploration-exploitation trade-off, and a baseline myopic design in [5] that myopically minimizes the current energy cost.

Fig. 3 presents the comparison of the normalized total costs when  $\gamma = 15$  dB. For fair comparison, identical constraints have been applied to all strategies and the overall cost is normalized to the cost at the initial time of the proposed strategy. The bursts at the beginning of each exploration cycle appear as a result of highly uncertain renewable energy generation as well as the exploration perturbation in step 5 of Algorithm 1. As seen in Fig. 3, the proposed strategy outperforms the baseline myopic design in [5]. Furthermore, the proposed strategy indicates smaller variations and better average reward with increasing time, as compared to that of the designs in [8] with fixed and with decreasing  $\epsilon$ , respectively. This is due to the negligence of the coupling effects among individual A-CPU in the nature of the design in [8], which provides poorer adaptation to the variations in the renewable energy generation as well as the real-time energy price in a decentralized scenario.

Fig. 4 illustrates the comparison of the proposed adaptive  $\epsilon$ -greedy method against a traditional hand-tuning fixed- $\epsilon$ -greedy method in [8] and a decreasing- $\epsilon$  method [12]. In order to demonstrate the advantage of the proposed adaptive method, all strategies are implemented with fixed real-time price of  $c_{[r]}^t = \mathcal{L}0.15/\text{W}$  and identical constraints in two extreme cases, i.e., a near-static environment and a highly uncertain environment. One can conclude from the average reward curves in Fig. 4(a) that the proposed adaptive method outperforms the fixed  $\epsilon$ -greedy method in the near-static environment. This is due to the fact that the probability of high-cost exploration in the proposed method is gradually decreased with the agent's increasing knowledge of the environment. Whilst as can be observed from Fig. 4(b), the average accumulative reward of the proposed method achieves an approximately 9% improvement as compared to the decreasing- $\epsilon$  method. This is because the latter method fails to adjust the exploration-exploitation trade-off according to the true learning progress in the highly uncertain



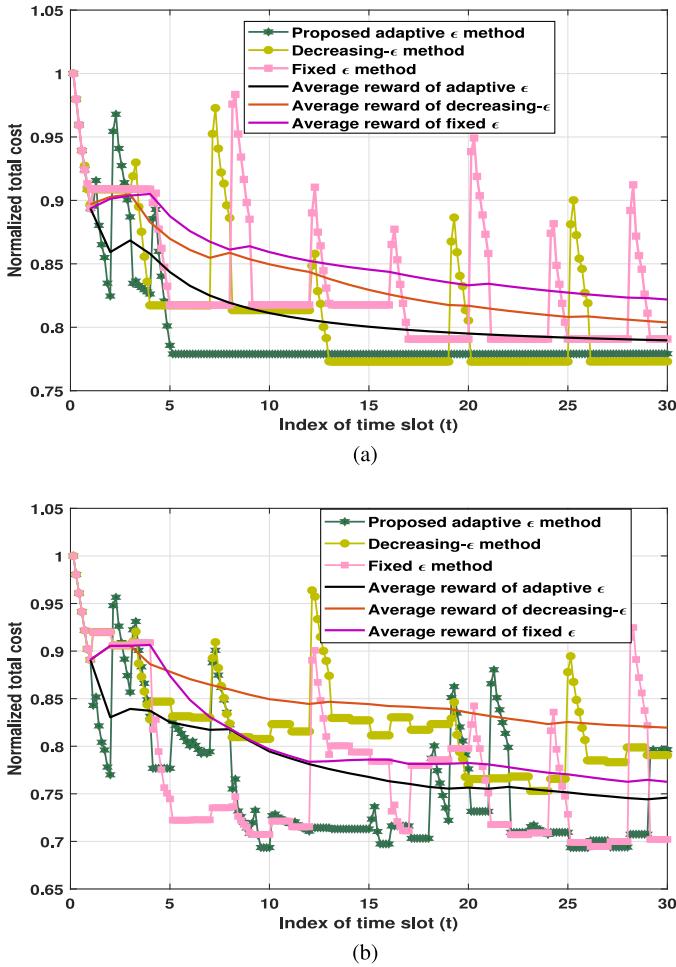


Fig. 4. Comparison of various  $\epsilon$ -greedy methods at  $\gamma = 15$  dB in (a) a near-static environment, (b) a highly uncertain environment.

environment, e.g., intermittent and highly irregular renewable energy generation.

## V. CONCLUSION

The variability of renewable sources introduces large ramps in the energy supply which can lead to increased overall cost as well as grid stability issues. A foresighted energy management strategy has been proposed in this paper for the ISO to minimize the long-term overall

cost of the network in the presence of uncertain renewable energy generation. The proposed strategy accounts for the impacts of the individual A-CPU's on the power network and reacts adaptively with the storage pre-charging amount in a distributed way. Simulation results confirm the effectiveness of the proposed foresighted strategy in achieving a significant performance gain over a recently proposed learning-based storage management design and a baseline myopic design.

## REFERENCES

- [1] D. Li, W.-Y. Chiu, H. Sun, and H. V. Poor, "Multiobjective optimization for demand side management program in smart grid," *IEEE Trans. Ind. Inform.*, vol. 14, no. 4, pp. 1482–1490, Apr. 2018.
- [2] Z. Fan *et al.*, "Smart grid communications: Overview of research challenges, solutions, and standardization activities," *IEEE Commun. Surv. Tut.*, vol. 15, no. 1, pp. 21–38, Firstquarter 2013.
- [3] Z. Zhu *et al.*, "A game theoretic optimization framework for home demand management incorporating local energy resources," *IEEE Trans. Ind. Inform.*, vol. 11, no. 2, pp. 353–362, May 2015.
- [4] W. N. S. F. Wan Ariffin *et al.*, "Sparse beamforming for real-time resource management and energy trading in green C-RAN," *IEEE Trans. Smart Grid*, vol. 8, no. 4, pp. 2022–2031, Jul. 2017.
- [5] J. Xu and R. Zhang, "Cooperative energy trading in CoMP systems powered by smart grids," *IEEE Trans. Veh. Technol.*, vol. 65, no. 4, pp. 2142–2153, Apr. 2016.
- [6] Y. Zhang and M. van der Schaar, "Structure-aware stochastic storage management in smart grids," *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 6, pp. 1098–1110, Dec. 2014.
- [7] X. Wang, Y. Zhang, T. Chen, and G. B. Giannakis, "Dynamic energy management for smart-grid-powered coordinated multipoint systems," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 5, pp. 1348–1359, May 2016.
- [8] X. Zhang, M. R. Nakhai, and W. N. S. F. Wan Ariffin, "Adaptive energy storage management in green wireless networks," *IEEE Signal Process. Lett.*, vol. 24, no. 7, pp. 1044–1048, May 2017.
- [9] Y. Xiao and M. van der Schaar, "Distributed demand side management among foresighted decision makers in power networks," in *Proc. Asilomar Conf. Signals, Syst. Comput.*, Nov. 2013, pp. 1383–1387.
- [10] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. Cambridge MA, USA: MIT Press, 2017.
- [11] M. Tokic, "Adaptive  $\epsilon$ -greedy exploration in reinforcement learning based on value differences," in *KI 2010: Advances in Artificial Intelligence*. Berlin, Germany: Springer, 2010.
- [12] S. Maghsudi *et al.*, "Channel selection for network-assisted D2D communication via no-regret bandit learning with calibrated forecasting," *IEEE Trans. Wireless Commun.*, vol. 14, no. 3, pp. 1309–1322, Mar. 2015.