Aggregate capacity of TCLs with cycling constraints

Austin R. Coffman*,†, Neil Cammardella*, Prabir Barooah*, and Sean Meyn*

Abstract—Thermostatically Controlled Loads (TCLs) such as air conditioners and water heaters typically maintain their temperature within a preset range using on/off actuation. These types of loads are inherently flexible: many different power consumption trajectories exist that can keep the temperature within range. Decades of research has shown that flexible loads can provide valuable grid services.

Quantifying the power and energy capacities of a collection of TCLs is a well-studied problem. However, most works focus on temperature constraints. In this work, we present a characterization of the capacity of a collection of TCLs that considers not only temperature, but also cycling and energy constraints. The characterization leads to a set of convex constraints. A grid operator can use this characterization to compute a feasible power consumption trajectory for an ensemble of TCLs that comes closest to what the operator needs to maintain demand-supply balance. Unlike prior attempts at capacity characterizations incorporating cycling constraints, our results are independent of the algorithm used to coordinate the TCLs.

I. INTRODUCTION

Currently, power balance in power grids is maintained mostly through supply-side actions, i.e., generators are ramped up and down to meet demand, resulting in negative economic and environmental impacts. These negative impacts motivate an active area of research: controlling flexible loads to provide grid support. Some examples of flexible loads that are suitable for grid support are Thermostatically Controlled Loads (TCLs) [1]–[6], HVAC systems in commercial buildings [7], heat pumps [8], and electric pumps for irrigation [9] and pool cleaning [10].

Flexible loads can alter their power consumption without violating Quality of Service (QoS) constraints. A grid operator or balancing authority (BA) can request an ensemble of flexible loads that they consume more or less power with respect to a baseline. Baseline refers to the power consumption that would have occurred without the BA's interference. From the perspective of the BA, an increase (or decrease) of power consumption is identical to the charging (or discharging) of a battery. Due to this similarity, these resources are often termed *Virtual Energy Storage* (VES) [11]. However, VES is cheaper than grid-scale batteries [12].

An extensive literature exists on reference tracking by collections of TCLs [2], [4], [5], [13], [14]. The BA computes

AC and PB are with the Dept. of Mechanical and Aerospace Engineering, University of Florida, Gainesville, FL 32601, USA. NC and SM are with the department of Electrical and Computer Engineering, University of Florida, Gainesville, Fl 32601, USA. The research reported here has been partially supported by the NSF through award 1646229 (CPS-ECCS).

a reference signal (say, in MW) for a collection of loads, and a coordination algorithm has to ensure that the total power consumption deviation of the collection (from the baseline) tracks this reference. However, in order for a BA to design a feasible reference signal, the capacity of the collection of TCLs must be known. Though no formal definition of capacity exists, in this context capacity denotes limits on aggregate power consumption deviation due to QoS constraints at the individual loads. For TCLs, there are at least three QoS constraints: (i) temperature, (ii) cycling rate, and (iii) total energy consumption. If a BA designs a reference signal with an incomplete notion of capacity, the BA must accept poor tracking or the TCL users must accept QoS violations. In both scenarios the long-term outlook is grim: either the BA views TCLs as an unreliable resource, or the TCL users view the BA as an authoritative monarch with unrealistic expectations.

A significant amount of research has focused on characterizing demand flexibility capacity of TCLs [15]–[20]. The impact of enforcing QoS constraints on the loads on reference tracking performance, when the reference is planned without considering capacity, has also been investigated [13], [21]. Early work in this area only accounted for temperature constraints [15], [16]. More recent work has included cycling constraints [17]–[19], [22], but their capacity characterization is limited to specific coordination algorithms, and furthermore not suitable for planning a feasible reference.

In this work we develop a capacity characterization that accounts for three QoS constraints at the individual TCLs: temperature, cycling, and energy. The key novelty of our characterization is its ability to account for cycling constraints. The characterization is independent of the algorithm used to control the ensemble of TCLs. Two, the capacity characterization can be used by a BA to compute the reference for an ensemble of TCLs by solving an optimization problem that is always *feasible* and *convex*. Together, these features ensure that the reference signal so planned can be tracked with any well-designed coordination algorithm that respects the three QoS constraints of each TCL.

The effectiveness of our capacity characterization is demonstrated in simulation experiments by comparing reference tracking performance with two distinct references: one planned with our method and another planned with a method that is representative of prior work that does not incorporate cycling constraints. The capacity characterization we develop requires making an approximate homogeneity assumption, but the numerical results show the method is robust to those assumptions.

^{*} University of Florida

[†] corresponding author, email: bubbaroney@ufl.edu.

Our work starts with the paradigm introduced in [15] of constructing an ensemble battery model. The discrete nature of a TCL's power consumption was ignored in [15] to develop an average temperature like quantity for the collection. Extending this framework to incorporate cycling constraints is challenging due to the on-off nature of TCLs control that leads to an integer valued constraints on number of cycles. We address this challenge by transferring the constraints into quantities that involve fraction of loads at various states, thereby eliminating integer-valued constraints.

This work extends our recent work [20] in the following two ways. The reference planning problem developed here is convex and independent of the coordination algorithm used, whereas the one in [20] is non-convex and only applicable to a specific coordination algorithm. A preliminary version of this paper is published in [23]. The results in [23] is for strictly homogeneous loads, which is extended to a class of heterogeneous loads here. A more extensive numerical investigation is included here to demonstrate the effectiveness of the method, especially with heterogeneous loads.

The paper proceeds as follows: Section II contains descriptions of individual TCL behavior, Section III contains descriptions of aggregate TCL behavior, and Section IV contains the derived aggregate capacity constraints. In Section V, the proposed reference planning method is described. Lastly, Section VI reports the results of numerical experiments.

II. THE INDIVIDUAL TCL

An on/off TCL is any device that turns on or off to maintain a temperature within a preset temperature deadband. Time is discrete, with a sampling period T_s , and is denoted by the index k. There are N TCLs, indexed by $j=1,\ldots,N$. The temperature of the j-th TCL at discrete time instant k is denoted by θ_k^j and its on/off status during the continuous time interval $[kT_s,(k+1)T_s)$ is denoted by m_k^j (=1 if on and 0 if off). We denote the rated electrical power consumption of the j-th TCL, the power consumed by it when on, by the constant P^j .

A. Modeling a TCL's temperature

As in much of prior work [4], [13] temporal evolution of the temperature θ_k^j of the *j*-th TCL is modeled in discrete time as a linear difference equation

$$\theta_{k+1}^{j} = a^{j} \theta_{k}^{j} + (1 - a^{j}) \left(\theta_{k}^{a} - R_{th}^{j} m_{k}^{j} \eta^{j} P^{j} \right)$$
 (1)

with $a^j \triangleq \exp\left(\frac{-T_s}{R_{th}^j C_{th}^j}\right)$, where R_{th}^j and C_{th}^j represent the thermal resistance to ambient temperature θ_k^a and thermal capacitance, respectively. For an air conditioner (AC) providing cooling, the term $\eta^j P^j$ is the thermal power rejected to the ambient by the TCL j when it is on, and η^j is its Coefficient of Performance (COP).

For later use we now define the analytical baseline demand of the j-th TCL, \bar{P}_k^j : it is the electrical power demand needed to maintain θ_k^j at the setpoint, θ_{Set}^j for all k. Because eventually we are interested in aggregate quantities over the

whole collection, it is common to ignore the binary nature of power consumption at this stage; see, e.g., [15]. The qualifier "analytical" is used to emphasize that this is a quantity introduced for analysis: such a demand cannot be observed for a single TCL. The analytical baseline demand can be computed by finding the value of P_k^j in (1) that ensures an equilibrium of (1), with $\theta_{k+1}^j = \theta_k^j = \theta_{\text{set}}^j$ for all k. It follows from straightforward calculations that

$$P_k^{j,b} = \frac{\theta_k^a - \theta_{\text{set}}^j}{\eta^j R_{th}^j}.$$
 (2)

The analytical baseline is a time varying quantity since the ambient temperature θ_k^a is time-varying.

B. QoS constraints for a TCL

The quality of service constraints (QoS) for the j^{th} TCL are:

QoS 1:
$$\left| \theta_k^j - \theta_{\text{set}}^j \right| \le \delta^j, \quad \forall \ k,$$
 (3)

QoS 2:
$$\sum_{i=0}^{\tau_{tcl}^{j}-1} \left| m_{k-i}^{j} - m_{k-1-i}^{j} \right| \le 1, \quad \forall \ k, \quad (4)$$

QoS 3:
$$T_s \left| \sum_{k=0}^{H_b} \left(m_k^j P^j - \hat{P}_k^j \right) \right| \le \tilde{E}^j$$
. (5)

The first constraint says that TCL j's temperature must be kept within $\pm \delta^j$ of the setpoint θ_{set} , where δ^j is a predetermined constant. For later reference, we note that the full width temperature deadband is denoted as $\Delta \triangleq 2\delta$. The second is the cycling constraint; it says that the device can only flip - from either on to off or from off to on - once within a specified period τ_{tel}^{j} . The third is a constraint on the energy consumed over the billing horizon H_b^j : it says the total energy consumed by the TCL over a horizon H_h^j cannot deviate from its (analytical) baseline by more than a specified amount, \tilde{E}^{j} (> 0). Just like the temperature deadband, the parameters H_h^j, \tilde{E}^j are design choices that depend on the j-th consumer's preference. For instance, if the consumer wishes that the energy use over 30 days do not vary by more than 10% of a baseline energy use of 1000 kWh, then $H_b^j = \frac{60}{5} \times 24 \times 30 = 8640$ (for a 5-minute sampling period) and $\tilde{E}^j = 100$ kWh.

The set of TCL-specific parameters that appear in (1),(3)-(5) is $Q_s^j \triangleq \{\theta_{\text{Set}}, \delta, \tau_{tcl}, \tilde{E}, H_b, R_{th}, C_{th}, \eta, a, P\}^j$. A subset of these specifies the QoS constraints of the consumer while the remaining describe mechanical/thermal properties of the hardware.

For later use, we now define variables to describe a TCL's state of flipping from on to off (or vice versa) state, and the state of being stuck in the on (or off) state. The "flip on" or "flip off" variables are defined as

(Flip on)
$$F_{k-1}^{\mathbf{on},j} \triangleq \begin{cases} 1, & \text{if } (m_k^j - m_{k-1}^j) = 1. \\ 0, & \text{otherwise.} \end{cases}$$
 (6)

(Flip off)
$$F_{k-1}^{\text{off},j} \triangleq \begin{cases} 1, & \text{if } (m_{k-1}^j - m_k^j) = 1. \\ 0, & \text{otherwise.} \end{cases}$$
 (7)

We say a TCL is stuck on (respectively, stuck off) at time k if it is off (respectively, on) at that time and has changed mode once in the past τ_{tcl} time instants, so that it is unable to switch mode at the current time k. We define the stuck on and off state as $S_k^{{\rm on},j}$ and $S_k^{{\rm off},j}$:

$$\begin{split} S_k^{\text{on},j} &\triangleq \begin{cases} 1, \text{ if } \sum_{i=0}^{\tau_{tcl}^j-1} \left| m_{k-i}^j - m_{k-1-i}^j \right| = 1, \ m_k^j = 1. \\ 0, \text{ otherwise.} \end{cases} \\ S_k^{\text{off},j} &\triangleq \begin{cases} 1, \text{ if } \sum_{i=0}^{\tau_{tcl}^j-1} \left| m_{k-i}^j - m_{k-1-i}^j \right| = 1, \ m_k^j = 0. \\ 0, \text{ otherwise.} \end{cases} \end{split}$$

III. AGGREGATE QUANTITIES AND ASSUMPTIONS

Section II was devoted to the individual TCL; we now define variables for a collection of N TCLs that are needed to pose the problem precisely. The maximum possible electrical demand of the collection of N TCLs is denoted by P^{agg} and the demand at k is denoted by P_k :

$$P^{\text{agg}} \triangleq \sum_{j=1}^{N} P^{j}, \qquad P_{k} \triangleq \sum_{j=1}^{N} P^{j} m_{k}^{j}.$$
 (8)

The quantity corresponding to P_k during baseline operation is denoted by P_k^b . Recall that the collection of TCLs provide VES service by varying the individual on/off status m_k^j so that the deviation of demand from the baseline tracks a grid supplied reference as closely as possible, without violating any individual's QoS constraints. The grid supplied VES reference is denoted by R_k , which is the desired value of the demand deviation from baseline, denoted by Y_k :

$$Y_k \triangleq P_k - P_k^b. \tag{9}$$

A related quantity that will be useful later is the analytical baseline demand of the aggregate, denoted by P_k^b :

$$\hat{P}_{k}^{b} \triangleq \sum_{j=1}^{N} \hat{P}_{k}^{j,b} = \sum_{j=1}^{N} \frac{\theta_{k}^{a} - \theta_{\text{set}}^{j}}{\eta^{j} R_{th}^{j}}.$$
 (10)

It is the analytical counterpart to P_k^b .

Our development uses the following "fractional" counterparts to aggregate quantities:

$$n_k^{\text{OR}} \triangleq \frac{\sum_{j=1}^N m_k^j}{N},\tag{11}$$

$$f_k^{\text{on}} \triangleq \frac{\sum_{j=1}^N F_k^{\text{on},j}}{N}, \quad f_k^{\text{off}} \triangleq \frac{\sum_{j=1}^N F_k^{\text{off},j}}{N}, \qquad (12)$$

$$s_k^{\text{on}} \triangleq \frac{\sum_{j=1}^N S_k^{\text{on},j}}{N}, \quad s_k^{\text{off}} \triangleq \frac{\sum_{j=1}^N S_k^{\text{off},j}}{N}. \qquad (13)$$

$$s_k^{\text{on}} \triangleq \frac{\sum_{j=1}^N S_k^{\text{on},j}}{N}, \quad s_k^{\text{off}} \triangleq \frac{\sum_{j=1}^N S_k^{\text{off},j}}{N}.$$
 (13)

The quantity f_k^{ON} is called the fraction at time k that decide to flip on at k+1, and s_k^{ON} is called the fraction that is stuck on at k, and similarly for the "off" fractions.

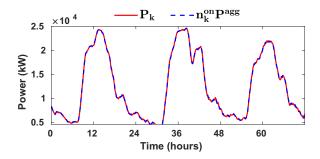


Fig. 1. Comparison of P_k and $n_k^{\rm On}P^{\rm agg}$, for a heterogeneous population of TCLs with thermostat control. Simulation parameters are described in Sec. VI.

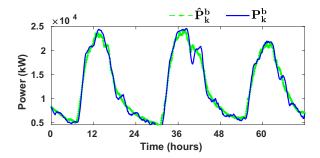


Fig. 2. Comparison of P_k^b and \hat{P}_k^b , for a heterogeneous population of TCLs with thermostat control. Simulation parameters are described in Sec. VI.

A. Role of heterogeneity

We limit ourselves to populations in which the following assumption holds, which we call quasi-heterogeneous populations.

Assumption 1.

$$(i): \quad P_k = n_k^{OR} P^{agg}. \tag{14}$$

(i):
$$P_k = n_k^{on} P^{agg}$$
. (14)
(ii): $\tau_{tcl}^1 = \tau_{tcl}^2 = \dots = \tau_{tcl}^N = \tau_{tcl}$. (15)

Assumption 1(i) means the fraction of loads on at k is equivalent to the total power consumption at that time, and Assumption 1(ii) means the lock-out constraint is the same for all the loads.

The assumption holds for a homogeneous population. Assumption 1(i) holds approximately for a heterogeneous population if the P^{j} 's are drawn from a uniform or Gaussian distribution, or for that matter, any symmetric uni-modal distribution. Results from one numerical experiment are shown in Figure 1; the quantities are nearly identical in this experiment. Details of these simulations are described in Section VI.

The reason for introducing this assumption is that the capacity will be characterized in terms of the fraction on, n_k^{on} , and related quantities introduced above, since they are easy to relate to the cycling constraint of individual TCLs. The Assumption 1(i) allows us to translate the ensemble's power demand P_k to n_k^{On} .

B. Role of the coordination algorithm

For a collection of TCLs to provide VES service, the aggregate power deviation Y_k of the collection has to track a reference R_k . A coordination algorithm is needed to perform this tracking. There are many ways to pose/design a coordination algorithm, see, for example, the references [2], [4], [5], [13], [14]. There are also potentially many metrics to deem a coordination algorithm well designed. While our results do not depend on a specific coordination algorithm, we do specify a requirement for a coordination algorithm to be considered well-designed. The requirement is that when the grid is not asking for VES, i.e., $R_k \equiv 0$, the coordination algorithm should mimic the baseline operation by the thermostat. This is stated formally as the following assumption.

Assumption 2. $P_k^b = \hat{P}_k^b$.

Assumption (2) is used in translating QoS constraints of the individuals into a tractable constraint for the ensemble providing VES, by allowing definitions such as (9) to be independent of the coordination algorithm. It also allows us to use predictions of the ambient temperature to predict the baseline power, which will be useful for reference planning in Section V. Again, the assumption needs to hold only approximately; and the reader can find numerical results justifying Assumption 2 in Figure 2.

IV. AGGREGATE CAPACITY CONSTRAINTS

Aggregate capacity constraints that will be derived in this section refers to constraints on aggregate quantities due to the temperature, cycling, and energy constraints at the individual TCL, i.e., (3)-(5). That is, if each TCL in a collection enforces (3)-(5) then the aggregate constraints will be satisfied. The contrapositive is: if the aggregate constraints are violated, there would exist at least a single TCL that violates its individual QoS constraints. Hence, violating the aggregate constraints means that it is *impossible* for every TCL to satisfy their own QoS; at least one TCL will violate its local constraints.

A. Aggregate scaled temperature

If the TCLs are homogeneous, the dynamics of the average temperature of the ensemble is the same as that of the individual, with aggregate values for the parameters in the model (1), but that is not the case in the heterogeneous case [24]. It is still possible to develop an aggregate model that has a connection to each individual TCL's temperature constraint (3), as done in [15], which we do next.

Consider the aggregate demand deviation from the analytical baseline demand:

$$\hat{Y}_k \triangleq P_k - \hat{P}_k^b. \tag{16}$$

We call \hat{Y}_k the analytical demand deviation, and it is the analytical counterpart of the actual demand deviation Y_k defined in (9). We have the following result.

Lemma 1 (Theorem 5 in [15]). For an arbitrary $\alpha > 0$, denote $\bar{a} = \exp(-T_s/\alpha)$, $\bar{b} = (1 - \bar{a})\alpha$, and define

$$Z_k = \bar{a}Z_{k-1} - \bar{b}\hat{Y}_{k-1} \tag{17}$$

with $Z_0 = 0$. If for all $j \in \{1, ..., N\}$ and for all $k \ge 0$ the constraint (3) is maintained with $\theta_0^j = \theta_{set}$, then $|Z_k| \le \tilde{C}$ for all k, where

$$\tilde{C} \triangleq \sum_{j=1}^{N} \left(1 + \left| 1 - \frac{R_{th}^{j} C_{th}^{j}}{\alpha} \right| \right) \frac{C_{th}^{j} \delta^{j}}{\eta^{j}}.$$
 (18)

Lemma 1 allows us to use the bound (18) as a necessary condition for each TCL to maintain the temperature constraint (3). The original proof of Lemma 1 in [15] is for a continuous time system with constant ambient temperature. A proof for the current setting is given in the Appendix.

Corollary 1. Let the ensemble of TCLs be homogeneous, denote the quantity,

$$g_k = \frac{C_{th}}{\eta} \sum_{i=1}^{N} (\theta_k^j - \theta_{set}), \tag{19}$$

and let $\alpha = C_{th}R_{th}$. Then $g_k = Z_k$ for all k.

That is, in the homogeneous case the quantity Z_k is proportional to the temperature deviation from the setpoint, with the unit of energy (kWh thermal). While it is hard to interpret the quantity Z_k in Lemma 1, it is trying to capture the sum in (19) for a heterogeneous ensemble. We refer to the quantity Z_k in the sequel as the *scaled temperature deviation* of the ensemble.

Comment 1. Lemma 1 holds with \hat{Y}_k replaced with Y_k . This follows immediately from Assumption (1) and the definitions (9) and (16).

B. Fraction of TCLs stuck in on or off mode

The fraction of TCL's stuck on, or off, evolves according to the following inventory model:

$$s_k^{\text{on}} = s_{k-1}^{\text{on}} + f_{k-1}^{\text{on}} - f_{k-1-\tau}^{\text{on}}.$$
 (20)

In words, the fraction that are stuck on, s_{k-1}^{OI} , increases by the fraction that flip on f_{k-1}^{OI} from k-1 to k and decreases by the fraction that had flipped on $k-1-\tau$ time instants in the past. Note that Assumption 1 is used here: if τ^j 's were distinct the equality will not hold. A similar relationship holds for the fraction stuck off:

$$s_k^{\text{off}} = s_{k-1}^{\text{off}} + f_{k-1}^{\text{off}} - f_{k-1-\tau}^{\text{off}}.$$
 (21)

C. Fraction of TCLs on

The fraction of TCLs on, $n_k^{\rm OR}$, is a particulary important quantity since the total electrical power consumption of the ensemble at k is proportional to it due to Assumption 1. We now derive a dynamic model of and constraints on $n_k^{\rm OR}$. This exercise does not have to be repeated for fraction off since that is completely determined by the fraction on.

1) Dynamics: An inventory equation - similar to (20) - couples dynamics of fraction on and fraction that flips:

$$n_k^{\text{on}} = n_{k-1}^{\text{on}} + f_{k-1}^{\text{on}} - f_{k-1}^{\text{off}}.$$
 (22)

In words, the fraction of on devices at time k is the fraction already on at k-1, plus the fraction that flipped on minus the fraction that flipped off from time k-1 to k.

2) Constraints: In the boundary case all N TCLs can be on at time k, which means that no TCLs were previously stuck off. In the case where some TCLs were previously stuck off, an upper bound for the fraction that can be on is $n_k^{\text{ON}} \leq 1 - s_{k-1}^{\text{Off}}$. Similarly, since those TCLs that are stuck on at k-1 must be kept on at k, we have $n_k^{\text{ON}} \geq s_{k-1}^{\text{ON}}$. Thus, we have the following constraint:

$$s_{k-1}^{\text{on}} \le n_k^{\text{on}} \le 1 - s_{k-1}^{\text{off}}.$$
 (23)

D. Aggregate capacity characterization

From Assumption (1), the demand deviation Y_k is related to the fraction of loads on $n_k^{\rm ON}$, which is also related to fraction stuck on/off and fraction flipped through (20)-(21) and (22), respectively. Each of these signals have constraints and some have dynamics, which were derived in previous sections. These are now collected to describe all the constraints on the signal Y_k in order to satisfy TCLs' local QoS. We first "lift" the signal Y_k , for $k=t+1,\ldots,t+H_p$ over a planning horizon H_p to a decision vector $\psi_t^{t+H_p-1}$ that is defined as

$$\psi_t^{t+H_p-1} \triangleq \left[\{ Z_k \}_{t+1}^{t+H_p}, \ \{ Y_k \}_{t+1}^{t+H_p}, \ \{ f_k^{\text{OI}} \}_t^{t+H_p-1}, \dots \right]$$

$$\left\{ f_k^{\text{OII}} \right\}_t^{t+H_p-1}, \ \left\{ s_k^{\text{OII}} \right\}_{t+1}^{t+H_p}, \ \left\{ s_k^{\text{OII}} \right\}_{t+1}^{t+H_p} \right].$$

$$(24)$$

The capacity of the ensemble, the admissible $\{Y_k\}_{k=t+1}^{t+H_p}$ is obtained in terms of the expanded signal ψ_t . Specifically, given a baseline demand \hat{P}_k^b over the same horizon, the capacity of the collection is the set of $\psi_t^{t+H_p-1}$'s that lie in the set $\Omega_t^{t+H_p-1}$, where

$$\Omega_t^{t+H_p-1} \triangleq \left\{ \psi_t^{t+H_p-1} \middle| Z_t = 0, \ s_t^{\text{off}} = 0, \ s_t^{\text{on}} = 0, \ (25) \right. \\
Y_t = 0, \text{ for all } k \in \{t, \dots, t + H_p - 1\}, \\
Z_{k+1} = \bar{a}Z_k - \bar{b}Y_k, \quad |Z_{k+1}| \le \tilde{C}, \quad (26) \\
n_k^{\text{on}} = \frac{1}{P^{\text{agg}}} (Y_k + \hat{P}_k^b), \quad (27)$$

$$s_k^{\text{on}} \le n_{k+1}^{\text{on}} \le 1 - s_k^{\text{off}},\tag{28}$$

$$s_{k+1}^{\text{on}} = s_k^{\text{on}} + f_k^{\text{on}} - f_{k-\tau}^{\text{on}},$$
 (29)

$$s_{k+1}^{\text{off}} = s_k^{\text{off}} + f_k^{\text{off}} - f_{k-\tau}^{\text{off}},\tag{30}$$

$$n_{k+1}^{\text{on}} = n_k^{\text{on}} + f_k^{\text{on}} - f_k^{\text{off}},$$
 (31)

$$n_k^{\text{on}}, s_k^{\text{on}}, s_k^{\text{off}}, f_k^{\text{on}}, f_k^{\text{off}} \in [0, 1],$$
 (32)

$$\sum_{k=t}^{t+H_p} Y_k = 0 \bigg\}. {33}$$

Recall that the constants \tilde{C} , $P^{\rm agg}$ are defined in (18), (8), and the signal \hat{P}_k^b in (10). Eq. (27) uses Assumptions 1 and 2. The last constraint (33) acts as a necessary condition for the QoS constraint (5) for any collection of positive numbers $\{\tilde{E}^j\}$ and any H_p that satisfiew $H_p \leq H_b^j$.

The following result is useful when the constraint set $\Omega_t^{t+H_p-1}$ is used to perform reference planning.

Lemma 2. The set $\Omega_t^{t+H_p-1}$ is convex for every t and $H_p \ge 1$. Suppose that for a given τ and H_p for all t, the following signal

$$\bar{\theta}_k^a \triangleq \rho \theta_k^a, \quad \text{with} \quad \rho \triangleq \sum_{j=1}^N \frac{1}{\eta^j R_{th}^j} \left(\sum_{j=1}^N P^j\right)^{-1}$$
 (34)

satisfies

$$\Theta_{k}^{-}(\tau) + \Gamma < \bar{\theta}_{k+1}^{a} < 1 - \Theta_{k}^{+}(\tau) + \Gamma, \tag{35}$$

for $k \in \{t, ..., t + H_p - 1\}$, where

$$\Theta_k^{-}(\tau) = \sum_{s=k-\tau+1}^{k} \max\{\bar{\theta}_s^a - \bar{\theta}_{s-1}^a, 0\}$$
 (36)

$$\Theta_k^+(\tau) = \sum_{s=k-\tau+1}^k \max\{\bar{\theta}_{s-1}^a - \bar{\theta}_s^a, 0\}, \quad and$$
 (37)

$$\Gamma = \sum_{i=1}^{N} \frac{\theta_{set}^{j}}{\eta^{j} R_{th}^{j}} \left(\sum_{i=1}^{N} P^{j} \right)^{-1}.$$
 (38)

Then the set $\Omega_t^{t+H_p-1}$ is non-empty for every t and $H_p \geq 1$.

Proof. See appendix.
$$\Box$$

The condition on the ambient temperature in Lemma 2 is technical: we have never run into a numerical example (with time varying θ_k^a) where the result of the Lemma does not hold. An example of an ambient temperature trajectory that satisfies this assumption is a constant trajectory. We emphasize that none of the results in this section require the ambient temperature to be constant.

V. REFERENCE PLANNING

Reference planning utilizes the aggregate capacity set from Section IV to plan a reference power deviation trajectory for an ensemble of TCLs to track so that the planned reference is within the TCL's capacity. At time t, this is done by projecting the BA's total desired demand deviation, $\{R_k^{BA}\}_{k=t}^{t+H_p-1}$, onto the aggregate capacity set $\Omega_t^{t+H_p-1}$ to obtain the optimal ψ^* . We need the following definition:

$$(\psi^{BA})_{t}^{t+H_{p}-1} \triangleq \left[\{0\}_{t+1}^{t+H_{p}}, \{R_{k}^{BA}\}_{t+1}^{t+H_{p}}, \{0\}_{t}^{t+H_{p}-1}, \{0\}_{t+1}^{t+H_{p}-1}, \{0\}_{t+1}^{t+H_{p}}, \{0\}_{t+1}^{t+H_{p}} \right].$$
(39)

The *reference planning* problem can be cast as the following convex optimization problem,

$$\psi^* = \arg\min_{\psi} J(\psi) = \|\psi^{BA} - \psi\|_{\Xi}^2$$
s.t. $\psi \in \Omega$ (40)

where sub/super-scripts are omitted from ψ, ψ^* to reduce clutter, Ξ is a symmetric positive definite (s.p.d.) weighting matrix of appropriate dimension, and for $x \in \mathbb{R}^n$, $\|x\|_Q^2 := x^T Q x$ for a s.p.d. $n \times n$ matrix Q.

The component $\{Y_k^*\}$ of ψ^* – see the definition (24) – is denoted by R_k^* in the sequel: it is the "largest" power deviation reference, aligned with the BA's needs, that the TCLs can track without any TCL having to violate its QoS constraints.

The objective function $J(\psi)$ is strictly convex since Ξ is a s.p.d matrix. Combining this with Lemma 2, we have that a solution to the reference planning problem always exists and is unique. In other words, for any ψ^{BA} there will always exist a unique reference signal that a collection of TCLs are ideally suited to track.

1) Information requirement: In order for a BA to solve the reference planning problem (40), it needs to know: (i) the parameters $P^{\rm agg}$, τ , \tilde{C} , \bar{a} , and \bar{b} (ii) the initial conditions $n_t^{\rm OI}$, $f_{t-1}^{\rm OI}$, ..., $f_{t-\tau}^{\rm OIf}$, $f_{t-\tau}^{\rm OIf}$, and (iii) forecasts of the signals θ_k^a , R_k^{BA} over the planning horizon H_p . The ambient temperature forecast can be obtained from weather services and the forecast of R_k^{BA} can be obtained from a prediction of the net load [11]. In the numerical simulations conducted later, we set the initial condition $n_t^{\rm OI}$ to $\hat{P}_k^b/P^{\rm agg}$ (which corresponds to $Y_t=0$ as prescribed in $\Omega_t^{t+H_p-1}$). The initial fraction of loads stuck on/off and the initial scaled temperature deviation are assumed to be zero, as specified in (25). Alternatively, the BA could obtain these quantities through measurements from the population of TCLs.

A. Alternative Method for Reference Planning

To compare with past literature we define a constraint set based on the constraints developed in [15] and the scaled aggregate temperature deviation model (17) for projection of R_k^{BA} . The disadvantage with this constraint set is that the aggregate power and scaled temperature deviation bounds developed in [15] do not account for the individual cycling (4) or energy (5) constraint. This alternative reference planning problem is posed as

$$\min_{\{Y_k\},\{Z_k\}} \sum_{k=t}^{t+H_p-1} (R_k^{BA} - Y_k)^2 \xi + \sum_{k=t+1}^{t+H_p} Z_k^2$$
 (41)

s.t.
$$\forall k \in \{t, ..., t + H_p - 1\}$$

$$Z_{k+1} = \bar{a}Z_k - \bar{b}Y_k, \quad Z_t = 0,$$
 (42)

$$|Z_{k+1}| \le \bar{C}, \quad -\hat{P}_k^b \le Y_k \le P_{\text{agg}} - \hat{P}_k^b,$$
 (43)

where ξ is a constant that specifies the relative importance of goals in the objective. If compared to the bounds developed in Section IV, the bounds for Z_k and Y_k in (43) assume that no TCLs have lock out constraints.

VI. NUMERICAL EXPERIMENTS

We survey here numerical experiments conducted with our proposed reference planning method, and compare the results with those from the alternative method that is representative

TABLE I SIMULATION PARAMETERS

Par.	Unit	value	Par.	Unit	value
N	thousand	60	η^j	N/A	2.5
\bar{C}	MWh	50	θ_k^a	°C	time var.
au	Mins.	20	θ_{set}^{j}	°C	U[21, 21.4]
$ au_{tcl}$	Mins.	10	δ^j	°C	U[0.75, 1]
R_{th}^j	°C/kW	U[2, 2.4]	T_s	Mins.	2
C_{th}^j	kWh/°C	U[2, 2.4]	P^{j}	kW	U[5.6, 7]
\tilde{E}^j	kWh	6.4	P^{agg}	MW	134.4

*U[a,b] represents uniform distribution on [a,b].

of the prior art. The simulated TCLs are residential air conditioner units (ACs). The alternative method is designed to satisfy indoor temperature constraint but *does not* account for cycling constraints of the TCLs. For a full description of the alternative method see [15], where a description of how to use it to plan reference trajectories was given in Section V-A.

Both the proposed method and the alternate method return a reference trajectory for the ensemble. Both methods involve the solution of a convex optimization problem, which is performed using CVX [25].

We also present closed loop simulations. The purpose is to illustrate that the trajectory computed with the proposed method is within the capacity of the TCLs, meaning that the TCLs can collectively track it without any TCL having to violate its OoS constraints. In contrast, we will show that the reference from the alternate method is beyond the capacity of the ensemble; some of the TCLs will have to violate their local QoS constraints in order to collectively track the reference. Alternately, if local OoS is enforced by a local controller, the ensemble will not be able to track the reference. This is demonstrated by performing closed loop simulation with a centralized controller to coordinate the TCLs to track the planned reference signal. The centralized coordinator is a priority stack controller: It is a modified version of the one presented in [15]. While the original one described in [15] enforces the temperature QoS of each TCL, i.e., (3), the modified coordinator presented here also enforces each TCL's cycling QoS (4), but not the energy QoS (5).

The closed loop simulation are performed for three reference tracking scenarios: (t-i) reference computed from the proposed method, (t-ii) reference computed from the alternative method, and (t-iii) reference from the alternative method, but the coordinator *does not* enforce the cycling constraint of the TCLs. We find that only in scenario (t-i) will the ensemble of TCLs be able to track the planned reference while each individual maintaining all three of its QoS constraints. Details are described next.

A. Reference Planning

For both reference planning methods the BA supplied reference, R_k^{BA} , is obtained from BPA, a Balancing Authority in the Pacific Northwest of the United States, and is shown in

TABLE II REFERENCE TRACKING ERRORS

Reference planning method	Tracking Error
Proposed method (Figure 4)	0.06 %
Alternative method (Figure 5)	24 %

Figure 3. A heterogeneous ensemble of loads is considered. The parameters for the loads are based on the values provided in [26] and these values are shown in Table I, along with other simulation parameters. The ambient air temperature is time varying; it is obtained from weatherunderground.com for a summer day in Gainesville, Fl. Each TCL experiences the same ambient temperature.

Figure 3 shows the reference signals planned by the two methods, the proposed method and the alternate one. We plan both references for one day, but only show a portion of the results in Figure 3 for clarity; tracking results in the next section are shown for the full horizon. The reference signal planned with the proposed method is noticeably less aggressive than the reference signal planned with the alternative method. That is, when cycling constraints are not taken into account higher ramp rates are asked from the collection of TCLs to get closer to the BA's requirement. As we will see shortly, this leads to either poor reference tracking, violation of individual TCL's QoS, or both.

B. Closed Loop Reference Tracking

- a) Scenario t-i: The closed loop output P_k is shown in Figure 4 along with the reference planned with the proposed method. The collection of AC units are able to track the planned reference signal with minimal tracking error (see Table II). The individual cycling QoS results are shown in Figure 4 (bottom). Every AC satisfies its cycling QoS: No units cycle faster than $\tau_{tcl}=10$ minutes and the majority of the cycling times concentrate near $\tau=20$ minutes.
- b) Scenario t-ii: The closed loop output Y_k is shown in Figure 4, along with the reference (planned by the alternative method that does incorporate cycling constraints). Since this reference is beyond the capacity of the TCLs, and the coordinator enforces cycling QoS at the individuals, the collection of AC units track the planned reference poorly. For comparison, the reference tracking error reported in Table II is two orders of magnitude higher than the error with our proposed method. This illustrates the need for TCL's cycling constraints to be incorporated in reference planning.
- c) Scenario t-iii: Results are shown in Figure 6: good reference tracking at the cost of excessive cycling. Roughly 20~% of the total switches occurring 2~minutes apart (the sampling time).

VII. SUMMARY AND CONCLUSION

The aggregate capacity characterization proposed here takes into account temperature, cycling, and energy use constraints at each individual TCL. The characterization is in the form of a set of constraints on aggregate quantities.

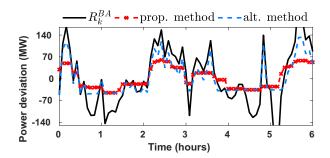
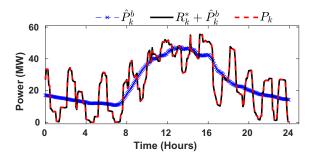


Fig. 3. BA signal (R_k^{BA}) and the reference trajectories (R_k^*) for a collection of 60,000 TCLs.



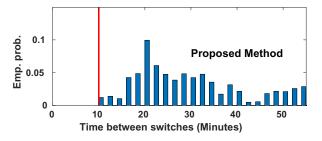
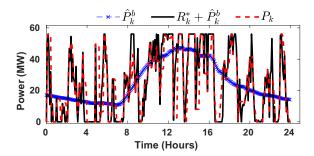


Fig. 4. Closed loop results in scenario t-i: reference planned from the proposed method. (Top): reference tracking results, (Bottom): individual TCL cycling QoS results. The vertical red line indicates τ_{tcl} .

These constraints can be thought of as necessary conditions: if the aggregate state variables for the collection violates these constraints, at least one TCL will have to violate its QoS constraints. Numerical experiments show that the cost of ignoring some of the QoS constraints in the capacity characterization - a feature of prior work - is high: the alternative characterization that does not include cycling constraints leads to tracking errors two orders of magnitude higher than the proposed one.

The information needed to set up the reference planning problem include parameters representing an average TCL such as its COP, allowable temperature bounds, etc. Numerical experiments indicate the results are quite robust to the quasi-homogeneity assumption. It remains to be explored how heterogeneous a collection has to be before the characterization provided is no longer useful.

The reference planning problem we examined here is



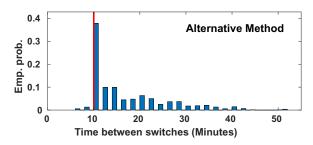
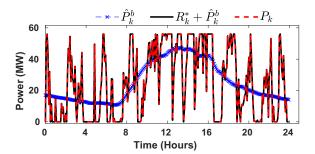


Fig. 5. Closed loop results in scenario t-ii: reference planned from the alternative method. (Top): reference tracking, (Bottom): individual TCL cycling QoS. The vertical red line indicates τ_{tcl} .

a short-term planning problem: its problem data includes prediction of mismatch between demand and supply (in MW). An open problem is capacity characterization of TCLs for long-term planning. An investigation for flexible loads that do not have cycling constraints is provided in [27].

REFERENCES

- [1] D. Callaway and I. Hiskens, "Achieving controllability of electric loads," *Proceedings of the IEEE*, vol. 99, no. 1, pp. 184–199, 2011.
- [2] Y. Chen, M. U. Hashmi, J. Mathias, A. Bušić, and S. Meyn, "Distributed control design for balancing the grid using flexible loads," in *IMA Volume on the Control of Energy Markets and Grids*, 2017, pp. 1–26.
- [3] A. Coffman, A. Bušić, and P. Barooah, "A study of virtual energy storage from thermostatically controlled loads under time-varying weather conditions," in 5th International Conference on High Performance Buildings, July 2018, pp. 1–10.
- [4] J. L. Mathieu, S. Koch, and D. S. Callaway, "State estimation and control of electric loads to manage real-time energy imbalance," *IEEE Transactions on Power Systems*, vol. 28, pp. 430–440, 2013.
- [5] W. Zhang, J. Lian, C.-Y. Chang, and K. Kalsi, "Aggregated modeling and control of air conditioning loads for demand response," *IEEE Transactions on Power Systems*, vol. 28, no. 4, pp. 4655–4664, 2013.
- [6] M. Liu, S. Peeters, D. S. Callaway, and B. J. Claessens, "Trajectory tracking with an aggregation of domestic hot water heaters: Combining model-based and model-free control in a commercial deployment," *IEEE Transactions on Smart Grid*, 2019.
- [7] H. Hao, A. Kowli, Y. Lin, P. Barooah, and S. Meyn, "Ancillary service for the grid via control of commercial building HVAC systems," in *American Control Conference*, June 2013, pp. 467–472.
- [8] Z. E. Lee, Q. Sun, Z. Ma, J. Wang, J. S. MacDonald, and K. Max Zhang, "Providing Grid Services With Heat Pumps: A Review," ASME Journal of Engineering for Sustainable Buildings and Cities, vol. 1, no. 1, 01 2020, 011007.
- [9] A. Aghajanzadeh and P. Therkelsen, "Agricultural demand response for decarbonizing the electricity grid," *Journal of Cleaner Production*, vol. 220, pp. 827 – 835, 2019.



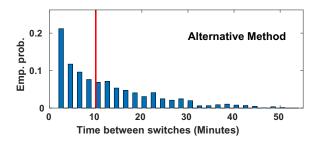


Fig. 6. Closed loop results in scenario t-iii: reference planned from the alternative method and the coordinator TCLs *does not* enforce TCL's cycling constraint. (Top): reference tracking, (Bottom): individual TCL cycling QoS. The vertical red line indicates τ_{tcl} .

- [10] Y. Chen, A. Bušić, and S. Meyn, "State estimation for the individual and the population in mean field control with application to demand dispatch," *IEEE Transactions on Automatic Control*, vol. 62, no. 3, pp. 1138–1149, March 2017.
- [11] P. Barooah, *Smart Grid Control: An Overview and Research Opportunities*. Springer Verlag, 2019, ch. Virtual energy storage from flexible loads: distributed control with QoS constraints, pp. 99–115.
- [12] N. J. Cammardella, R. W. Moye, Y. Chen, and S. P. Meyn, "An energy storage cost comparison: Li-ion batteries vs Distributed load control," in 2018 Clemson University Power Systems Conference (PSC), Sep. 2018, pp. 1–6.
- [13] A. Coffman, A. Bušić, and P. Barooah, "Virtual energy storage from TCLs using QoS preserving local randomized control," in 5th ACM International Conference on Systems for Built Environments (BuildSys), November 2018, p. 10.
- [14] M. Liu and Y. Shi, "Model predictive control of aggregated heterogeneous second-order thermostatically controlled loads for ancillary services," *IEEE Transactions on Power Systems*, vol. 31, no. 3, pp. 1963–1971, May 2016.
- [15] H. Hao, B. M. Sanandaji, K. Poolla, and T. L. Vincent, "Aggregate flexibility of thermostatically controlled loads," *IEEE Transactions on Power Systems*, vol. 30, no. 1, pp. 189–198, Jan 2015.
- [16] L. Zhao, W. Zhang, H. Hao, and K. Kalsi, "A geometric approach to aggregate flexibility modeling of thermostatically controlled loads," *IEEE Transactions on Power Systems*, vol. 32, no. 6, pp. 4721–4731, Nov 2017.
- [17] C. Ziras, S. You, H. W. Bindner, and E. Vrettos, "A new method for handling lockout constraints on controlled TCL aggregations," in 2018 Power Systems Computation Conference (PSCC), June 2018, pp. 1–7.
- [18] B. M. Sanandaji, T. L. Vincent, and K. Poolla, "Ramping rate flexibility of residential HVAC loads," *IEEE Transactions on Sustainable Energy*, vol. 7, no. 2, pp. 865–874, April 2016.
- [19] D. Cheng, W. Zhang, and K. Wang, "Hierarchical reserve allocation with air conditioning loads considering lock time using Benders decomposition," *International Journal of Electrical Power & Energy* Systems, vol. 110, pp. 293 – 308, 2019.
- [20] A. R. Coffman, A. Bušić, and P. Barooah, "Aggregate capacity for TCLs providing virtual energy storage with cycling constraints," in IEEE Conference on Decision and Control, December 2019.
- [21] Y. Chen, "Markovian demand dispatch design for virtual energy

storage to support renewable energy integration," Ph.D. dissertation, Ph. D. dissertation, University of Florida, Gainesville, FL, USA, 2016.

- [22] P. Wang, D. Wu, and K. Kalsi, "Flexibility estimation and control of thermostatically controlled loads with lock time for regulation service," *IEEE Transactions on Smart Grid*, pp. 1–1, 2020.
- [23] A. R. Coffman, N. Cammardella, P. Barooah, and S. Meyn, "Flexibility capacity of thermostatically controlled loads with cycling/lock-out constraints," in 2020 American Control Conference (ACC), July 2020, pp. 527–532.
- [24] Z. Guo, A. R. Coffman, J. Munk, P. Im, T. Kuruganti, and P. Barooah, "Aggregation and data driven identification of building thermal dynamic model and unmeasured disturbance," *Energy and Buildings*, September 2020, in press, available online Sept 2020.
- [25] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 1.21," http://cvxr.com/cvx, Feb. 2011.
- [26] J. Mathieu, M. Dyson, and D. Callaway, "Using residential electric loads for fast demand response: The potential resource and revenues, the costs, and policy recommendations," in 2012 ACEEE Summer Study on Energy Efficiency in Buildings, 2012.
- [27] A. R. Coffman, Z. Guo, and P. Barooah, "Capacity of flexible loads for grid support: statistical characterization for long term planning," in 2020 American Control Conference (ACC), July 2020, pp. 533–538.

APPENDIX

A. Proof of Lemma 1

The proof roughly follows the one found in [15], with slight modification to handle time varying weather and the discrete time dynamics. By construction, the discrete time dynamics (17) is the discrete time equivalent of the ode

$$\dot{z}(t) = -\alpha z(t) - \tilde{y}(t) \tag{44}$$

with zero-order-hold, where $Z_k=Z(t_k)$ and $\tilde{Y}_k=\tilde{Y}(t_k)$. A similar observation is true for the recursion (1). That is, the discrete time dynamics (1) is the discrete-time equivalent of the ode

$$\dot{\theta}(t) = \frac{1}{RC} \left(\theta^a(t) - \theta(t) \right) + \frac{\eta}{C} P(t), \tag{45}$$

with zero-order hold, where $\theta_k = \theta(t_k)$, $\theta_k^a = \theta^a(t_k)$, and $P^j m_k^j = P(t_k)$. Now if we define,

$$Z^{j}(t) := \frac{C_{th}^{j}}{\eta^{j}} (\theta^{j}(t) - \theta_{\text{set}}^{j})$$
 (46)

then this quantity evolves as,

$$\dot{Z}^{j}(t) = -a^{j}Z^{j}(t) - \tilde{P}^{j}(t), \quad \tilde{P}^{j}(t) = P(t) - \hat{P}^{j,b}(t), \tag{47}$$

where $a^j = R^j_{th} C^j_{th}$, and

$$\hat{P}^{j,b}(t) = \frac{\theta^a(t) - \theta_{\text{set}}^j}{\eta^j R_{tb}^j}.$$
(48)

Now taking the Laplace transform of both the continuous time odes we have,

$$Z(s) = -\frac{1}{s+\alpha}\tilde{Y}(s)$$
, and $Z^{j}(s) = -\frac{1}{s+a^{j}}\tilde{P}^{j}(s)$. (49)

Where we have assumed Z(0)=0 and used $\theta^j(0)=\theta_{\mbox{set}},$ so that $Z^j(0)=0.$

Now, by their respective definitions we have that $\tilde{Y}(s) = \sum_{j=1}^{N} \tilde{P}^{j}(s)$ so that,

$$Z(s) = \sum_{j=1}^{N} -\frac{1}{s+\alpha} \tilde{P}^{j}(s),$$
 (50)

$$= \sum_{j=1}^{N} \frac{s+a^{j}}{s+\alpha} \frac{-1}{s+a^{j}} \tilde{P}^{j}(s)$$
 (51)

$$= \sum_{i=1}^{N} \frac{s+a^{j}}{s+\alpha} Z^{j}(s).$$
 (52)

Now taking the inverse Laplace transform of the equation (52) and applying the bound $\|y(t)\|_{\infty} \le \|h(t)\|_1\|u(t)\|_{\infty}$ for the inverse transforms of the relation Y(s) = H(s)U(s) we have,

$$||Z(t)||_{\infty} \le \sum_{j=1}^{N} \left(1 + \left| 1 - \frac{R_{th}^{j} C_{th}^{j}}{\alpha} \right| \right) ||Z^{j}(t)||_{\infty}.$$
 (53)

Since the above is valid for any $t \in \mathbb{R}$, we evaluate it at the point t_k to get,

$$||Z_k||_{\infty} \le \sum_{j=1}^N \left(1 + \left|1 - \frac{R_{th}^j C_{th}^j}{\alpha}\right|\right) ||Z_k^j||_{\infty},$$
 (54)

which is valid for any sequence of times $\{t_k\}_k$ that satisfy $t_k = t_{k-1} + T_s$ with $t_0 = 0$. Now, by assumption in the Lemma the quantity $\|Z_b^j\|_{\infty} < \bar{C}^j$ so that

$$|Z_k| \le \sum_{j=1}^N \left(1 + \left| 1 - \frac{R_{th}^j C_{th}^j}{\alpha} \right| \right) \bar{C}^j, \quad \forall \ k,$$
 (55)

which is the desired result. \square

Note that Lemma 1 here appears in [15] in continuous time. In our proof we use a connection to continuous time, and the fact that our recursion (17) is an exact discretization of a certain ode. Additionally, in [15] the result in Lemma 1 is done for a time invariant ambient temperature. As we see from the proof here, the ambient temperature can be time varying and this will not effect the result.

B. Proof of Lemma 2

To show convexity, we use the fact that the intersection of a finite number of convex sets is convex. Each constraint in $\Omega_t^{t+H_p-1}$ is convex as the inequality constraints are convex sets and the equality constraints are affine. Thus, $\Omega_t^{t+H_p-1}$ is convex as it is the finite intersection of convex sets.

To show feasibility consider the baseline scenario. In this scenario $Y_k \equiv 0$, which together with the initial condition $Z_t = 0$ produces $Z_k \equiv 0$. Hence constraints (26) and (33) are satisfied. From the constraint (27) we have that n_k^{On} will equal

$$n_k^{\text{on}} = \bar{\theta}_k^a - \sum_{j=1}^N \frac{\theta_{\text{set}}^j}{\eta^j R_{th}^j} \left(\sum_{j=1}^N P^j\right)^{-1} = \bar{\theta}_k^a - \Gamma,$$
 (56)

and the difference satisfies $n_k^{\text{OR}} - n_{k-1}^{\text{OR}} = \bar{\theta}_k^a - \bar{\theta}_{k-1}^a$. The constraint (31) is satisfied by

$$f_k^{\text{off}} = \max\{\bar{\theta}_{k-1}^a - \bar{\theta}_k^a, 0\}, \quad \text{and}$$
 (57)

$$f_k^{\text{on}} = \max\{\bar{\theta}_k^a - \bar{\theta}_{k-1}^a, 0\},$$
 (58)

by definition. Upon substituting these choices in the constraints (29) and (30) and using the initial conditions, we have

$$s_k^{\text{on}} = \sum_{s=k-\tau+1}^k \max\{\bar{\theta}_s^a - \bar{\theta}_{s-1}^a, 0\} = \Theta_k^-(\tau)$$
 (59)

$$s_k^{\text{off}} = \sum_{s=k-\tau+1}^k \max\{\bar{\theta}_{s-1}^a - \bar{\theta}_s^a, 0\} = \Theta_k^+(\tau)$$

so that by hypothesis, we have

$$s_k^{\text{on}} + \Gamma \le \bar{\theta}_{k+1}^a \le 1 - s_k^{\text{off}} + \Gamma, \tag{60}$$

which implies that

$$s_k^{\text{on}} \le n_{k+1}^{\text{on}} \le 1 - s_k^{\text{off}},\tag{61}$$

and hence the constraint (28) is satisfied. Additionally, by construction n_k^{OI} satisfies (32) and since f_k^{OI} and f_k^{OI} are the positive difference of successive values of n_k^{OI} , they too will satisfy (32). By construction s_k^{OI} and s_k^{OI} are non-negative. Further from the constraint (28) holding we have $s_k^{\text{OI}} \leq 1$. Since the fraction of loads stuck on and off satisfy $s_k^{\text{OI}} + s_k^{\text{OI}} \leq 1$ we have that $s_k^{\text{OI}} \leq 1$. Hence, both s_k^{OI} and s_k^{OII} satisfy (32).

The above argument, for all of the above constraints, works for any starting index t and any positive planning horizon H_p . \Box

C. VES constraint

The BA requires the constraint (33) to ensure that the collection of TCLs do not act as generators. We repeat this constraint here for t=0:

$$\sum_{k=0}^{H_p} Y_k = 0. {(62)}$$

We now show that this constraint is a necessary condition for the individual TCLs energy constraint (5). We assume that $H_p = H_b$, which loses no generality as H_p is arbitrary and would already be a function of H_b . Summing (5) over the j index and expanding the absolute value,

$$-\sum_{j=1}^{N} \tilde{E}^{j} \le T_{s} \sum_{j=1}^{N} \sum_{k=0}^{H_{p}} (m_{k}^{j} P - \hat{P}_{k}^{j}) \le \sum_{j=1}^{N} \tilde{E}^{j}.$$
 (63)

$$\implies -\sum_{j=1}^{N} \tilde{E}^{j} \leq T_{s} \sum_{k=0}^{H_{p}} \sum_{j=1}^{N} (m_{k}^{j} P - \hat{P}_{k}^{j}) \leq \sum_{j=1}^{N} \tilde{E}^{j},$$

$$\implies -\sum_{j=1}^{N} \tilde{E}^{j} \leq T_{s} \sum_{k=0}^{H_{p}} Y_{k} \leq \sum_{j=1}^{N} \tilde{E}^{j}.$$
(64)

Converting back to absolute value, the aggregated version of (5) is

$$T_s \left| \sum_{k=0}^{H_p} Y_k \right| \le \sum_{j=1}^N \tilde{E}^j, \tag{65}$$

which due to (62) will be true for all values of \tilde{E}^j , as the RHS term in (65) is defined to be greater than or equal to zero. If (65) is not satisfied, then it can be shown through the law of the contrapositive that there would exist at least a single TCL that does not satisfy (5). In the scenario that the individual TCLs do not have symmetric energy constraints, then the aggregate version of (5) would resemble (64); The constraint (62) still enforces this.