Aggregate Flexibility Capacity of TCLs With Cycling Constraints

Austin R. Coffman , Neil Cammardella, Prabir Barooah , Senior Member, IEEE, and Sean Meyn , Fellow, IEEE

Abstract—Thermostatically Controlled Loads (TCLs), e.g., an air conditioner, 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. Grid operators need tools for quantifying demand flexibility of a collection of TCLs in order to use them as a resource. However, computationally tractable characterization of flexibility capacity obtained so far has considered temperature constraints alone, ignoring their cycling/lock-out constraints: the length of time a TCL must stay in "on" stage before switching to "off," or vice versa. In this work, we present a tractable characterization of the flexibility capacity of a collection of TCLs that incorporates not only temperature but also cycling and total energy consumption constraints. Unlike prior attempts at capacity characterizations incorporating cycling constraints, our results are independent of the algorithm used to coordinate the TCLs. The characterization leads to a set of convex constraints. A grid operator can use this characterization to compute a power consumption trajectory for an ensemble of TCLs that comes closest to what the operator needs to maintain demand-supply balance. Numerical results are provided to showcase the effectiveness of the proposed characterization.

Index Terms—Aggregates, batteries, load modeling, mathematical model.

I. INTRODUCTION

ANY electric loads have flexibility in their electricity demand since distinct power demand profiles can provide the same level of consumer satisfaction. With appropriate control algorithms, flexible loads can vary their demand over their baseline demand - within a limit - without adversely affecting consumers' quality of service (QoS). Baseline demand refers to the power consumption that would have occurred without such intervention. From the point of view of the grid, such intentional

Manuscript received 16 April 2021; revised 11 August 2021, 3 November 2021, and 1 February 2022; accepted 12 February 2022. Date of publication 22 March 2022; date of current version 22 December 2022. This work was supported by NSF under Awards 1646229 (CPS) and 2122313 (ECCS). The work of Austin R. Coffman was supported by the University of Florida Graduate School Preeminence Award. Paper no. TPWRS-00618-2021. (Corresponding author: Austin R. Coffman.)

Austin R. Coffman and Prabir Barooah are with the Department of Mechanical and Aerospace Engineering, University of Florida, Gainesville, FL 32601 USA (e-mail: bubbaroney@ufl.edu; pbarooah@ufl.edu).

Neil Cammardella and Sean Meyn are with the Department of Electrical and Computer Engineering, University of Florida, Gainesville, Fl 32601 USA (e-mail: neilcammardella@gmail.com; meyn@ece.ufl.edu).

Color versions of one or more figures in this article are available at https://doi.org/10.1109/TPWRS.2022.3160071.

Digital Object Identifier 10.1109/TPWRS.2022.3160071

variation appears like charging and discharging of a battery. Such demand deviation is therefore called Virtual Energy Storage (VES) [1]. VES can be often less expensive than batteries [2].

Thermostatically controlled loads (TCLs), in particular, have been the subject of intense study as a valuable source of VES. An extensive literature exists on designing algorithms to coordinate a collection of TCLs so that their aggregate demand deviation tracks a reference for demand deviation from the baseline demand [3]–[7]. In this framework, it is assumed that a reference (in kW or MW) for demand deviation of the collection is available in real-time or with look-ahead to the coordination algorithm. Such a reference can be provided by the BA or a load aggregator. A control algorithm has to then coordinate actions of the TCLs so that the deviation of the aggregate demand from the baseline tracks the reference.

This paper is *not* on design of coordination algorithms, but on a related and equally important topic: that of the flexibility capacity of a collection. One can define the flexibility capacity of a collection of loads as the set of feasible reference signals for the collection. A reference signal is feasible if it is possible for the collection to track it while simultaneously each TCL is able to maintain its local QoS constraints. The geometry of such sets is complex, so approximations are sought [8].

For TCLs, and particularly air conditioners, there are at least three QoS constraints: (i) temperature, (ii) cycling rate, and (iii) total energy consumption. The cycling constraint is sometimes also referred to as the "lock-out constraint": once a TCL switches from "on" to "off," or vice versa, it needs to stay in that state for some time before switching back to prevent damage or performance degradation [3]. Finally, the total energy consumed by each TCL should not increase in providing grid service. If a reference signal is designed with an incomplete notion of capacity, the BA or load aggregator must accept poor tracking or the TCL users must accept QoS violations. In both scenarios the long-term outlook is grim: either the BA or load aggregator views TCLs as an unreliable resource, or the TCL users view the BA or load aggregator as an authoritative monarch with unrealistic expectations.

Some works on flexibility capacity characterization of TCL collections only account for temperature constraints [8]–[10]. Some works have included cycling constraints [11]–[15], but their characterizations are limited to specific coordination algorithms, or are not suitable for computing a feasible reference for the collection. Another body of work falls in between: it examines the impact of enforcing cycling constraints at the loads on the reference tracking performance of the collection [6], [16].

0885-8950 © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information. In this paper we present a flexibility capacity characterization - as a set of convex constraints on aggregate quantities - that accounts for all three QoS constraints at the individual TCLs: indoor temperature, compressor cycling, and energy consumption. These constraints act as a set of necessary conditions: if the aggregate quantities, such as the aggregate power consumption, do not satisfy these constraints then it is impossible to simultaneously realize the two goals: the collection tracks the aggregate power consumption and every TCL satisfies its local QoS.

This work makes two key contributions over prior art: the capacity characterization provided here not only incorporates cycling constraints at the individual but it is also independent of the algorithm used to coordinate the TCLs. Note that the prior works that do include cycling constraints provide results that are valid only for specific coordination algorithms. In addition, the proposed characterization can be used by a BA or a load aggregator to compute an "optimal reference" for an ensemble of TCLs - one that satisfies aggregate necessary conditions for the collection and comes closest to the grid's needs - by solving an optimization problem that is always feasible and convex.

Our work develops necessary conditions on aggregate quantities for individual TCLs to satisfy local constraints. Sufficient conditions solely in terms of aggregate quantities cannot be developed since no matter how reasonable an aggregate reference trajectory is, a poorly designed coordination algorithm can cause at least one TCL to violate at least one of its local constraints. If the coordination algorithm is also a part of the reference design problem, sufficient conditions are at-least possible since the coordination prescribes how each TCLs state will evolve, e.g., the work [17]. However, references designed in this way are tailored to the specific coordination algorithm and may not generalize to other algorithms. In contrast, the work here utilizes numerical simulations with a simple centralized coordination algorithm to show references planned with our method (that relies on aggregate necessary conditions and no dependence on the coordination algorithm) are tracked by TCLs without any TCL violating any of its local constraints.

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 without considering 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.

For ease of exposition of the numerical results, we assume the reference is computed by some central entity - such as the BA or a load aggregator - and provided to the coordination algorithm. The main contribution of the paper - the set of constraints that determine the capacity of a collection - is not dependent on the regulatory or economic framework that enables loads to participate in grid support services.

A preliminary version of this paper is published in [18]. The results in [18] are for strictly homogeneous loads, which is extended to a class of heterogeneous loads here. Furthermore, this paper provides a more streamlined development of the capacity characterization, and presents a more extensive numerical investigation 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 [3], [6] temporal evolution of the temperature θ_k^j of the j-th TCL is modeled in discrete time as a linear difference equation

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

with $a^j \triangleq \exp(\frac{-T_g}{R_{th}^j C_{th}^j})$, where R_{th}^j and C_{th}^j represent the thermal resistance to ambient temperature $\theta_k^{a,j}$ 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, $\hat{P}_k^{j,b}$: it is the electrical power demand needed to maintain θ_k^j at the setpoint, $\theta_{\rm set}^j$ for all k (superscript b to denote baseline). 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., [8]. 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^j in (1) that ensures an equilibrium of (1), with $\theta_{k+1}^j = \theta_k^j = \theta_{\rm set}^j$ for all k. It follows from straightforward calculations that the analytical baseline demand is

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

The analytical baseline is a time varying quantity since the ambient temperature $\theta_k^{a,j}$ is time-varying.

B. QoS Constraints for a TCL

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

QoS1:
$$\left|\theta_k^j - \theta_{\text{Set}}^j\right| \le \delta^j, \quad \forall k,$$
 (3)

QoS2:
$$\sum_{i=0}^{\tau_{tel}^{j}-1} \left| m_{k-i}^{j} - m_{k-1-i}^{j} \right| \le 1, \forall k,$$
 (4)

QoS3:
$$T_s \left| \sum_{k=0}^{H_b^j} \left(m_k^j P^j - \hat{P}_k^{j,b} \right) \right| \le \bar{E}^j$$
. (5)

The first constraint says that TCL j's temperature must be kept within $\pm \delta^{j}$ of the setpoint θ_{set}^{j} , 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_b^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 , \bar{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 $\tilde{E}^{j} = 100 \text{ kWh}$.

The set of TCL-specific parameters that appear in (1), (3)–(5) is $Q_s^j \triangleq \{\theta_{Set}, \delta, \tau_{tel}, \tilde{E}, H_b, R_{th}, C_{th}, \eta, a, P\}^j$. 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)
$$F_{k-1}^{\text{on},j} \triangleq \begin{cases} 1, & \text{if } (m_k^j - m_{k-1}^j) = 1. \\ 0, & \text{otherwise.} \end{cases}$$
 (6)

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

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

$$S_k^{\mathrm{cm},j} \triangleq \begin{cases} 1, \text{ if } \sum_{t=0}^{\tau_{tet}^j-1} \left| m_{k-t}^j - m_{k-1-t}^j \right| = 1, \ m_k^j = 1. \\ 0, \text{ otherwise.} \end{cases}$$

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

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 P_k^b denotes the power consumption during baseline operation, i.e., when the population of TCLs make on/off decisions according to their thermostat controller. Further detail on this quantity is given in Section III-C. 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 :

$$Y_k \triangleq P_k - P_k^b$$
. (9)

A related quantity that will be useful later is the analytical baseline demand of the aggregate, denoted by \hat{P}_k^b :

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

It is the analytical counterpart to P_k^b .

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

$$n_k^{\text{on}} \triangleq \frac{\sum_{j=1}^{N} m_k^j}{N}, \quad (11)$$

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

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

A. Role of Heterogeneity

We limit ourselves to populations in which the following assumption holds, which we call *quasi-homogeneous* populations. Assumption 1:

(i):
$$P_k = n_k^{on} P^{agg}$$
. (14)

(ii):
$$\tau_{tcl}^1 = \tau_{tcl}^2 = \dots = \tau_{tcl}^N \stackrel{=}{\nabla} \tau_{tcl}$$
. (15)

(iii):
$$H_b^1 = H_b^2 = \cdots = H_b^N \stackrel{=}{\nabla} H_b$$
 (16)

All other quantities are allowed to vary across each TCL. Assumption 1(i) means the fraction of loads on at k is equivalent to the total power consumption at that time, Assumption 1(ii) means the lock-out constraint is the same for all the loads, and Assumption 1(iii) means the energy billing horizon is the same for all the devices.

We point out that the main motivation for the quasihomogeneous assumption is tractability. For example, without

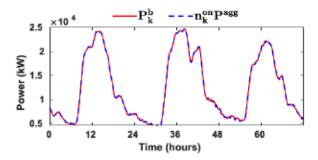


Fig. 1. Comparison of P_k^b and $n_k^{on}P^{agg}$, for a heterogeneous population of TCLs with thermostat control. Simulation parameters are described in Section VI.

the assumption that the lock-out time is uniform over all TCLs, deriving evolution equations for the various fractions become intractable.

Assumption 1(iii) is likely to hold, as consumers billing periods typically occur on a monthly basis. In any case, the heterogeneity between consumers can be accounted for by the \bar{E}^{j} term.

The assumption holds for a homogeneous population. Assumption 1(i) replaces each P^j with the average in computation of the total power and it 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 Fig. 1; the quantities are nearly identical in this experiment. Details of these simulations are described in Section VI.

The reason for introducing Assumption 1(i) is that the capacity will be characterized in terms of the fraction on, $n_k^{\rm on}$, 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 $n_k^{\rm on}$.

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 [3]-[7]. 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 the coordination algorithm must enforce local QoS constraints. The reasoning for this requirement is our focus on aggregate level conditions for reference planning. These aggregate level conditions are low dimensional translations of each TCLs QoS constraints, so that the BA or load aggregator can easily perform resource allocation of TCLs. In this setting, without a local controller ensuring QoS these aggregate conditions do not have the intended effect.

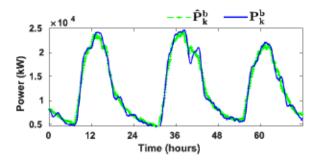


Fig. 2. Comparison of P_k^b and \hat{P}_k^b , for a heterogeneous population of TCLs with thermostat control. Simulation parameters are described in Section VI.

C. Computing the Baseline

In order to design aggregate power deviation trajectories, knowledge of the baseline power consumption, P_k^b , is required. This quantity can be computed by simulating a collection of TCLs under thermostat control. We assume that the BA or load aggregator utilizes the approximate model (10) for computation of the baseline in reference planning, i.e., use \hat{P}_k^b instead of P_k^b ; see (28). The ambient temperature is a strong predictor of baseline power consumption for air conditioners (see [19] and reference within). This introduces a slight plant-model mismatch in the sense that our model for the baseline, \hat{P}_k^b may have some discrepancies from the actual baseline P_k^b . As shown in Fig. 2, the two baseline power trajectories are similar for a typical trajectory of ambient temperature. The ambient temperature utilized to produce the results in Fig. 2 is data collected from Gainesville, FL during a hot summer day.

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 [20]. It is still possible to develop an aggregate model that has a connection to each individual TCL's temperature constraint (3), as done in [8], which we do next.

Consider the aggregate demand deviation from the analytical baseline demand:

$$\hat{Y}_k \triangleq P_k - \hat{P}_k^b$$
. (17)

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 [8]): For an arbitrary $\alpha > 0$, denote $\bar{a} = \exp(-T_s/\alpha)$, $\bar{b} = (1 - \bar{a})\alpha$, and define

$$Z_k = \bar{a} Z_{k-1} - b \hat{Y}_{k-1} \tag{18}$$

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

$$\tilde{C} \triangleq \sum_{j=1}^{N} \left(1 + \left| 1 - \frac{R_{th}^{j} C_{th}^{j}}{\alpha} \right| \right) \frac{C_{th}^{j} \delta^{j}}{\eta^{j}}.$$
 (19)

Lemma 1 allows us to use the bound (19) as a necessary condition for each TCL to maintain the temperature constraint (3). The original proof of Lemma 1 in [8] 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,

$$g_k = \frac{C_{th}}{\eta} \sum_{j=1}^{N} (\theta_k^j - \theta_{set}), \qquad (20)$$

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 (20) for a heterogeneous ensemble. We refer to the quantity Z_k in the sequel as the scaled temperature deviation of the ensemble.

B. Fraction of TCLs Stuck in on or Off Mode

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

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

In words, the fraction that are stuck on, s_{k-1}^{on} , increases by the fraction that flip on f_{k-1}^{on} 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:

$$s_k^{\text{off}} = s_{k-1}^{\text{off}} + f_{k-1}^{\text{off}} - f_{k-1-\tau}^{\text{off}}.$$
 (22)

C. Fraction of TCLs on

The fraction of TCLs on, n_k^{on} , is a particularly 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^{on} . This exercise does not have to be repeated for fraction off since that is completely determined by the fraction on.

 Dynamics: An inventory equation - similar to (21) - couples dynamics of fraction on and fraction that flips:

$$n_k^{\text{on}} = n_{k-1}^{\text{on}} + f_{k-1}^{\text{on}} - f_{k-1}^{\text{off}}.$$
 (23)

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:

$$s_{k-1}^{\text{on}} \le n_k^{\text{on}} \le 1 - s_{k-1}^{\text{off}}.$$
 (24)

D. Aggregate Capacity Characterization

From Assumption (1), the analytical demand deviation \hat{Y}_k is related to the fraction of loads on n_k^{on} , which is also related to fraction stuck on/off and fraction flipped through (21)–(22) and (23), 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 \hat{Y}_k in order to be consistent with TCLs' local QoS. We first "lift" the signal \hat{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

$$\psi_{t}^{t+H_{p}-1} \triangleq \left[\left\{ Z_{k} \right\}_{t+1}^{t+H_{p}}, \left\{ \hat{Y}_{k} \right\}_{t+1}^{t+H_{p}}, \left\{ f_{k}^{\text{on}} \right\}_{t}^{t+H_{p}-1}, \dots \right. \\ \left. \left\{ f_{k}^{\text{off}} \right\}_{t}^{t+H_{p}-1}, \left\{ s_{k}^{\text{on}} \right\}_{t+1}^{t+H_{p}}, \left\{ s_{k}^{\text{off}} \right\}_{t+1}^{t+H_{p}} \right]. \tag{25}$$

The capacity of the ensemble, the admissible $\{\hat{Y}_k\}_{k=t+1}^{t+H_p}$ 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

$$\Omega_t^{t+H_p-1} \triangleq \left\{ \psi_t^{t+H_p-1} \middle| Z_t = 0, \ s_t^{\text{off}} = 0, \ s_t^{\text{on}} = 0, \right.$$
(26)

 $Y_t = 0$, and for all $k \in \{t, ..., t + H_p - 1\}$,

$$Z_{k+1} = \bar{a}Z_k - \bar{b}\hat{Y}_k, |Z_{k+1}| \le \bar{C},$$
 (27)

$$n_k^{\text{on}} = \frac{1}{P^{\text{agg}}} (\hat{Y}_k + \hat{P}_k^b),$$
 (28)

$$s_k^{\text{on}} \le n_{k+1}^{\text{on}} \le 1 - s_k^{\text{off}},$$
 (29)

$$s_{k+1}^{\text{on}} = s_k^{\text{on}} + f_k^{\text{on}} - f_{k-\tau}^{\text{on}},$$
 (30)

$$s_{k+1}^{\text{off}} = s_k^{\text{off}} + f_k^{\text{off}} - f_{k-\tau}^{\text{off}},$$
 (31)

$$n_{k+1}^{\text{on}} = n_k^{\text{on}} + f_k^{\text{on}} - f_k^{\text{off}},$$
 (32)

$$n_k^{\text{on}}, s_k^{\text{on}}, s_k^{\text{off}}, f_k^{\text{on}}, f_k^{\text{off}} \in [0, 1],$$
 (33)

$$\sum_{k=t}^{t+H_p} \hat{Y}_k = 0$$
 (34)

Recall that the constants \tilde{C} , P^{agg} are defined in (19), (8), and the signal \hat{P}_k^b in (10). (28) uses Assumption 1. The last constraint (34) acts as a necessary condition for the QoS constraint (5) for any collection of positive numbers $\{\tilde{E}^j\}$ and any H_p that satisfies $H_p \leq H_b$ (see Appendix C).

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

is used to perform reference planning. Lemma 2: The set $\Omega_t^{t+H_p-1}$ is convex for every t and $H_p \geq 1$. Suppose that for a given τ and H_p for all t, the following signal

$$\bar{\theta}_k^a \triangleq \sum_{j=1}^N \frac{\theta_k^{a,j}}{\eta^j R_{th}^j} \left(\sum_{j=1}^N P^j\right)^{-1} \tag{35}$$

satisfies

$$\Theta_{k}^{-}(\tau) + \Gamma \le \bar{\theta}_{k+1}^{a} \le 1 - \Theta_{k}^{+}(\tau) + \Gamma,$$
 (36)

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

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

$$\Theta_k^+(\tau) = \sum_{s=k-\tau+1}^k \max\{\bar{\theta}_{s-1}^a - \bar{\theta}_s^a, 0\}, \text{ and } (38)$$

$$\Gamma = \sum_{i=1}^{N} \frac{\theta_{\text{set}}^{j}}{\eta^{j} R_{th}^{j}} \left(\sum_{i=1}^{N} P^{j} \right)^{-1}.$$
 (39)

Then the set $\Omega_t^{t+H_p-1}$ is non-empty for every t and $H_p \ge 1$. *Proof:* See appendix.

The condition on the ambient temperature in Lemma 2 is technical: we have never run into a numerical example (with time varying $\theta_k^{a,j}$) 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 total desired demand deviation of the BA or load aggregator, $\{R_k^{BA}\}_{k=t}^{t+H_p-1}$, onto the aggregate capacity set $\Omega_t^{t+H_p-1}$ to obtain the optimal ψ^* . We need the following definition:

$$(\psi^{BA})_{t}^{t+H_{p}-1} \triangleq \left[\{0\}_{t+1}^{t+H_{p}}, \{R_{k}^{BA}\}_{t+1}^{t+H_{p}}, \{0\}_{t}^{t+H_{p}-1}, \{0\}_{t}^{t+H_{p}-1}, \{0\}_{t+1}^{t+H_{p}}, \{0\}_{t+1}^{t+H_{p}} \right]. \tag{40}$$

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. $\psi \in \Omega$ (41)

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 $\{\hat{Y}_k^*\}$ of ψ^* – see the definition (25) – is denoted by R_k^* in the sequel: it is the "largest" power deviation reference, aligned with the needs of the BA or load aggregator, 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.

The objective function in (41) is aligned with the needs of the BA or load aggregators, that a collection of TCLs track a grid-supplied reference signal. This choice has the intuitive interpretation: by tracking the reference as closely as possible within their capacity, the collection of TCLs help reduce supply-demand mismatch. In general, the objective function is a design choice. In other market based scenarios, the BA or load aggregator may have alternative needs that can lead to a different objective function choice. For example, by including generation resources and transforming the objective to have units of currency (e.g. USD) yields the economic dispatch problem. However, we reiterate that the constraints present in (41) are not a design choice, and they are required to ensure the capacity of the collection of TCLs no matter the choice of objective.

 Information Requirement: In order for a BA or load aggregator to solve the reference planning problem (41), it needs to know: (i) the parameters P^{agg} , τ , \bar{C} , \bar{a} , and \bar{b} (ii) the initial conditions $n_t^{\text{on}}, f_{t-1}^{\text{on}}, \dots, f_{t-\tau}^{\text{on}}, f_{t-1}^{\text{off}}, \dots, f_{t-\tau}^{\text{off}}$, and (iii) forecasts of the signals $\theta_k^{a,j}$ (to compute the analytical baseline \hat{P}_k^b), 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 [1]. In the numerical simulations conducted later, we set the initial condition n_t^{on} to P_k^b/P^{agg} (which corresponds to $Y_t=0$ as prescribed in $\Omega_r^{t+H_p-1}$). The initial fraction of loads stuck on/off and the initial scaled temperature deviation can be obtained or estimated in several ways. For instance, the BA or load aggregator could obtain all of the initial state values through measurements from the population of TCLs. Sampling a small subset of TCLs should be enough [21]. Since many coordination algorithms will require exchange of such data with the TCLs, the required sensors are likely to be available. Otherwise, the BA or load aggregator could run a simulation with thermostat control for a time period starting from a past time and ending at the initial time for reference planning, and then use the simulation data to estimate these initial fractions. In our numerical simulations, we have assumed initial conditions for s^{off} , s^{on} to be zero for the sake of illustration, as specified in (26). This leads to the initial condition $f_{t-1}^{\text{on}} = \cdots = f_{t-\tau}^{\text{on}} = f_{t-1}^{\text{off}} = \cdots = f_{t-\tau}^{\text{off}} = 0$, i.e., since no TCLs are currently stuck then none must have previously flipped. In this case the estimated capacity is larger than the true capacity since some TCLs are likely to be stuck on or off at the initial times, which has been ignored.

A. Alternative Method for Reference Planning

To compare with past literature we define a constraint set based on the constraints developed in [8] and the scaled aggregate temperature deviation model (18) for projection of R_k^{BA} . The disadvantage with this constraint set is that it does not account for the individual cycling (4) or energy (5) constraint. This alternative reference planning problem is posed as

$$\min_{\{\hat{Y}_k\}, \{Z_k\}} \sum_{k=t}^{t+H_p-1} (R_k^{BA} - \hat{Y}_k)^2 \xi + \sum_{k=t+1}^{t+H_p} Z_k^2$$
 (42)

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

$$Z_{k+1} = \bar{a}Z_k - \bar{b}\hat{Y}_k, \quad Z_t = 0,$$
 (43)

$$|Z_{k+1}| \le \bar{C}, -\hat{P}_k^b \le \hat{Y}_k \le P^{agg} - \hat{P}_k^b,$$
 (44)

where ξ is a constant that specifies the relative importance of goals in the objective.

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 of the prior art. The simulated TCLs are residential air conditioner units (ACs). The alternative method, described in Section V-B, is designed to satisfy the indoor temperature constraint but does not account for cycling constraints of the TCLs. For a full description of the constraints in the alternative method see [8].

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 [22].

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 QoS 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 QoS is enforced by a local controller, the ensemble will not be able to track the reference. This is demonstrated by performing a 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 [8]. While the original one described in [8] 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,

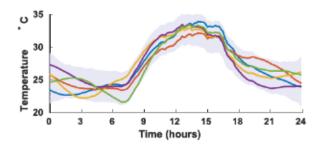


Fig. 3. Band of the ambient temperature trajectories and five sample paths. Each TCLs distinct ambient temperature trajectory lies within the shaded region.

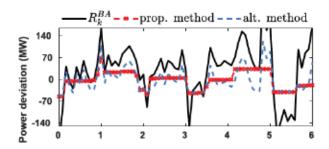


Fig. 4. BA signal (R_k^{BA}) and the reference trajectories (R_k^*) for a collection of 60,000 TCLs.

TABLE I SIMULATION PARAMETERS

Par.	Unit	value	Par.	Unit	value
N	N/A	6×10^{4}	η^j	N/A	2.5
Ĉ	MWh	50	$\theta_k^{a,j}$	°C	time var.
$T_s\tau$	Mins.	20	θ ^j set δ ^j	°C	U[21, 22]
$T_s \tau_{tcl}$	Mins.	10	δ^j	°C	U[0.75, 1]
R_{th}^{j}	°C/kW	U[1.5, 3]	T_s	Mins.	2
C_{th}^{j}	kWh/°C	U[1.5, 3]	P^j	kW	in [2.1 3.8]
\tilde{E}^{j}	kWh	6.4	P^{agg}	MW	134.4

 $^{^{*}}U[a,b]$ represents uniform distribution on [a,b]. Since the rated power P^{j} of the j-th TCL is chosen to enable maintenance of indoor climate, it is a function of the building's thermal properties. Hence, the distribution of the random variable P^{j} is complex.

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 or load aggregator supplied reference, R_k^{BA} , is obtained from BPA, a Balancing Authority in the Pacific Northwest of the United States, and is shown in Fig. 4. A heterogeneous ensemble of loads is considered. The thermal parameters for the loads are based on the values provided in [23] and these values are shown in Table I, along with other simulation parameters. The rated power P^j for the j-th air conditioner is chosen as $P^j = \frac{1}{\eta^j R_{th}^j} (35 - \theta_{sct}^j)$, so that the air conditioner is powerful enough to keep the indoor temperature at the setpoint on the hottest day, when

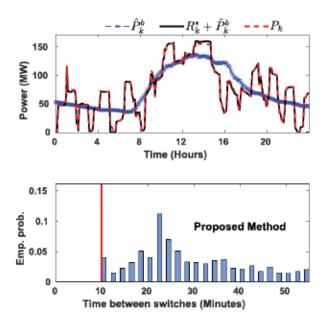


Fig. 5. Closed loop results in scenario t-i: reference planned from the proposed method. (Top): reference tracking results, (Bottom): individual TCL cycling QoS results. The vertical red line indicates τ_{tcl} .

TABLE II REFERENCE TRACKING ERRORS

Reference planning method	Tracking Error	
Proposed method (Figure 5)	0.1 %	
Alternative method (Figure 6)	25 %	

the outside temperate is 35° C, following ASHRAE guidelines for sizing of residential air conditioners [24, Chapter 17]. The ambient air temperature is time varying; it is obtained from wunderground.com for a summer day in Gainesville, FL and shown in Fig. 3. Each TCL experiences a distinct ambient temperature trajectory, as shown in Fig. 3. The unique ambient temperature sequence for each TCL is obtained by adding a random Gaussian sequence to the data collected.

Fig. 4 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 Fig. 4 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 requirements of the BA or load aggregator. 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 Fig. 5 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 Fig. 5 (bottom). Every AC

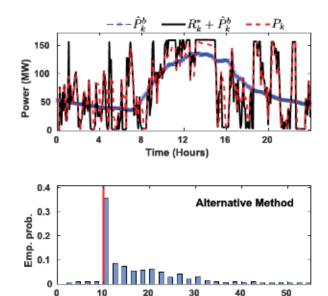


Fig. 6. Closed loop results in scenario t-ii: reference planned from the alternative method, and each TCL locally enforces its cycling constraint. (Top): reference tracking, (Bottom): individual TCL cycling QoS. The vertical red line indicates τ_{tcl} .

Time between switches (Minutes)

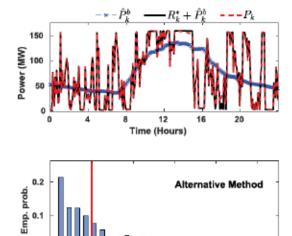


Fig. 7. Closed loop results in scenario t-iii: reference planned from the alternative method, and TCLs do not enforce their local cycling constraints. (Top): reference tracking, (Bottom): individual TCL cycling QoS. The vertical red line indicates τ_{tol} .

Time between switches (Minutes)

30

40

50

20

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 P_k is shown in Fig. 6, along with the reference (planned by the alternative method that does not 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 Fig. 7: good reference tracking at the cost of excessive cycling. Roughly 20 % of the total mode flips 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. 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 or aggressive cycling that is likely to damage compressors.

Even though sufficient conditions on aggregate quantities are not possible, there is an open problem that needs to be tackled in future work: The existence of coordination algorithms that can track aggregate references that are planned with the proposed method without violating TCLs' local QoS constraints. This question is especially pertinent for decentralized coordination algorithms. The successful tracking observed in our numerical studies was obtained by a centralized coordinator.

The reference planning method proposed here can be used by the grid operator to compute a suitable reference signal for a collection of TCLs. In practice, this can be one part of the overall resource allocation by the grid operator among multiple controllable resources to balance demand and supply. 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 robust to the quasi-homogeneity assumption: the closed loop simulations presented are for a heterogeneous population of TCLs whereas reference design utilized 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 a shortterm 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. Some work on this problem for flexible loads that do not have cycling constraints is provided in [25].

APPENDIX

A. Proof of Lemma 1

The proof roughly follows the one found in [8], with slight modification to handle time varying weather and the discrete time dynamics. By construction, the discrete time dynamics (18) is the discrete time equivalent of the ode

$$\dot{z}(t) = -\alpha z(t) - \tilde{y}(t) \tag{45}$$

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

$$\dot{\theta}^{j}(t) = \frac{1}{R_{th}^{j}C_{th}^{j}} \left(\theta^{a,j}(t) - \theta^{j}(t)\right) + \frac{\eta^{j}}{C_{th}^{j}}P^{j}(t),$$
 (46)

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

$$Z^{j}(t) \triangleq \frac{C_{th}^{j}}{n^{j}}(\theta^{j}(t) - \theta_{set}^{j})$$
 (47)

then this quantity evolves as,

$$\dot{Z}^{j}(t) = -a^{j}Z^{j}(t) - \tilde{P}^{j}(t), \quad \tilde{P}^{j}(t) = P^{j}(t) - \hat{P}^{j,b}(t),$$
(48)

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

$$\hat{P}^{j,b}(t) = \frac{\theta^{a,j}(t) - \theta_{set}^{j}}{\eta^{j}R_{th}^{j}}$$
 (49)

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

$$Z(s) = -\frac{1}{s+\alpha}\bar{Y}(s)$$
, and $Z^{j}(s) = -\frac{1}{s+\alpha^{j}}\bar{P}^{j}(s)$. (50)

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

$$Z(s) = \sum_{j=1}^{N} -\frac{1}{s+\alpha} \bar{P}^{j}(s),$$
 (51)

$$=\sum_{j=1}^{N} \frac{s+a^{j}}{s+\alpha} \frac{-1}{s+a^{j}} \tilde{P}^{j}(s)$$
 (52)

$$=\sum_{j=1}^{N} \frac{s+a^{j}}{s+\alpha} Z^{j}(s). \tag{53}$$

Now taking the inverse Laplace transform of the equation (53) 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,

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

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

$$||Z_k||_{\infty} \le \sum_{j=1}^{N} \left(1 + \left|1 - \frac{R_{th}^j C_{th}^j}{\alpha}\right|\right) ||Z_k^j||_{\infty},$$
 (55)

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 $||Z_k^j||_{\infty} \leq \frac{C_{th}^j \delta^j}{\eta^j}$ so that

$$|Z_k| \le \sum_{j=1}^N \left(1 + \left|1 - \frac{R_{th}^j C_{th}^j}{\alpha}\right|\right) \frac{C_{th}^j \delta^j}{\eta^j}, \quad \forall k,$$
 (56)

which is the desired result.

Note that Lemma 1 here appears in [8] in continuous time. In our proof we use a connection to continuous time, and the fact that our recursion (18) is an exact discretization of a certain ode. Additionally, in [8] 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 $\hat{Y}_k \equiv 0$, which together with the initial condition $Z_t = 0$ produces $Z_k \equiv 0$. Hence constraints (27) and (34) are satisfied. From the constraint (28) we have that n_k^{on} will equal

$$n_k^{\text{on}} = \bar{\theta}_k^a - \sum_{j=1}^N \frac{\theta_{\text{sct}}^j}{\eta^j R_{th}^j} \left(\sum_{j=1}^N P^j\right)^{-1} = \bar{\theta}_k^a - \Gamma, \quad (57)$$

and the difference satisfies $n_k^{\text{on}} - n_{k-1}^{\text{on}} = \bar{\theta}_k^a - \bar{\theta}_{k-1}^a$. The constraint (32) is satisfied by

$$f_{k}^{\text{off}} = \max\{\bar{\theta}_{k-1}^{a} - \bar{\theta}_{k}^{a}, 0\}, \text{ and } (58)$$

$$f_k^{on} = \max{\{\bar{\theta}_k^a - \bar{\theta}_{k-1}^a, 0\}},$$
 (59)

by definition. Upon substituting these choices in the constraints (30) and (31) and using the initial conditions, we have

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

$$s_k^{\text{off}} = \sum_{s=k-\tau+1}^k \max\{\bar{\theta}_{s-1}^a - \bar{\theta}_s^a, 0\} = \Theta_k^+(\tau)$$
 (60)

so that by hypothesis, we have

$$s_{k}^{\text{on}} + \Gamma \le \bar{\theta}_{k+1}^{a} \le 1 - s_{k}^{\text{off}} + \Gamma,$$
 (61)

which implies that

$$s_k^{\text{on}} \le n_{k+1}^{\text{on}} \le 1 - s_k^{\text{off}},$$
 (62)

and hence the constraint (29) is satisfied. Additionally, by construction $n_k^{\rm on}$ satisfies (33) and since $f_k^{\rm off}$ and $f_k^{\rm on}$ are the positive difference of successive values of $n_k^{\rm on}$, they too will satisfy (33). By construction $s_k^{\rm off}$ and $s_k^{\rm on}$ are non-negative. Further from the constraint (29) holding we have $s_k^{\rm on} \leq 1$. Since the fraction of loads stuck on and off satisfy $s_k^{\rm off} + s_k^{\rm on} \leq 1$ we have that $s_k^{\rm off} \leq 1$. Hence, both $s_k^{\rm on}$ and $s_k^{\rm off}$ satisfy (33).

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

C. VES Constraint

The BA or load aggregator requires the constraint (34) to ensure that the collection of TCLs do not act as generators. We repeat this constraint here for t=0:

$$\sum_{k=0}^{H_p} \hat{Y}_k = 0. {(63)}$$

We now show that this constraint is a necessary condition for the individual TCLs energy constraint (5). Summing (5) over the j index and expanding the absolute value,

$$-\sum_{j=1}^{N} \bar{E}^{j} \leq T_{s} \sum_{j=1}^{N} \sum_{k=0}^{H_{b}} (m_{k}^{j} P^{j} - \hat{P}_{k}^{j,b}) \leq \sum_{j=1}^{N} \bar{E}^{j}. \quad (64)$$

$$\Rightarrow -\sum_{j=1}^N \tilde{E}^j \leq T_s \sum_{k=0}^{H_b} \sum_{j=1}^N (m_k^j P^j - \hat{P}_k^{j,b}) \leq \sum_{j=1}^N \tilde{E}^j,$$

$$\Rightarrow -\sum_{j=1}^{N} \bar{E}^{j} \le T_{s} \sum_{k=0}^{H_{b}} \hat{Y}_{k} \le \sum_{j=1}^{N} \bar{E}^{j}.$$
 (65)

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

$$T_s \left| \sum_{k=0}^{H_b} \hat{Y}_k \right| \le \sum_{j=1}^{N} \bar{E}^j.$$
 (66)

If $H_b = H_p$, then due to (63) the above will be true for all values of \tilde{E}^j , as the RHS term in (66) is defined to be greater than or equal to zero. If $H_p < H_b$ then this means the planning horizon is shorter then the billing horizon. In reality, the planning problem is solved indefinitely: after the first segment of length H_p it is solved again for another segment of length H_p . Since the constraint (63) is enforced in each segment, the LHS of (66) will also be zero.

If (66) 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 (65); The constraint (63) still enforces this.

REFERENCES

- P. Barooah, Smart Grid Control: An Overview and Research Opportunities. Berlin, Germany: Springer-Verlag, 2019.
- [2] 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 Proc. Clemson Univ. Power Syst. Conf., 2018, pp. 1–6.
- [3] J. L. Mathieu, S. Koch, and D. S. Callaway, "State estimation and control of electric loads to manage real-time energy imbalance," *IEEE Trans. Power Syst.*, vol. 28, no. 1, pp. 430–440, Feb. 2013.
- [4] 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* Control Energy Markets Grids, 2017, pp. 1–26.
- [5] W. Zhang, J. Lian, C.-Y. Chang, and K. Kalsi, "Aggregated modeling and control of air conditioning loads for demand response," *IEEE Trans. Power* Syst., vol. 28, no. 4, pp. 4655–4664, Nov. 2013.

- [6] A. Coffman, A. Bušić, and P. Barooah, "Virtual energy storage from TCLs using QoS preserving local randomized control," in *Proc. 5th ACM Int. Conf. Syst. Built Environments*, 2018, pp. 93–102.
- [7] M. Liu and Y. Shi, "Model predictive control of aggregated heterogeneous second-order thermostatically controlled loads for ancillary services," *IEEE Trans. Power Syst.*, vol. 31, no. 3, pp. 1963–1971, May 2016.
- [8] H. Hao, B. M. Sanandaji, K. Poolla, and T. L. Vincent, "Aggregate flexibility of thermostatically controlled loads," *IEEE Trans. Power Syst.*, vol. 30, no. 1, pp. 189–198, Jan. 2015.
- [9] L. Zhao, W. Zhang, H. Hao, and K. Kalsi, "A geometric approach to aggregate flexibility modeling of thermostatically controlled loads," *IEEE Trans. Power Syst.*, vol. 32, no. 6, pp. 4721–4731, Nov. 2017.
- [10] S. Barot and J. A. Taylor, "A concise, approximate representation of a collection of loads described by polytopes," *Int. J. Elect. Power Energy* Syst., vol. 84, pp. 55–63, 2017.
- [11] C. Ziras, S. You, H. W. Bindner, and E. Vrettos, "A new method for handling lockout constraints on controlled TCL aggregations," in *Proc. Power Syst. Computation Conf.*, 2018, pp. 1–7.
- [12] B. M. Sanandaji, T. L. Vincent, and K. Poolla, "Ramping rate flexibility of residential HVAC loads," *IEEE Trans. Sustain. Energy*, vol. 7, no. 2, pp. 865–874, Apr. 2016.
- [13] D. Cheng, W. Zhang, and K. Wang, "Hierarchical reserve allocation with air conditioning loads considering lock time using benders decomposition," *Int. J. Elect. Power Energy Syst.*, vol. 110, pp. 293–308, 2019.
- [14] P. Wang, D. Wu, and K. Kalsi, "Flexibility estimation and control of thermostatically controlled loads with lock time for regulation service," *IEEE Trans. Smart Grid*, vol. 11, no. 4, pp. 3221–3230, Jul. 2020.
- [15] A. R. Coffman, A. Bušić, and P. Barooah, "Aggregate capacity for TCLs providing virtual energy storage with cycling constraints," in *Proc. IEEE Conf. Decis. Control*, 2019, pp. 4208–4215.

- [16] Y. Chen, Markovian Demand Dispatch Design for Virtual Energy Storage to Support Renewable Energy Integration. Gainesville, FL, USA: Univ. florida press, 2016.
- [17] A. R. Coffman, A. Bušić, and P. Barooah, "A unified framework for coordination of thermostatically controlled loads," pp. 1–21, 2021, arXiv:2108.05840.
- [18] A. R. Coffman, N. Cammardella, P. Barooah, and S. Meyn, "Flexibility capacity of thermostatically controlled loads with cycling/lock-out constraints," in *Proc. Amer. Control Conf.*, 2020, pp. 527–532.
- [19] N. Mahdavi, J. H. Braslavsky, and C. Perfumo, "Mapping the effect of ambient temperature on the power demand of populations of air conditioners," *IEEE Trans. Smart Grid*, vol. 9, no. 3, pp. 1540–1550, May 2018.
- [20] 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 Buildings*, vol. 231, Jan. 2021, Art. no. 110500.
- [21] Y. Chen, A. Bušić, and S. P. Meyn, "State estimation for the individual and the population in mean field control with application to demand dispatch," *IEEE Trans. Autom. Control*, vol. 62, no. 3, pp. 1138–1149, Mar. 2017.
- [22] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 1.21," Feb. 2011.
- [23] 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 *Proc. ACEEE Summer Study Energy Efficiency Buildings*, 2012, pp. 1-189–1-203.
- [24] American Society of Heating, Refrigerating and Air-Conditioning Engineers, "The ASHRAE handbook fundamentals (SI Edition)," Ch. 17, Feb. 2011. [Online]. Available: http://cvxr.com/cvx
- [25] A. R. Coffman, Z. Guo, and P. Barooah, "Characterizing capacity of flexible loads for providing grid support," *IEEE Trans. Power Syst.*, vol. 36, no. 3, pp. 2428–2437, May 2021.