

A Packetized Energy Management Macromodel With Quality of Service Guarantees for Demand-Side Resources

Luis A. Duffaut Espinosa , *Member, IEEE*, and Mads Almassalkhi , *Senior Member, IEEE*

Abstract—Using distributed energy resources (DERs), such as thermostatically controlled loads (TCLs), electric vehicles (EVs), and energy storage systems (ESSs) as a way to manage demand has been known for decades. A demand management scheme that explicitly considers the individual DER's local quality of service (QoS) is known as demand dispatch. Packetized energy management (PEM) is a demand dispatch paradigm that borrows packet-based concepts from wireless communications to dynamically manage fleets of DER at-scale and in realtime via small, discrete fixed-duration/fixed-power energy packets. PEM addresses QoS in a bottom-up fashion by having a coordinator authorize/deny incoming requests from DERs to consume energy packets. This manuscript extends prior work on modeling a large-scale population (i.e., macro-model) of homogeneous TCLs and ESSs operating under the PEM paradigm. In particular, we extend the macro-model methodology to include deferrable loads (DLs), such as EVs, together with analysis of QoS guarantees. Comparisons between an agent-based (micro-model) simulation and the proposed macro-model are presented to validate modeling accuracy and QoS guarantees.

Index Terms—Distributed energy resources, packetized energy management, demand dispatch, relay control, modeling.

ACRONYMS

CMC	Controlled Markov chain
DER	Distributed energy resource
DL	Deferrable load
DR	Driving mode
EWB	Electric water heater
ESS	Energy storage system
EV	Electric vehicle
MTTR	Mean time-to-request
PEM	Packetized energy management
PRP	Poisson rectangular pulse
PV	Photo-voltaic

Manuscript received June 27, 2019; revised January 19, 2020; accepted March 7, 2020. Date of publication March 17, 2020; date of current version August 24, 2020. This work was supported in part by the U.S. Department of Energy's Advanced Research Projects Agency - Energy (ARPA-E) award DE-AR0000694 and in part by NSF Grant CMMI-1839387. Paper no. TPWRS-00921-2019. (Corresponding author: Luis Duffaut Espinosa.)

The authors are with the Department of Electrical and Biomedical Engineering, University of Vermont, Burlington, VT 05405 USA (e-mail: lduffaut@uvm.edu; malmassa@uvm.edu).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TPWRS.2020.2981436

QoS	Quality of service
SB	Standby mode
SoC	State of charge
TCL	Thermostatically controlled load
VPP	Virtual power plant

NOMENCLATURE

0_N	Zero matrix of size N -by- N .
I_N	Identity matrix of size N -by- N .
$\mathbf{1}_N$	Column vector of ones of length N .
β_h	Proportion of requests accepted for $h \in \{c, d\}$.
β_h^-	Proportion of population in $h \in \{c, d\}$ switching to standby.
ϕ_n	Switching mode of n -th DER.
m_R	Frequency that defines MTTR at the set point.
n_r^c, n_r^d	Number of charge (c) or discharge requests (d).
$\eta_{sl,n}, \eta_{c,n}, \eta_{d,n}$	Standby, charging, and discharging energy loss parameters of n -th DER.
N_i	i -th Poisson process with parameter λ_i .
N_e	Total number of DERs.
$\mathcal{N}(\mu, \sigma)$	Normal distribution with mean μ and standard deviation σ .
p_{ij}^h	Transition probability from bin i to bin j for a specific $h \in \{c, sb, d\}$.
$p_i^{\text{req},h}$	Request probability of standby states for $h \in \{c, d\}$.
$P_{c,n}^{\text{rate}}, P_{d,n}^{\text{rate}}$	Energy transfer rates of the n -th DER when charging (c) or discharging (d).
P_{dem}	Aggregated demand power.
P_{ref}	Balancing signal or reference power.
q_h	Vector of population percentages for $h \in \{c, sb, d\}$.
$u_{\phi_n,n}$	Power input of n -th DER.
w_n	End-user event process of n -th DER.
\mathcal{X}_h	Finite set of states with elements x_h^i for $h \in \{c, sb, d\}$.
$x_{p,h}$	Timer state for charging ($h = c$) or discharging ($h = d$) modes.
z_n	Dynamic state of n -th DER.
z_n^+	Dynamic state update of n -th DER.
$\underline{z}_n, \bar{z}_n$	Lower and upper dynamic state boundaries.

z_n^{set}	Deadband set point of n -th DER.
$z_{\text{SB}}, z_{\text{DR}}$	Standby (SB) and driving (DR) mode for EVs.

I. INTRODUCTION

SINCE the early 1980s, aggregated DERs have been known to be capable of significant actuation in bulk power systems [1]. Yet, since then, the technology deployment on this front has been underwhelming. However, recently, demand management has become the centerpiece of bold renewable portfolio standards as the means to integrate large-scale, intermittent renewable generation. In [2], the authors illustrate how fleets of DERs can be employed in transmission and distribution system operations to manage the variability from renewable generation and to provide relevant grid services. A fleet of DERs in this context may consists of TCLs, such as electric water heaters (EWHs), bidirectional ESSs, such as Enphase's AC Batteries, and DLs, such as EVs. These principles were expanded upon in [3] where a state-bin transition (macro) model was developed for a fleet of TCLs. The TCLs then transition probabilistically between ON and OFF based on a broadcasted control signal. While this framework depends on solving a challenging state-estimation problem and may not always be observable [4], it has been analyzed and extended to include interesting use cases [5]–[7]. Related works with state bin transition models of TCLs has also focused on higher order models [8], and compressor constraints [9], and analyzing the aggregation abstraction error for populations of TCLs [10].

Ineffective management of QoS will drive DER owners (i.e., humans) to permanently opt-out of the scheme, which reduces the availability of flexible resources and limits the long-term viability of DER coordination programs. A demand management scheme that explicitly considers QoS is known as *demand dispatch* [11]. The work based on [4] employs a novel mean-field model that via linearizations can be well-approximated by a single-input, single-output (SISO) model for fleets of pool pumps, fleets of TCLs, and fleets of ESSs [12]. The demand dispatch approach then broadcasts a single scalar control signal that perturbs the transition probabilities of all DERs of the same type (from a given baseline) and uses measured power of the fleet as feedback. This line of work has since been expanded to include opt-out control to improve QoS [13] and device-level filter design to consider heterogeneity in the fleet [14].

This manuscript focuses on a demand dispatch framework that uses low-bandwidth, bidirectional communications between a device and the coordinator and also includes QoS guarantees via opt-out control. It is called *packetized energy management* (PEM) [15]–[18]. PEM leverages packet-based strategies from random access communication channels that have previously been applied to the distributed management of wireless sensor networks (i.e., similar to ALOHA protocol, but with multiple channels). PEM enables the delivery of energy to DERs via multiple fixed-duration/fixed-power (charging or discharging) *energy packets*, similar to how digital communication networks enable the transmission of small, kB-scale data packets rather than bulky files. Unlike other approaches, PEM is device-driven

and does not broadcast a control signal to all DERs. Instead, PEM, in *bottom-up* fashion, is designed to have each DER probabilistically *request* an energy packet from the coordinator based on the DER's local *need for energy*. This gives the coordinator the ability to respond in realtime to incoming (asynchronous) packet request based on grid and/or market conditions. Other work related to packet-based coordination is [19], where a packet control algorithm is proposed that requires just binary information from each DER at each time instant (in bottom-up fashion) and with the drawbacks of synchronized packet acceptances and the need for continuously queuing packet requests, which serves as memory but creates delays in service. This manuscript presents a complete macro-model for PEM for diverse demand-side resources. Specifically, the technical contributions are:

- 1) Generalized modeling of DER end-user events is presented in the form of Poisson rectangular pulses (PRPs) and analysis provides key event statistics that are used to improve modeling of transition probabilities.
- 2) Development of a generalized bin transition macro-model of a fleet of diverse DERs under PEM. Unlike the authors' prior work in [16], this macro-model is now able to describe a fleet of packetized EV chargers by including a new explicit embedding for uncontrollable transitions between modes of operation, such as "driving" and "charging."
- 3) Analysis of QoS guarantees are provided by including opt-out dynamics and bounds on uncontrollable end-user events.
- 4) Simulation-based analysis provides validation of the macro-model with deferrable loads and QoS results.

The manuscript is organized as follows. Section II provides a review of PEM fundamentals. In Section III, a description of end-user events for DERs is given and it is followed by the development of a bin transition Markov model for PEM. Section IV presents analysis of QoS guarantees under PEM with particular emphasis on the case of DLs, such as EVs. The conclusions are given in the final section.

II. FUNDAMENTALS OF PACKETIZED ENERGY MANAGEMENT

The following high-level description summarizes the bottom-up approach that is PEM and is detailed in [15], [16]:

- 1) A DER estimates its local state of charge, SoC.
- 2) If the SoC is within a predefined range of comfort, the DER, based on its state, probabilistically requests to consume energy from the grid at a fixed rate (e.g., 4 kW) for a pre-specified epoch (e.g., 5 minutes) to beget an energy packet (e.g., 0.33 kWh). If the SoC exceeds a local pre-defined range (e.g., too low), the DER automatically opts out of PEM (to guarantee QoS) and reverts to a default control mode (e.g., charges) until the SoC is returned within limits when it opts back into to PEM.
- 3) If a request is received, the coordinator or *Virtual Power Plant* (VPP) either accepts or denies the DER's packet request based on grid or market conditions. If the request is denied, go to *i*). If the request is accepted, consume the energy packet and then go to *i*).

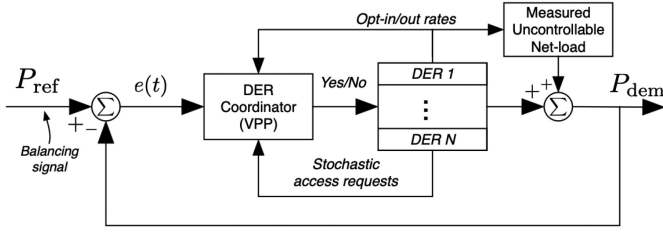


Fig. 1. Closed-loop feedback system for PEM with P_{ref} provided by the grid or market operator and the aggregate net-load P_{dem} measured by VPP.

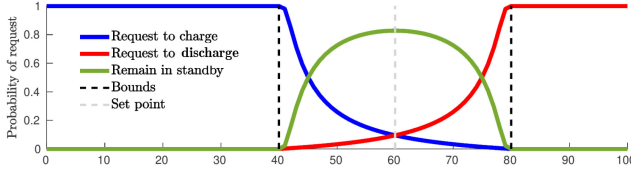


Fig. 2. Illustrating the charge/discharge energy packet request rates and MTTR for a generic packetized DER. Note that (3) is represented by the blue line (left to right top plot). Top plot gives the effect of local state z_n (e.g., state-of-charge) on the packet request probabilities and bottom plot provides the corresponding MTTR of a packetized DER under PEM.

The above scheme can ensure consumer's QoS for a heterogeneous fleet of electric water heaters by including the opt-out control when the SoC falls below a certain pre-defined threshold. At the same time, randomization is injected to the request rule based on the local SoC, which limits synchronization and promotes equitable access to the grid. Fig. 1 illustrates the closed-loop system under PEM.

In a fleet of diverse DERs, the general discrete-time dynamic model for the n -th DER having SoC z_n is given by

$$z_n^+ = f_n(z_n, \phi_n, P_{c,n}^{\text{rate}}, P_{d,n}^{\text{rate}}, w_n), \quad (1)$$

where f_n is a one-dimensional mapping (usually linear or bilinear), w_n is the parameter mapping end-consumer usage to the energy state, $P_{c,n}^{\text{rate}}$ and $P_{d,n}^{\text{rate}}$ are the energy transfer rates of the n -th DER when charging (c) or discharging (d), respectively, and ϕ_n is the hybrid state of the DER dynamics. Here ϕ_n take values in the set $\{c, sb, d\}$ that corresponds to the $\{\text{charge, standby, discharge}\}$ modes, respectively [15], [16]. In this manuscript, the focus is on EWHs, ESSs, and EVs (as deferrable loads). The latter represents the first instance of EVs for a PEM macromodel.

Since EWHs were presented in [15], [16], consider (1) for an ESS with background power usage $w_n \in \mathbb{R}$ and charge (discharge) rate limits of $P_{c,n}^{\text{rate}}$ ($P_{d,n}^{\text{rate}} > 0$). Then,

$$z_n^+ = \eta_{sl,n} z_n + u_{\phi_n,n} \eta_{\phi_n,n} + w_n, \quad (2)$$

where $u_{\phi_n,n}$ is $P_{c,n}^{\text{rate}}$ for $\phi_n = c$, $-P_{d,n}^{\text{rate}}$ for $\phi_n = d$ and 0 for $\phi_n = sb$ in [kW], and $\eta_{sl,n}$, $\eta_{c,n}$, and $\eta_{d,n}$ are the standing losses, charging, and discharging parameters, respectively. Herein, an ESS is modeled as a bidirectional battery, in which, simultaneous charging and discharging are not possible. Fig. 2 illustrates the ESS probabilistic request mechanism.

In a similar manner, an EV is modeled as an ESS, however, driving is assumed the only means of discharging.¹ Therefore, EVs are modeled as in (2), where $\phi_n \in \{c, sb, d\}$ is the hybrid state corresponding to charge/standby/driving, respectively, $u_{\phi_n,n} = P_{c,n}^{\text{rate}}$ [kW] is the control input equal to the EV's charging rate when $\phi_n = c$ and $u_{\phi_n,n} = 0$, otherwise. A key difference between ESSs and EVs is that w_n is dependent on ϕ_n . That is, $w_n = 0$ for $\phi_n = \{c, sb\}$ and $w_n \neq 0$ is the power consumed by the EV's battery when driving ($\phi_n = d$). Given that the timescale of interest in this work is hours and minutes, EVs in this manuscript assumes no standing losses (i.e., $\eta_{sl,n} = 1$). Since we will be modeling a population of EVs, we employ the simplifying assumption that the (average) discharging rate, w_n , is a constant value based on the notion of the (average) driving speed of EVs. For example, at a speed of 50mph, the discharging rate is approximately 7kW, which provides 3 hours (or ca. 150 miles) of continuous driving for a battery with capacity of 22.5kWh, which is reasonable on average [20]. This assumption can be relaxed by considering a distribution of driving rates rather and is a topic for ongoing work and outside the scope of this manuscript.

In this context, the discrete-time implementation of PEM assigns a probability of requesting access to the grid to the packetized load n based on its local SoC $z_n[k] \in [\underline{z}_n, \bar{z}_n]$ and desired set-point $z_n^{\text{set}} \in (\underline{z}_n, \bar{z}_n)$ during time-step k (over interval Δt). This request probability has been defined by the cumulative exponential distribution function, $P(z_n[k]) := 1 - e^{-\mu(z_n[k])\Delta t}$, where the rate parameter $\mu(z_n[k]) > 0$ is dependent on the SoC. Denoting by $P_k^h(n|Q)$ the probability that DER n requests a packet for consumption ($h = c$) or injection ($h = d$) given condition Q is satisfied. While any request probability function would suffice, the key is a mapping of the SoC to the request probability that considers the boundary conditions:

- i) $P_k^c(n|z_n[k] \leq \underline{z}_n) = 1 \quad \wedge \quad P_k^d(n|z_n[k] \geq \bar{z}_n) = 0$,
 - ii) $P_k^d(n|z_n[k] \leq \underline{z}_n) = 0 \quad \wedge \quad P_k^c(n|z_n[k] \geq \bar{z}_n) = 1$,
- from which i) gives rise to the following helpful design of $\mu(z_n[k])$ for *consuming* a packet:

$$\mu(z_n[k]) = \begin{cases} 0, & \text{if } z_n[k] \geq \bar{z}_n \\ m_R \left(\frac{\bar{z}_n - z_n[k]}{z_n[k] - \underline{z}_n} \right) \cdot \left(\frac{z_n^{\text{set}} - \underline{z}_n}{\bar{z}_n - z_n^{\text{set}}} \right), & \text{if } z_n[k] \in (\underline{z}_n, \bar{z}_n) \\ \infty, & \text{if } z_n[k] \leq \underline{z}_n \end{cases} \quad (3)$$

where $m_R > 0$ [Hz] is a design parameter that defines the mean time-to-request (MTTR) at z_n^{set} . For example, if one desires a MTTR of 5 minutes when $z_n[k] \equiv z_n^{\text{set}}$ then $m_R = \frac{1}{300}$ Hz. The design of $\mu(z_n[k])$ for *injecting* a packet is described in similar fashion, but with boundary conditions ii) above. Fig. 2 maps boundary conditions to charging and discharging packet requests. Next, we present a general model for DER end-use events, which is then embedded in a state bin transition model for PEM for a large population of DERs.

¹ We do not consider Vehicle-to-Grid (V2G) capability in this model, but could be included in future work.

III. STATE TRANSITION MODEL UNDER END-USER EVENTS

This section develops a state bin transition macro-model for a large population of *packetized DERs*, which explicitly captures the unique packet request-notification dynamics inherent to PEM. In particular, the cases for TCLs, ESSs and EVs are provided. A macro-model for a diverse population of multiple DER types with charging and discharging is comprised of a finite number of homogeneous populations of DERs coordinated under the same VPP. However, each class of DERs is affected by different types of end-user events, which makes the aggregation of homogeneous DERs behave differently depending upon the DER class. Therefore, the discussion is initially focused on the modeling of end-user events: hot water usage, unscheduled power consumption/injection, and driving behavior.

A. Modeling End-User Events for DERs

The end-user events are uncontrollable and modeled employing a simple birth/death stochastic differential equation for the process, $w_n(t)$. In this regard, the main assumption for choosing a user model is that water consumption starts with certain probability and stops with another. If one thinks of starting (stopping) water events as independent from each other, then a reasonable assumption is that these occur with an exponentially distributed inter-arrival time. This amounts to a Poisson process for starting water events and another for stopping water events. The parameters of these two processes can be chosen so that the average time between start and stop events is related to the average historical usage at some time during a 24-hour period. These assumptions permit to formulate a model for a process of this kind in a manner that the aggregate statistics of the aggregation of a number of these processes can be computed analytically. For the sake of simplicity, the 24-hour variation is neglected in our simulations, however, the intensity of usage is modeled with an appropriate random variable whose mean is fixed. To clarify notation, the subscript n is omitted hereafter as this section focuses on a single DER at first (and later we extend to a population average). Assume that there exists an appropriate probability space (Ω, P, \mathcal{F}) , where Ω is the set of events, \mathcal{F} a filtration, and P the probability measure of elements in \mathcal{F} . For this purpose, a *Poisson rectangular pulse* (PRP) stochastic differential model is employed [21]. That is,

$$dw(t) = (v(t) - w(t)) dN_1(t) - w(t) dN_2(t), \quad (4)$$

where N_1 (N_2) is an independent, stationary Poisson point process with constant rate parameter λ_1 (λ_2), representing the starting (stopping) of a random end-user event and $v(t)$ is a random variable independent of N_1 and N_2 that describes the intensity of the end-user event and is appropriate for the type of DER under study. For example, with an ESS, v may have a symmetric probability density function with mean approximately zero and with an EV, driving behavior can be approximated as (4) with v as a random variable corresponding to the driving speed of the EV. However, EVs are considered a special type of DER in that they become unavailable to PEM when driving. This means that when an EV end-user event occurs the EV's hybrid state change from sb to d, and when driving concludes d

to sb. This mobility is a fundamental feature that differentiates EVs from other DERs.

The statistics of the aggregated behavior of end-user events are employed in the next section for the computation of the (average) transition probabilities for a fleet of DERs.

1) *TCL and ESS Populations*: A reasonable assumption is that the end-user event for each DER are independent and identically distributed random processes. Denote the expected value of the random process w as $\bar{w}(t) := E[w(t)]$. Due to the assumed independence of the processes N_1 , N_2 and v in time, one can compute the expected end-user event for each DER as

$$\frac{d\bar{w}(t)}{dt} = (\bar{v}(t) - \bar{w}(t))\lambda_1 + \bar{w}(t)\lambda_2. \quad (5)$$

The solution of (5) when $w(0) = 0$ is

$$\bar{w}(t) = E[v] \frac{\lambda_1}{\lambda_1 + \lambda_2} (1 - \exp(-(\lambda_1 + \lambda_2)t))$$

The expected event reaches steady state as t goes to infinity. Hence, the mean of end-user event in steady state is

$$\bar{w}_{\text{sst}} := \lim_{t \rightarrow \infty} \bar{w}(t) = \frac{E[v]\lambda_1}{\lambda_1 + \lambda_2}. \quad (6)$$

The next theorem describes the probability distribution of these events as the number of devices increases.

Theorem 1: The steady-state aggregation of individual end-user events, w , is distributed as $\mathcal{N}(\mu_w, \sigma_w/\sqrt{N_e})$, where N_e is the total number of end-user event processes and μ_w and σ_w are the corresponding expected value and standard deviation of the process w in steady state.

Proof: The proof is based on deriving the differential equation for the characteristic function of w in (4) from a direct application of the Itô chain rule for jump processes [22]. Let $F_\kappa(w) = e^{i\kappa w}$, then

$$d e^{i\kappa w} = (F_\kappa(v) - F_\kappa(w)) dN_1 + (1 - F_\kappa(w)) dN_2.$$

By definition, the characteristic function of w is given by $E[F_\kappa(w)] =: \Psi_w(\kappa, t)$ and $E[N_i(t)] = \lambda_i$. It then follows that

$$\frac{d\Psi_w(\kappa, t)}{dt} = \Psi_v(\kappa, t)\lambda_1 + \lambda_2 - \Psi_w(\kappa, t)(\lambda_1 + \lambda_2). \quad (7)$$

In steady state, $\frac{d\Psi_w(\kappa, t)}{dt} = 0$. Thus,

$$\Psi_w(\kappa, \infty) = \frac{\lambda_2 + \Psi_v(\kappa)\lambda_1}{(\lambda_1 + \lambda_2)}.$$

Clearly, the moments of w in steady state can be obtained by computing $E[w^n] = (-i)^n d\Psi_w(\kappa, \infty)/d\kappa|_{\kappa=0}$. A direct application of the central limit theorem for i.i.d random variables completes the proof given that in steady state all end-user events are independent of each other and identically distributed with the distribution associated to the solution of (7). Hence one can consider, on average, that a single DER is driven by a process $w \sim \mathcal{N}(\mu_w, \sigma_w)$. ■

The previous theorem simply states that the aggregation of PRP realizations behaves on average as a Gaussian process. It also allows the computation of PRP aggregation statistics, which is illustrated next in Example 1.

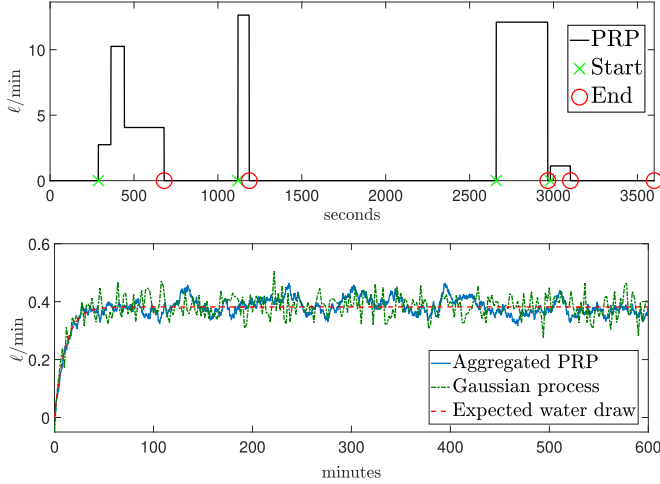


Fig. 3. PRP simulation. (Top) Realization of one PRP using (4). (Bottom) Average of 2000 end-user events modeled by PRPs compared against the aggregation given in Theorem 1 with same parameters as in Example 1.

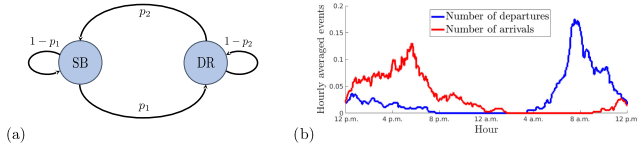


Fig. 4. a) Simple driving model. z_{SB} is the state in which an EV is in standby mode, and z_{DR} is the state representing when an EV is in driving mode. b) NHTS data indicating the average number of arrivals and departures over a 24 hour period [20].

Example 1: If $v \sim \exp(\lambda)$, then $\mu_w = \lambda p$ and $\sigma_w = \lambda \sqrt{2p - p^2}$, where $p := \frac{\lambda_1}{\lambda_1 + \lambda_2}$. In Fig. 3, the average of 2000 water usage profiles generated following (4) is compared against the aggregation model given by Theorem 1 with $\lambda = 2.1$ liters per minute, $\lambda_1 = 1/3600 \text{ sec}^{-1}$ and $\lambda_2 = 1/800 \text{ sec}^{-1}$. In this particular case, the mean and standard deviation of the aggregated PRP are 0.3868 and 0.0382, respectively. Using Theorem 1 gives a mean of 0.3818 and standard deviation of 0.0369, which is very close for a small population of 2000. Hence one can consider, on average, that a single DER is driven by a process $w \sim \mathcal{N}(\mu_w, \sigma_w)$. ■

2) EV Population: The modeling of the aggregated behavior of EVs includes the (mobility) transitions from standby to driving ($SB \rightarrow DR$) and driving to standby ($DR \rightarrow SB$). The EV driver model consists of a two-state Markov chain and is illustrated in Fig. 4 with p_1 the probability of going from SB to DR and p_2 the probability of going from DR to SB. It is assumed that the driving model is independent of the energy state of the EV population, which simplifies integration of this model together with the population model described in the subsequent section. More complex driver models could be developed and integrated with a PEM macromodel, but such detailed EV driver models are considered outside the scope of this manuscript.

Many useful and simple EV driving metrics, such as average standby and driving times, can be derived from this model. In this regard, a discrete-time model for only the driving behavior

of EVs, with time step Δt , is given by

$$\begin{pmatrix} z_{SB}[k+1] \\ z_{DR}[k+1] \end{pmatrix} = \begin{pmatrix} 1-p_1 & p_2 \\ p_1 & 1-p_2 \end{pmatrix} \begin{pmatrix} z_{SB}[k] \\ z_{DR}[k] \end{pmatrix}. \quad (8)$$

Clearly, this model permits a non-trivial, unique stationary distribution for probabilities p_1 and p_2 . The stationary distribution provides the *averaged occupancy* of each state, which is the percentage of the EVs that, on average, are either in z_{SB} or z_{DR} . Moreover, p_1 and p_2 can be chosen from driving data [20], as in Fig. 4b). For example, from this data set, the average driving duration in urban cities is about 30 minutes, which is what we use in this manuscript. That is, if $\pi = (\pi_{z_{SB}}, \pi_{z_{DR}})^T$ denotes such stationary distribution, then one has that $\pi_{z_{SB}} = \frac{p_2}{p_1 + p_2}$ and $\pi_{z_{DR}} = \frac{p_1}{p_1 + p_2}$. Here the occupancy of the standby states is provided by $\pi_{z_{SB}}$, and the occupancy of the driving states is provided by $\pi_{z_{DR}}$. The number of time steps that an EV spends driving (also known as *sojourn time*) is computed from the transition probability in (8). Denote with $\psi(z_{DR})$ the expected number of time steps needed to reach state z_{SB} given that one starts in z_{DR} and $\psi(z_{SB})$ if one were to start in state z_{SB} . Forcing the state z_{SB} to be absorbing, it follows that $\psi(z_{SB}) = 0$ and $\psi(z_{DR}) = 1 + (1-p_2)\psi(z_{DR})$, which provides $\psi(z_{DR}) = 1/p_2$. Thus, p_2 describes how many time steps an average EV stays in driving mode, and the actual expected time spent in driving mode is trivially $t_{z_{DR}} = \psi(z_{DR})\Delta t = \Delta t/p_2$.

Example 2: Consider an EV population where the data shows that the expected time spent driving is 30 minutes, then $p_2 = 1/120$ for $\Delta t = 15$ sec. Since the occupancy of the driving state is given by $\pi_{z_{DR}}$, one has that $p_1 = p_2(1/\pi_{z_{DR}} - 1) \approx 0.00092$ for an occupancy of the driving state of $\pi_{z_{DR}} = 0.1$. ■

B. Dynamics of State Bin Transitions

In this section, end-user event models are embedded into the state bin transition description. Consider a population of DERs obeying some dynamical equation (1) and with common underlying state space. To create a finite state abstraction (i.e., a macro-model) of the entire population's evolution, the state space is discretized in a manner that the main features of the system are preserved and the system as a whole is such that the effects of an individual DER is negligible with respect to the average behavior [23]. Clearly, the spatial and temporal discretization strongly affects the modeling the aggregated dynamics [24], [25]. Therefore, in this manuscript, an appropriate discretization time step Δt (used to obtain the discrete model (1)) was chosen so that it is ensured that only contiguous bin transitions occur [26]. Note that the approach described below has the capability of overcoming such restriction. The focus of this manuscript is on DERs that have hybrid one dimensional dynamics as in (1). More specifically, an interval $[z, \bar{z}]$ within Z is divided into N consecutive bins each corresponding to a *bin state* in \mathcal{X} , where $x_i \in \mathcal{X}$ corresponds to the interval $[z_{i-1}, z_i] \subset [z, \bar{z}]$. Since (1) includes three types of dynamics (charge/standby/discharge)², the state space for the

²The ON/OFF dynamics of (1) can be seen as charging/standby dynamics with a disconnected/inaccessible and trivial discharging dynamics, and the driving dynamics of EVs simply corresponds to uncontrollable discharging dynamics.

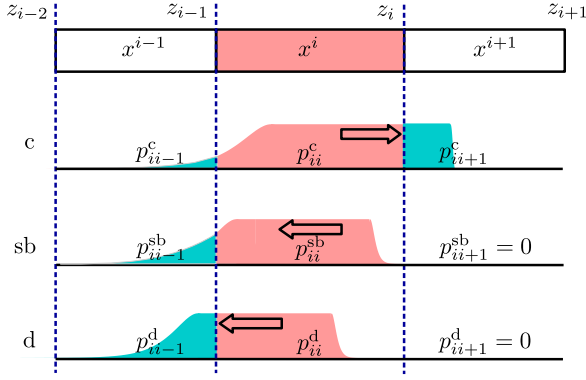


Fig. 5. Typical transition rate calculation for charging, standby and discharging states.

system consists of the union of the full state space given by $\mathcal{X} = \mathcal{X}_c \cup \mathcal{X}_{sb} \cup \mathcal{X}_d$. At time k , the probability mass function of the system is $q^\top = (q_c^\top, q_{sb}^\top, q_d^\top)$ with $q_h = (q_h^1, \dots, q_h^N)^\top$ for $h = \{c, sb, d\}$. Note that q contains the percentage of the population in each state of \mathcal{X} . For example, if N_e is the total number of DERs and $N_{e,c}^i$ is the number of devices in state x_c^i , then $N_{e,c}^i = q_c^i N_e$. Similarly, the percentage of N_e that is charging and discharging, and the total power of the system are

$$y_c = c_c q, \quad y_d = c_d q, \quad \text{and } y = c q, \quad (9)$$

where $c_c = (\mathbf{1}_N^\top, 0 \dots 0) \in \mathbb{R}^{3N}$, $c_d = (0 \dots 0, \mathbf{1}_N^\top) \in \mathbb{R}^{3N}$, $c = N_e P^{\text{rate}}(c_c - c_d) \in \mathbb{R}^{3N}$, $\mathbf{1}_N = (1, \dots, 1)^\top \in \mathbb{R}^N$, and P^{rate} is the average power consumption by the DERs. If the transition probability between bins q_i and q_j in q is denoted as p_{ij} and $M = \{p_{ij}\}_{i,j=1,\dots,3N}$, it then follows that

$$q[k+1] = M q[k], \quad (10)$$

which represents the dynamics of a Markov chain.

The transition rates p_{ij} are computed explicitly by considering how the dynamic state interval corresponding to a particular bin state is altered by the DER hybrid dynamics. Recall that the main factor affecting these transition rates is the background usage of DERs by the end-user as modeled in Section III-A with a generic birth and death process. For EVs, the Markov chain (8) is embedded into M in (10) by setting the transition $x_{sb}^i \in \mathcal{X}_{sb} \rightarrow x_d^i \in \mathcal{X}_d$ to p_1 and $x_d^i \in \mathcal{X}_d \rightarrow x_{sb}^i \in \mathcal{X}_{sb}$ to p_2 for all i . The above assumes that driving patterns are independent of the EV's SoC, however, one can add energy-dependent transitions, $p_1(x_{sb}^i)$ and $p_2(x_d^i)$, based on available data.

Fig. 5 provides a simple illustration of how a (uniformly distributed) probability mass in an arbitrary state shifts after Δt time as a function of a DER's dynamics for the hybrid states, c, sb and d, which in this manuscript follow either (1) or (2) and with respect to the average end-user event of the corresponding population. More specifically and dropping the subscript n , denote the solution of (1) with respect to the hybrid state $h \in \{c, sb, d\}$ and initial condition z_0 at time k by the mappings $\Phi_c^h(k, w) = \Phi_{z_0}(k, w)_{h=c}$, $\Phi_{sb}^h(k, w) = \Phi_{z_0}(k, w)_{h=sb}$ and $\Phi_d^h(k, w) = \Phi_{z_0}(k, w)_{h=d}$. Let W be the average end-user event resulting from the aggregation of processes satisfying

(4). One can show that W is a process comprised of normally distributed random variables having parameters μ_W (mean) and σ_W (standard deviation). In addition, if $P_W(w)$ denotes the occurrence probability of $W = w$ with $w \in [\underline{w}, \bar{w}]^3$, then $P(z = \Phi_{z_0}^h(t, w)) = P_W(w)$. Therefore, one can map any $z_0 \in [z_{i-1}, z_i]$ (interval corresponding to bin i) through $\Phi_{z_0}^h(t, w)$ with $h \in \{c, sb, d\}$ for all $w \in [\underline{w}, \bar{w}]$. The mapping gives as an outcome the dynamic state of the DER after naturally evolving for $k\Delta t$ seconds. Repeating this procedure for all $z_0 \in [z_{i-1}, z_i]$ and normalizing the resulting histogram of dynamic state outcomes gives the t -seconds-ahead distribution of the i th bin. The transition probability from bin i to bin j at a specific hybrid state $h \in \{c, sb, d\}$, p_{ij}^h , is then simply the probability mass that started entirely in bin i and after t seconds now overlaps with bin j . This procedure makes the bin transition probabilities a function of the statistics of the end-user events. The work in [16], [17] lack such feature and assumed that transitions only occurred between contiguous bin intervals. Therefore, there was an apparent mismatch from what the aggregated model produced in relation to what an agent based simulation provided under stressed conditions. Observe that a DER population does not transition to higher states, as expected, when $h = \{sb, d\}$ since no energy is injected into DERs while in these modes. For instance, (1) and (2) are driven by non-negative energy losses or zero-mean bounded damping terms, thus the only way these can increase their SoC is by charging.

C. The State Bin Model for Conventional DER Dynamics

Without coordination schemes, such as PEM, the DERs nominally operate based on (conventional) decentralized control logic that is specific to each DER type. For example conventional TCLs operate under hysteretic control, which is based on keeping the local state variable (e.g., temperature) within a dead-band $[\underline{z}, \bar{z}]$ of width z_{DB} and set point $z_{set} \in [\underline{z}, \bar{z}]$. More precisely, a conventional TCL transitions to the c state only when $z \leq \underline{z}$, transitions to sb from a d state only when $z \leq \underline{z}$, and transitions to sb from a c state when $z \geq \bar{z}$. Clearly, the TCLs' discharging states are unreachable since a TCL cannot actively inject power into the grid. Similarly, conventional EV or ESS control logic can be mapped to a state bin model to consider an EVs' charge-upon-arrival rule and average driving behavior and an ESSs' solar PV net-metering tariff. Thus, the associated Markov transition matrix M for a fleet of DERs in (10) can be described by the following nominal, autonomous (uncoordinated) dynamics:

$$M := \begin{pmatrix} M_c & M_{c,sb} & 0_N \\ M_{sb,c} & M_{sb} & M_{d,sb} \\ 0_N & M_{sb,d} & M_d \end{pmatrix}, \quad (11)$$

where 0_N denotes the N -dimensional zero matrix and M_h , for $h \in \{c, sb, d\}$, is a multi-diagonal matrix containing the probabilities of staying, going to higher energy states, and going to lower energy states. Similarly, $M_{c,sb}$ and $M_{sb,c}$ are responsible for transferring DERs that exceeds \bar{z} from c to sb and any DERs

³In reality, W is lower (upper) bounded by some fixed values \underline{w} (\bar{w}).

that fall below \underline{z} from sb to c, respectively. Finally, $M_{d, sb}$ and $M_{sb, d}$ provide the transition probabilities from d to sb and from sb to d including the probabilities of uncontrollable transition events following some model comparable to the one described for EVs in Section III-A. Observe that, by design, the Markov chain associated with M is irreducible since one can reach any state from any arbitrarily chosen initial state and is aperiodic due to every state having self-loops. It follows then that this abstraction possesses a unique invariant distribution since \mathcal{X} is finite dimensional. Next, we augment the hysteretic control scheme with the probabilistic transitions and opt-out control inherent to PEM as discussed in Section II.

D. The State Bin Model for Packetized Energy Management

Under PEM, a DER can only switch to charging/discharging modes for an epoch if the corresponding charging/discharging packet request is accepted by the coordinator (i.e., VPP). To capture the unique nature of PEM's fixed packet duration and the VPP's role in authorizing/denying packet requests, we leverage prior literature on fault-tolerant recovery logic [27] and TCL modeling with compressor lockout periods [9]. In this subsection, earlier work on modeling PEM in [16] is adapted and extended to consider EVs. PEM coordination can be described as a *controlled Markov chain*.

Definition 1: Let $\{u_k\}_{k \geq 0}$ be a sequence of real valued functions taking values on a set U . A Markov chain $\{X_k\}_{k \geq 0}$ is said to be a *controlled Markov chain* (CMC) if its transition matrix $M(u) := \{q_{ij}(u)\}_{1 \leq i, j \leq N}$ satisfies

$$\begin{aligned} P(X_{n+1} = x_{i_{n+1}} | X_n = x_{i_n}, \dots, X_0 = x_{i_0}, u_n, \dots, u_0) \\ = P(X_{n+1} = x_{i_{n+1}} | X_n = x_{i_n}, u_n) = p_{i_{n+1} i_n}(u_n). \end{aligned}$$

Note that the resulting matrix $M(u)$ must be a (column) stochastic matrix for any choice of $u \in U$. As usual, the probability mass function of a CMC is computed similarly using $q[k+1] = M(u[k])q[k]$ given an initial distribution $q[0]$ and control policy $u(x)[k] : \mathcal{X} \rightarrow U$ for $k = 0, 1, \dots$. The underlying transition matrix over which PEM is implemented is (11), but with $M_{sb, c} = M_{c, sb} = 0_N$ and $M_{sb, d}$ and $M_{d, sb}$ accounting only for uncontrollable transitions. In this section, our model assumes that any DER in the top/bottom states in \mathcal{X}_c , \mathcal{X}_{sb} and \mathcal{X}_d that transitions in the next time step outside of $[\underline{z}, \bar{z}]$ will remain in those top/bottom states. This assumption is relaxed in the next section where the opt-out mechanism is introduced in order to avoid using absorbing states.

Before detailing the PEM request coordination mechanism, consider the following CMC with controlled transition rates $\beta_h = \text{diag}\{\beta_h^1, \dots, \beta_h^N\}$ with $\beta_h^i \in [0, 1]$ and $h \in \{c, d\}$ as the percentage of the standby population in state x_{sb}^i that transitions to charge/discharge and $\beta_{h, sb} = \text{diag}\{\beta_{h, sb}^1, \dots, \beta_{h, sb}^N\}$ with $\beta_{h, sb}^i \in [0, 1]$ and $h \in \{c, d\}$ the percentage of the charging/discharging population in state x_h^i that transitions to standby. The relative transition rates of charging, discharging and standby devices to a different state in q is then given by the transformation:

$$\bar{q}[k] = \bar{M}(\beta[k], \beta_{sb}[k]) q[k], \quad (12)$$

where $\beta := (\beta_c, \beta_d)^\top$, $\beta_{sb} := (\beta_{c, sb}, \beta_{d, sb})^\top$, and

$$\bar{M}(\beta, \beta_{sb}) := \begin{pmatrix} I_N - \beta_{c, sb} & \beta_c & 0_N \\ \beta_{c, sb} & I_N - \beta_c - \beta_d & \beta_{d, sb} \\ 0_N & \beta_d & I_N - \beta_{d, sb} \end{pmatrix}, \quad (13)$$

where I_N denotes the N -dimensional identity matrix. Once $\bar{M}(\beta, \beta_{sb})$ has switched some DERs to a new charge/standby/discharge mode, the matrix M makes the DERs in \bar{q} evolve with the natural dynamics inside each mode of operation. It then follows that

$$q[k+1] = M\bar{q}[k] = M\bar{M}(\beta, \beta_{sb})q[k], \quad (14)$$

which is a CMC as shown by the next theorem.

Theorem 2: Let $\beta[k], \beta_{sb}[k] \in \mathbb{R}^{2N \times N}$ be defined as in (12) $\forall k \geq 0$. The sequence $\{X_k\}_{k \geq 0}$ of random variables X_k taking values in \mathcal{X} and probability distribution satisfying (14) is a controlled Markov chain as described by Definition 1 with input $u[k] = (\mathbf{1}_{2N}^\top \beta[k], \mathbf{1}_{2N}^\top \beta_{sb}[k])^\top \in \mathbb{R}^{4N}$.

Proof: The proof is straightforward since matrices (11) and (13) are stochastic for any choice of β, β_{sb} , and the product of stochastic matrices is a stochastic matrix. ■

The details of PEM model are provided next in a manner that, when applying Theorem 2, it is concluded that the resulting PEM model is a CMC. Unlike the CMC in (14), the PEM scheme is based on charging and discharging requests coming from the standby population as a function of bin state (e.g., based on temperature for TCLs and state-of-charge for ESSs and EVs). Thus, the number of charging and discharging requests are paramount for modeling DERs in PEM. Define

$$q_h^+[k] := T_{\text{req}, h} q_{sb}[k] \quad (15)$$

where $T_{\text{req}, h} = \text{diag}\{p_1^{\text{req}, h}, \dots, p_N^{\text{req}, h}\}$, $p_i^{\text{req}, h} := 1 - e^{-\mu_h(Z_i^m) \Delta t}$ is the request probability assigned to x_{sb}^i by (3) with respect to the mid-point of state bin i and $h = \{c, d\}$. The number of charging/discharging requests received by the VPP is then $n_r^h[k] := \mathbf{1}_N^\top q_h^+[k]$. It is assumed that each bidirectional ESS cannot request to both charge and discharge at the same time. This implies that if both packet types were requested during time-step k , they cancel each other out and no request is made. Therefore, for each individual DER, since a charging and a discharging request occur independently: $T_{\text{req}, c}$ and $T_{\text{req}, d}$ are replaced in (15) by

$$T_{\text{req}, c, \bar{c}} = T_{\text{req}, c}(I_N - T_{\text{req}, d}) \text{ and}$$

$$T_{\text{req}, d, \bar{d}} = T_{\text{req}, d}(I_N - T_{\text{req}, c}),$$

respectively. Thus, under PEM, the VPP determines the proportion of accepted charging/discharging packets ($\beta_c[k]$ for charging and $\beta_d[k]$ for discharging). Upon a packet being accepted by the VPP, the DER transitions to the new state.

Due to the pre-determined duration of packets, the model needs to capture the dynamics of the active and *expiring packets* from the charging and discharging populations, which introduces two sets of timer states. That is, given a packet epoch δ , the sampling time step Δt , and two timer states vectors $x_{p, h} \in \mathbb{R}^{n_p}$ with $n_p = \lfloor \delta / \Delta t \rfloor$ and $h = \{c, d\}$, the *timer dynamics* are given

by

$$x_{p,h}[k+1] = M_{p,h}x_{p,h}[k] + C_{p,h}\beta_h q_h^+[k], \quad (16)$$

where $C_{p,h} \in \mathbb{R}^{n_p \times N}$ is responsible for allocating the new charge/discharge population into their corresponding charge/discharge timer states. In reference to the matrix $C_{p,h}$, for any DER whose packet is accepted, there is a state $z_c(z_d)$ such that $\Phi_{z_c}^c(\delta) = \bar{z}(\Phi_{z_d}^d(\delta) = \underline{z})$. Therefore $C_{p,c}(C_{p,d})$ interrupts packets to prevent exceeding \bar{z} (falling below \underline{z}). That is, if $z_{i+1} < z_c(z_{i+1} > z_d)$, $C_{p,c}(C_{p,d})$ allocates all DERs requesting charging packets from bin $[z_i, z_{i+1}]$ into the timer state $x_{p,c}^1(x_{p,d}^1)$. Otherwise, it allocates the DER with $z_j > z_c(z_j < z_d)$ in the timer state $x_{p,c}^j(x_{p,d}^j)$ with $j = \lfloor (\delta - t_j^c)/\Delta t \rfloor$ ($j = \lfloor (\delta - t_j^d)/\Delta t \rfloor$) and $t_j^c(t_j^d)$ the time that the DER takes to move its state from z_j to $\bar{z}(\underline{z})$ – this captures the PEM concept of *interrupted packets*. The timers provide a formula for the percentage of DERs whose packet expires. That is, $\beta_h := x_{p,h}^{(n_p)} / \sum_{i=1}^{n_p} x_{p,h}^{(i)}$, where $x_{p,h}^{(i)}$ is the i -th component of $x_{p,h}$. Abusing of the notation, consider the particular β and β_{sb} in (14) given as $\beta = (\beta_c T_{\text{req},c,d}, \beta_d T_{\text{req},c,d})^\top$, $\beta_{sb} = (\beta_c^- I_N, \beta_d^- I_N)^\top$ where β_c and β_d are now scalars with values in $[0,1]$. Therefore, a simple algebraic procedure yields the PEM population dynamics and represents a CMC:

$$\begin{aligned} q[k+1] &= M(I + M_{\beta[k]}^+ - M_{\beta_{sb}[k]}^-)q[k] \\ &= \bar{M}(\beta[k], \beta_{sb}[k])q[k], \end{aligned} \quad (17)$$

where

$$\begin{aligned} M_{\beta_{sb}}^- &:= \begin{pmatrix} \beta_c^- I_N & 0_N & 0_N \\ -\beta_c^- I_N & 0_N & -\beta_d^- I_N \\ 0_N & 0_N & \beta_d^- I_N \end{pmatrix}, \\ M_{\beta}^+ &:= \begin{pmatrix} 0_N & \beta_c T_{\text{req},c,d} & 0_N \\ 0_N & -\beta_c T_{\text{req},c,d} - \beta_d T_{\text{req},c,d} & 0_N \\ 0_N & \beta_d T_{\text{req},c,d} & 0_N \end{pmatrix}. \end{aligned}$$

A full schematic diagram for the PEM dynamics is given in Fig. 7.

This section concludes with an illustrative simulation of a population of 2000 EWHs aimed at validating the state bin transition model developed for PEM. The model for the n -th EWH is given by

$$z_n^+ = z_n + \Delta t \left(\frac{P_n^{\text{rate}} z_n}{c\rho L_n \eta} - \frac{z_n - z_{\text{amb}}}{\tau_n} - \frac{z_n - z_{\text{in}}}{60L_n} w_n \right), \quad (18)$$

where the parameters of (18) are provided in Table I and the end-user events, w_n , are PRPs with the same parameters as in Example 1. The simulation initially accepts all charging requests ($\beta_c = 1$) until all incoming requests are denied after minute 120 ($\beta_c = 0$). Fig. 6 shows the result of such experiment and compares the resulting power output of the agent based (running each individual TCL and then aggregating the outcomes) and macro-model simulations. On average the error in power between these simulations never exceeds 8% throughout the entire simulation time of 10 hours. For QoS, the error in SoC was found to be less than 0.1°C , which amount to less

TABLE I
EWH SIMULATION PARAMETERS

Parameter	Value	Unit
Simulation period	600	mins
Sampling period, Δt	15	s
Specific heat capacity (Water), c	4.186	kJ / (kg- $^\circ\text{C}$)
Water density, ρ	0.99	kg / liter
Ambient insulation losses, τ_n	150	hr
Heater Capacity, L_n	250	liters
Set-point temperature, z_n^{set}	52	$^\circ\text{C}$
Dead-band temperature, $z_n^{\text{set,DB}}$	$0.12z_n^{\text{set}}$	$^\circ\text{C}$
PEM temperature bounds, $z_n^{\text{set,PEM}}$	$0.08z_n^{\text{set}}$	$^\circ\text{C}$
PEM request parameter m_R	$\frac{1}{300}$	Hz
Input heat transfer rate, $P_{c,n}^{\text{rate}}$	4.5	kW
Heating efficiency, η	100	%
Ambient temperature, z_{amb}	14	$^\circ\text{C}$
Inlet temperature, z_{in}	14	$^\circ\text{C}$

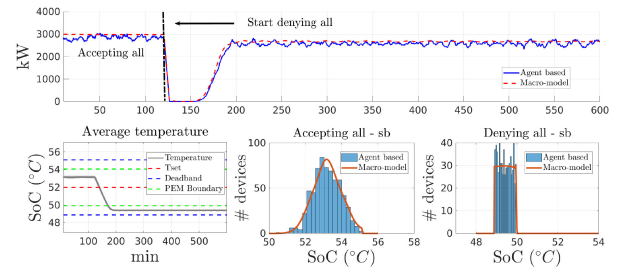


Fig. 6. Comparing macro-model and a realization of a micro-model simulation of 2000 EWHs for a 10-hour accept-all/deny-all VPP experiment. Power response and distributions of the standby populations (accepting and denying) are provided.

than 1% relative error with respect to the agent-based average SoC. Furthermore, the DER standby distributions for accepting all requests (at minute 100) and denying all (at minute 400) for both the agent base and macro-model simulations are presented to illustrate how close the distributions are for both simulations.

IV. QoS GUARANTEES AND DIVERSE DERs

When managing demand, it is critical to be cognizant of end-consumer QoS. For example, when coordinating EWHs, people will opt out *en masse* from water heater DR programs, the first time they experience cold showers. However, before discussing *QoS guarantees* for EWHs, ESSs, and EVs consider the following definition.

Definition 2: A coordinator (or VPP) providing grid services is said to *guarantee QoS* if for a pre-specified SoC range and set-point z_{set} , there exist conditions under which the average SoC of the DER population is greater than or equal to z_{set} .

One way to guarantee QoS is with *opt-out control*, which has been explored in the context of demand dispatch, e.g., see [15], [28], but not for a PEM-based macro-model. The opt-out control mechanism for PEM is described at the beginning of Section II in *ii*). Thus, DERs whose dynamic state are lower than \underline{z} exit PEM (move to charge or exit ON) and join a new set of energy states constituting the *Opt-Out* mode (denoted by \oplus). On the other hand, if the dynamic state is too high, packet interruptions provided by the timer's matrix $C_{p,h}$ for $h = \{c,d\}$ in (16) avoid

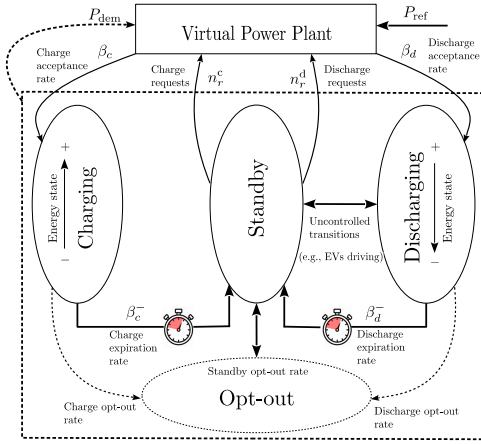


Fig. 7. Transition diagram of a DER population under PEM with opt-out control.

the need for a separate opt out (i.e., exit OFF). Interestingly, adding opt-out operation to the PEM macro-model only requires a simple augmentation of states with their corresponding transition rates as shown in Fig. 7. That is, q is redefined as $q^\top = (q_\oplus^\top, q^\top)$ with

$$q[k+1] = M_{\text{exit}}(I + M_{\beta[k]}^+ - M_{\beta_{\text{sb}}[k]}^-)q[k] \text{ and } y[k] = cq[k], \quad (19)$$

where $I + M_{\beta}^+ - M_{\beta_{\text{sb}}}^-$ adds a diagonal block identity matrix and uses zeros elsewhere since q_{opt} are unaffected by β , β_{sb} . Note that

$$M_{\text{exit}} := \begin{pmatrix} M_{\text{pem}}^\oplus & M_{\text{pem}}^\ominus \\ M_{\text{pem}}^\oplus & M \end{pmatrix},$$

where M^\oplus is a sub-matrix of M that has all rows and columns corresponding to states higher than the pre-specified PEM re-entry bound removed. Finally, M_{pem}^\ominus (M_{pem}^\oplus) provides the transition probabilities of exiting (re-entering) PEM. A depiction of the transition diagram for a DER population under PEM with opt-out control is provided in Fig. 7.

A. QoS for EWHs and ESSs

For this type of DERs, Definition 2 implies that there must exist some β such that $cq(\beta) \geq z_{\text{set}}$, where $q(\beta)$ is the invariant distribution associated to β . Naturally by making $\beta = (\beta_c, \beta_d) = (1, 0)$ for all times, the system reaches its maximum average dynamic state, and the system is in equilibrium, which also fixes β_{sb} . Since end-user events for EWHs and ESS do not alter their hybrid state, QoS guarantees are provided by adding the opt-out dynamics as in (19).

B. QoS for EVs

The end-user events for EVs (i.e., driving) do change the hybrid state. This renders some EVs unavailable to PEM, which couples QoS guarantees to the (average) EV driving model. Next, we employ the simple driving model introduced in Section III-A and formulate a condition under which EVs can

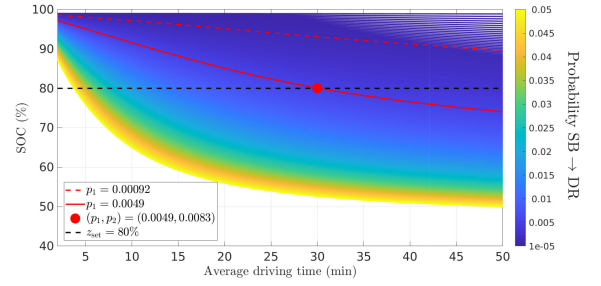


Fig. 8. Average SOC for EVs as a function of driving time and the probability of going from standby to a driving state (SB \rightarrow DR). The solid and dashed red lines indicate two level sets for constant p_1 . The red dot indicates the value for p_1 and p_2 for the maximum occupancy of driving states for which the fleet's QoS is guaranteed to reach 80% SOC.

also guarantee QoS. The condition will be in terms of the average drive time, which is related to p_2 and the probability of going from driving (discharge) to standby.

Fig. 8 shows the relationship between the average state of charge of a population of EVs as a function of departure rate p_2 and the probability of going from standby to driving (equal to p_1), where one can see that when the departure rate surpasses a threshold for an specific p_1 guaranteeing QoS is not possible. Recall from the discussion about the driving model that the probabilities p_1 and p_2 are independent of each other and that they can be chosen so that they follow data such as that from NHTS [20] (see Fig. 4b). For a fixed driving state occupancy, it is possible to compute a bound for the maximum average driving time that a fleet of EVs should have so that QoS is guaranteed, which amounts to a bound on p_2 . Recall that by fixing the driving states occupancy and p_2 , p_1 is automatically fixed. Moreover, setting $\beta = (0, 1)$ and assuming that cars return to standby from driving independent of their energy state, the Markov transition matrix for EVs has the form $A + p_2B$, where A is an irreducible and aperiodic column stochastic matrix and B is such that its columns add to zero. The invariant distribution of the evolution equation $q[k+1] = (A + p_2B)q[k]$, for a fixed p_2 , is computed by solving $\tilde{A}q^* = \tilde{b}$ with

$$\tilde{A} := \begin{pmatrix} I - \frac{(A + p_2B)}{\mathbf{1}_N^\top} \\ \mathbf{1}_N^\top \end{pmatrix} q^* \text{ and } \tilde{b} := \begin{pmatrix} 0_{N \times 1} \\ 1 \end{pmatrix}.$$

This system of algebraic equations has $3N + 1$ equations with $3N$ unknowns, where one equation is redundant due to the fact that the dimension of the nullity of $I - (A + p_2B)$ is one. Therefore, a least square procedure provides a unique solution q^* that is the desired stationary distribution. Specifically,

$$\begin{aligned} q^* &= (\tilde{A}^\top \tilde{A})^{-1} \tilde{A}^\top \tilde{b} \\ &= ((I - (A + p_2B))^\top (I - (A + p_2B)) + \mathbf{1}_N \mathbf{1}_N^\top)^{-1} \mathbf{1}_N. \end{aligned}$$

Since $p_2 \ll 1$ for any practical scenario, then

$$q^* \approx Q_1^{-1} \mathbf{1}_N - p_2 Q_1^{-1} Q_2 Q_1^{-1} \mathbf{1}_N, \quad (20)$$

where $Q_1 = \mathbf{1}_N \mathbf{1}_N^\top + I_N - (A + A^\top - A^\top A)$ and $Q_2 = A^\top B + B^\top A - B - B^\top$. The condition that must be satisfied for guaranteeing QoS is $z_{\text{set}} \leq cq^*$, where q^* is related to p_2

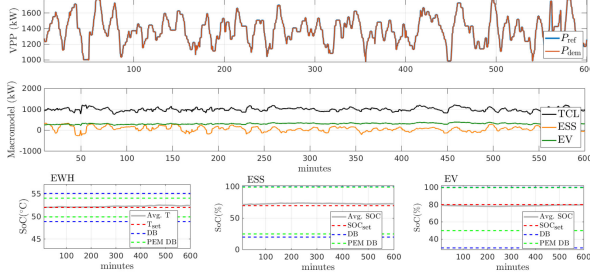


Fig. 9. (Top) The result of a fleet comprised of 1000 EWHs, 1000 ESSs and 250 EVs tracking a regulation reference signal (Middle) Individual contribution of EWHs, ESSs and EVs to balance the regulation signal (Bottom) The average SoC for EWHs ($^{\circ}\text{C}$), ESSs (%) and EVs (%).

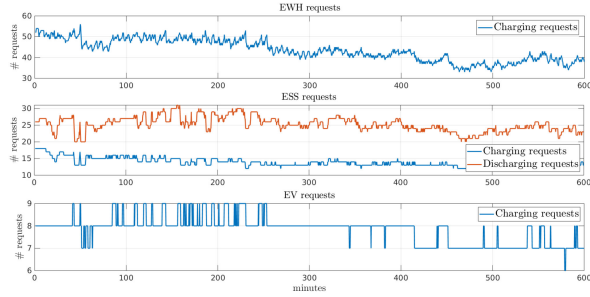


Fig. 10. Number of charging and discharging requests from the EWH, ESS and EV standby populations.

via (20) and the dependence on β is omitted given that it was previously fixed. Thus, a bound for p_2 is given by:

$$p_2 \leq \frac{c Q_1^{-1} \mathbf{1}_N - z_{\text{set}}}{c Q_1^{-1} Q_2 Q_1^{-1} \mathbf{1}_N}$$

Observe that $Q_1 = \tilde{A}^{\top} \tilde{A}|_{p_2=0}$ is always invertible because $\tilde{A}|_{p_2=0}$ has full column rank. From the discussion at the end of Section III-A, an occupancy of the driving states of $\pi_{\text{zDR}} = 0.1$ with $p_1 = 0.00092$ gives the exact threshold $p_2 < 0.00053$, whereas the approximation yields $p_2^{\text{approx}} < 0.00051$. In other words, the exact calculation says that one can guarantee QoS when the average drive is less than 471.72 minutes (for $\Delta t = 15$ sec), and the approximation gives 488.8 minutes as the driving time threshold. Fig. 8 shows the curve corresponding to the parameters above as well as the curve that intersects the set point at 30 min driving time. These parameters are in agreement with a fleet of EVs at off-peak driving hours of the day which is when EVs become a real flexible resource and are well beyond the 30 min average driving time assumed for the simulations in this manuscript. For instance, the driving state occupancy for guaranteeing QoS with 80% of average charge (accepting all charging requests and rejecting all discharging requests) and 30 min average driving time is approximately 37% (see Fig. 8).

C. Illustrative Simulation With Tracking and QoS Awareness

A fleet of 1000 EWHs, 1000 ESSs, and 250 EVs are modeled using the macro-model developed in Section III-B and presented in Figs. 9 and 10. The EWHs for this simulation have the

same parameters shown in Table I. The ESS models here are representative of Tesla's PowerWalls (2.0), which have battery capacity of 13.5 kWh, charge and discharge efficiency of around 95% (roundtrip of 92%), and a maximum (continuous) power rating ($P_{c,n}^{\text{rate}} = P_{d,n}^{\text{rate}}$) of 5.0kW. It is assumed that the battery owner charges or discharges the battery based on a Gaussian random walk with a minimum power draw of 1.5 kW in either direction. This could be representative of excess or deficit solar PV production. EVs, on the other hand, are assumed to have an electric driving range of 150 miles and an electric driving efficiency of 7 miles-per-kWh. The PEM system has the task to track a detrended and scaled regulation signal [29]. The most important observation is that under the conditions for guaranteeing QoS one can construct a rule for acceptance such that the three populations work together to balance their output power with respect to the given reference in a manner that average SoC is close to the predefined set points. The specifics of this tracking problem and how the populations work in tandem are detailed in [30]. In addition, Fig. 10 shows the number of requests from the three populations as a function of time. The VPP chooses a percentage of these charge and discharge packet requests in order to balance the regulation signal provided by the system operator as a reference. Note also that the QoS for each population is maintained around its predefined set point even though the populations are providing power balancing dynamically (i.e., without predictive optimization). As can be seen in Fig. 9, the EWHs and EVs effectively provide the bias while the ESS provide the corrective (together with EWHs) for tracking the regulation signal. Furthermore, to achieve desired tracking of the reference regulation signal, the ESS population's average SOC deviates slightly ($< 5\%$) from the desired set-point, which increases the number of ESS discharging requests. The internal feedback offered by the population's packet request mechanism drives the availability upward/downward flexibility. In this case, the ESS population alone can offer more downward flexibility (discharge) than upward flexibility (charge) when the VPP receives more discharge requests than charge requests. Of course, if the reference signal was biased downward, the populations would have to deviate from their SoC set-point to achieve satisfactory tracking performance, which would effectively discharge the populations over time and lead to an increase in the number of charge requests. After discharging for sufficiently long, the opt-out mechanism built into PEM would override the request-response mechanism and devices would opt-out and tracking performance would be negatively affected. The coupling between discharge/charge duration and tracking performance is the subject of ongoing work and has led to development of improved PEM-VPP controller designs [30] and energy-based modeling of PEM population to capture the battery-like, energy-power relations.

V. CONCLUSION AND FUTURE WORK

This manuscript presented a macro-model for the aggregation of a system comprised by DERs. The approach was based on a bottom-up DER coordination methodology called PEM. The macro-model was described as a controlled Markov chain

that included the mechanics of accepting, active, and expiring packets with the help of two timers to differentiate charging and discharging packet requests. Finally, QoS guarantees were given for TCLs and ESS with an opt-out mechanism while QoS guarantees for EVs were provided in terms of EVs' average arrival and departure rates.

Future work involves addressing heterogeneity of the macro-model either by clustering or by a set-based Markov model. Moreover, the dependence of end-user event rates and magnitudes on opt-out conditions is currently being explored by the authors to study time-of-day changes in dispatchable demand. Finally, incorporating live grid conditions into the PEM macro-model is of interest to grid operators and aggregators.

ACKNOWLEDGMENT

M. Almassalkhi is co-founder of startup Packetized Energy, which is actively commercializing packetized energy management.

REFERENCES

- [1] F. C. Schweppe, R. D. Tabors, J. L. Kirtley, H. R. Outhred, F. H. Pickel, and A. J. Cox, "Homeostatic utility control," *IEEE Trans. Power App. Syst.*, vol. PAS-99, no. 3, pp. 1151–1163, May 1980.
- [2] D. S. Callaway and I. A. Hiskens, "Achieving controllability of electric loads," *Proc. IEEE*, vol. 99, no. 1, pp. 184–199, Jan. 2011.
- [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] S. P. Meyn, P. Barooah, A. Bušić, Y. Chen, and J. Ehren, "Ancillary service to the grid using intelligent deferrable loads," *IEEE Trans. Autom. Control*, vol. 60, no. 11, pp. 2847–2862, Nov. 2015.
- [5] M. Vrakopoulou, J. Mathieu, and G. Andersson, "Stochastic optimal power flow with uncertain reserves from demand response," in *Proc. IEEE 47th Hawaii Int. Conf. Syst. Sci.*, 2014, pp. 2353–2362.
- [6] G. S. Ledva and J. Mathieu, "A linear approach to manage input delays while supplying frequency regulation using residential loads," in *Proc. Amer. Control Conf.*, May 2017, pp. 1–7.
- [7] J. L. Mathieu, M. Kamgarpour, J. Lygeros, and D. S. Callaway, "Energy arbitrage with thermostatically controlled loads," in *Proc. IEEE Eur. Conf. Circuit Theory Des.*, Jul. 2013, pp. 2519–2526.
- [8] W. Zhang, J. Lian, C. Y. Chang, K. Kalsi, and Y. Sun, "Reduced order modeling of aggregated thermostatic loads with demand response," in *Proc. 51st IEEE Conf. Decis. Control*, 2012, pp. 5592–5597.
- [9] 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.
- [10] S. Soudjani and A. Abate, "Aggregation and control of populations of thermostatically controlled loads by formal abstractions," *IEEE Trans. Control Syst. Technol.*, vol. 23, no. 3, pp. 975–990, May 2015.
- [11] A. Brooks, E. Lu, D. Reicher, C. Spirakis, and B. Wehl, "Demand dispatch," *IEEE Power Energy Mag.*, vol. 8, no. 3, pp. 20–29, May 2010.
- [12] Y. Chen, A. Bušić, and S. Meyn, "Estimation and control of quality of service in demand dispatch," *IEEE Trans. Smart Grid*, vol. 9, no. 5, pp. 5348–5356, May 2018.
- [13] Y. Chen, A. Bušić, and S. Meyn, "State estimation for the individual and population in mean field control with application to demand dispatch," *IEEE Trans. Autom. Control*, vol. 62, no. 3, pp. 1138–1149, Mar. 2017.
- [14] Y. Chen, M. U. Hashmi, J. Mathias, A. Bušić, and S. Meyn, *Distributed Control Design for Balancing the Grid Using Flexible Loads*. New York, NY, USA: Springer, 2018, pp. 383–411.
- [15] M. Almassalkhi, J. Frolik, and P. Hines, "Packetized energy management: Asynchronous and anonymous coordination of thermostatically controlled loads," in *Proc. Amer. Control Conf.*, May 2017, pp. 1431–1437.
- [16] L. Duffaut Espinosa, M. Almassalkhi, P. Hines, and J. Frolik, "Aggregate modeling and coordination of diverse energy resources under packetized energy management," in *Proc. 56th IEEE Conf. Decis. Control*, Dec. 2017, pp. 1394–1400.
- [17] L. Duffaut Espinosa, M. Almassalkhi, P. Hines, and J. Frolik, "System properties of packetized energy management for aggregated diverse resources," *Power Syst. Comput. Conf.*, pp. 1–7, Jun. 2018.
- [18] P. Rezaei, J. Frolik, and P. D. H. Hines, "Packetized plug-in electric vehicle charge management," *IEEE Trans. Smart Grid*, vol. 5, no. 2, pp. 642–650, Mar. 2014.
- [19] B. Zhang and J. Baillieul, "Control and communication protocols based on packetized direct load control in smart building microgrids," *Proc. IEEE*, vol. 104, no. 4, pp. 837–857, Apr. 2016.
- [20] Federal Highway Administration. National Household Travel Survey, U.S. Department of Transportation, Washington, DC, 2017. [Online]. Available: <https://nhts.ornl.gov>
- [21] S. G. Buchberger and L. Wu, "Model for instantaneous residential water demands," *J. Hydraul. Eng.*, vol. 121, pp. 232–246, 1995.
- [22] P. E. Protter, *Stochastic Integration and Differential Equations*. Berlin, Germany: Springer, 2005.
- [23] M. Huang, P. E. Caines, and R. P. Malhame, "Individual and mass behaviour in large population stochastic wireless power control problems: Centralized and Nash equilibrium solutions," in *Proc. 42nd IEEE Conf. Decis. Control*, Dec. 2003, vol. 1, pp. 98–103.
- [24] M. Nazir and I. Hiskens, "Noise and parameter heterogeneity in aggregate models of thermostatically controlled loads," *IFAC-PapersOnLine*, vol. 50, no. 1, pp. 8888–8894, 2017.
- [25] S. Kundu, N. Sinityn, S. Backhaus, and I. Hiskens, "Modeling and control of tcls," in *Proc. Power Syst. Comput. Conf.*, Jun. 2011, pp. 1–7.
- [26] L. Duffaut Espinosa, M. Almassalkhi, P. Hines, S. Heydari, and J. Frolik, "Towards a macromodel for packetized energy management of resistive water heaters," in *Proc. IEEE Conf. Inf. Sci. Syst.*, Mar. 2017, pp. 1–7.
- [27] H. Zhang, W. S. Gray, and O. R. Gonzalez, "Performance analysis of digital flight control systems with rollback error recovery subject to simulated neutron-induced upsets," *IEEE Trans. Control Syst. Technol.*, vol. 16, no. 1, pp. 46–59, Jan. 2008.
- [28] Y. Chen, A. Bušić, and S. Meyn, "Estimation and control of quality of service in demand dispatch," *IEEE Trans. Smart Grid*, vol. 9, no. 5, pp. 5348–5356, Sep. 2018.
- [29] ISO New England, "Simulated Automatic Generator Control (AGC) Setpoint Data: Energy Neutral," [Online]. Available: <https://www.iso-ne.com/isoexpress/web/reports/grid/-/tree/simulated-agc>, Accessed on: 2017-10-18.
- [30] L. Duffaut Espinosa, M. Almassalkhi, and A. Khurram, "Reference-tracking control policies for packetized coordination of diverse DER populations," *IEEE Trans. Control Syst. Technol.*, to be published.



research interests include modeling, control and estimation of nonlinear systems, stochastic processes, and algebraic combinatorics.



DERs, energy optimization in power systems, and multienergy systems.

Luis A. Duffaut Espinosa (Member, IEEE) received the B.S. degree in physics from the Universidad Nacional de Ingeniería, Lima, Perú, in 2003, the M.S. degree in mathematics with mention in stochastic processes from Pontificia Universidad Católica del Perú, Lima, Perú, in 2005, and the Ph.D. degree in electrical and computer engineering from Old Dominion University, Norfolk, VA, USA, in 2009. He is currently an Assistant Professor with the Department of Electrical and Biomedical Engineering, the University of Vermont, Burlington, VT, USA. His

Mads Almassalkhi (Senior Member, IEEE) received the B.S. degree in electrical engineering with a dual major in applied mathematics from the University of Cincinnati, Cincinnati, OH, USA, in 2008, and the M.S. degree in electrical engineering and systems and the Ph.D. degree from the University of Michigan, Ann Arbor, MI, USA, in 2010 and 2013, respectively. He is currently an Assistant Professor with the Department of Electrical and Biomedical Engineering, the University of Vermont, Burlington, VT, USA. His research interests include multitime-scale control of