

PORA: Predictive Offloading and Resource Allocation in Dynamic Fog Computing Systems

Xin Gao, *Student Member, IEEE*, Xi Huang, *Student Member, IEEE*, Simeng Bian, *Student Member, IEEE*, Ziyu Shao*, *Member, IEEE*, Yang Yang, *Fellow, IEEE*

Abstract—In multi-tiered fog computing systems, to accelerate the processing of computation-intensive tasks for real-time IoT applications, resource-limited IoT devices can offload part of their workloads to nearby fog nodes, whereafter such workloads may be offloaded to upper-tier fog nodes with greater computation capacities. Such hierarchical offloading, though promising to shorten processing latencies, may also induce excessive power consumptions and latencies for wireless transmissions. With the temporal variation of various system dynamics, such a trade-off makes it rather challenging to conduct effective and online offloading decision making. Meanwhile, the fundamental benefits of predictive offloading to fog computing systems still remain unexplored. In this paper, we focus on the problem of dynamic offloading and resource allocation with traffic prediction in multi-tiered fog computing systems. By formulating the problem as a stochastic network optimization problem, we aim to minimize the time-average power consumptions with stability guarantee for all queues in the system. We exploit unique problem structures and propose PORA, an efficient and distributed predictive offloading and resource allocation scheme for multi-tiered fog computing systems. Our theoretical analysis and simulation results show that PORA incurs near-optimal power consumptions with queue stability guarantee. Furthermore, PORA requires only mild-value of predictive information to achieve a notable latency reduction, even with prediction errors.

Index Terms—Internet of Things, fog computing, workload offloading, resource allocation, Lyapunov optimization, predictive offloading.

I. INTRODUCTION

In the face of the proliferation of real-time IoT applications, fog computing has come as a promising complement to cloud computing by extending cloud to the edge of the network to meet the stringent latency requirements and intensive computation demands of such applications [1].

A typical fog computing system consists of a set of geographically distributed fog nodes which are deployed at the network periphery with elastic resource provisioning such as storage, computation, and network bandwidth [2]. Depending on their distance to IoT devices, fog nodes are often organized in a hierarchical fashion, with each layer as a *fog tier*. In such a way, resource-limited IoT devices, when heavily loaded, can delegate workloads via wireless links to nearby fog nodes, *a.k.a.*, *workload offloading*, to reduce the power consumption and accelerate workload processing; meanwhile, each fog node can offload workloads to nodes in its upper fog tier. However, along with all the benefits come the extended latency and extra power consumption. Given such a power-latency tradeoff, two interesting questions arise. One is *where* and *how much workloads* to offload between successive fog tiers. The other

is how to *allocate resources* for workload processing and offloading. The timely decision making regarding these two questions is critical but challenging, due to temporal variations of system dynamics in wireless environment, uncertainty in the resulting offloading latency, and the unknown traffic statistics.

We summarize the main challenges of dynamic offloading and resource allocation in fog computing as follows:

- ◊ **Characterization of system dynamics and the power-latency tradeoff:** In practice, a fog system often consists of multiple tiers, with complex interplays between fog tiers and the cloud, not to mention the constantly varying dynamics and intertwined power-latency tradeoffs therein. A model that accurately characterizes the system and tradeoffs is the key to the fundamental understanding of the design space.
 - ◊ **Efficient online decision making:** The decision making must be computationally efficient, so as to minimize the overheads. The difficulties often come from the uncertainties of traffic statistics, online nature of workload arrivals, and intrinsic complexity of the problem.
 - ◊ **Understanding the benefits of predictive offloading:** One natural extension to online decision making is to employ predictive offloading to further reduce latencies and improve quality of service. For example, Netflix preloads videos onto users' devices based on user behavior prediction [3]. Despite the wide applications of such approaches, the fundamental limits of predictive offloading in fog computing still remain unknown.
- In this paper, we focus on the workload offloading problem for multi-tiered fog systems. We address the above challenges by developing a fine-grained queueing model that accurately depicts such systems and proposing an efficient online scheme that proceeds the offloading on a time-slot basis. To the best of our knowledge, we are the first to conduct systematic study on predictive offloading in fog systems. Our key results and main contributions are summarized as follows:
- ◊ **Problem Formulation:** We formulate the problem of dynamic offloading and resource allocation as a stochastic optimization problem, aiming at minimizing the long-term time-average expectation of total power consumptions of fog tiers with queue stability guarantee.
 - ◊ **Algorithm Design:** Through a non-trivial transformation, we decouple the problem into a series of subproblems over time slots. By exploiting their unique structures, we propose PORA, an efficient scheme that exploits predictive scheduling to make decisions in an online manner.
 - ◊ **Theoretical Analysis and Experimental Verification:** We conduct theoretical analysis and trace-driven simulations to evaluate the effectiveness of PORA. The

X. Gao, X. Huang, S. Bian, Z. Shao and Y. Yang are with the School of Information Science and Technology, ShanghaiTech University, Shanghai 201210, China. (E-mail: {gaoxin, huangxi, biansm, shaozy, yangyang}@shanghaitech.edu.cn) (*Corresponding author: Ziyu Shao)

TABLE I: Comparisons of related works

	D2D-enabled IoT	IoT-Fog ¹	Fog-Fog ²	Fog-Cloud ³	Dynamic	Prior Arrival Distribution	Prediction
[1]		✓	✓	✓		–	
[4]		✓		✓		–	
[5]		✓		✓	✓	Poisson	
[6]		✓				–	
[7]		✓			✓	Poisson	
[8]		✓			✓	Not Required	
[9]		✓			✓	Not Required	
[10]	✓	✓			✓	Not Required	
[11]		✓			✓	Not Required	
Ours			✓	✓	✓	Not Required	✓

^{1,2,3} “IoT-Fog” means offloading from IoT devices to fog, “Fog-Fog” means offloading between fog tiers, while “Fog-Cloud” means offloading from fog to cloud.

results show that PORA achieves a tunable power-latency tradeoff while effectively reducing the average latency with only mild-value of predictive information, even in the presence of prediction errors.

- ◇ **New Degree of Freedom in the Design of Fog Computing Systems:** We systematically investigate the fundamental benefits of predictive offloading in fog computing systems, with both theoretical analysis and numerical evaluations.

We organize the rest of the paper as follows. Section II discusses the related work. Next, in Section III, we provide an example that motivates our design for dynamic offloading and resource consumption in fog computing systems. Section IV presents the system model and problem formulation, followed by the algorithm design of PORA and performance analysis in Section V. Section VI analyzes the results from trace-driven simulations, while Section VII concludes the paper.

II. RELATED WORK

In recent years, a series of works have been proposed to optimize the performance fog computing systems from various aspects [1], [4]–[13]. Among such works, the most related are those focusing on the design of effective offloading schemes. For example, by adopting alternating direction method of multipliers (ADMM) methods, Xiao *et al.* [1] and Wang *et al.* [4] proposed two offloading schemes for cloud-aided fog computing systems to minimize average task duration and average service response time under different energy constraints, respectively. Later, Liu *et al.* [5] took the social relationships among IoT users into consideration and developed a socially aware offloading scheme by advocating game theoretic approaches. Misra *et al.* [6] studied the problem in software-defined fog computing systems and proposed a greedy heuristic scheme to conduct multi-hop task offloading with offloading path selection. Lei *et al.* [7] considered the joint minimization of delay and power consumption over all IoT devices; they formulated the problem under the settings of continuous-time Markov decision process and solved it via approximate dynamic programming techniques. The above works, despite their effectiveness, generally assume the availability of the statistical information on task arrivals in the systems which is usually unattainable in practice with highly time-varying system dynamics [14].

In the face of such uncertainties, a number of works have applied stochastic optimization methods such as Lyapunov optimization techniques to online and dynamic offloading scheme design [8]–[11]. For instance, Mao *et al.* [8] investigated

the tradeoff between the power consumption and execution delay, then developed a dynamic offloading scheme for energy-harvesting-enabled IoT devices. Chen *et al.* [9] designed an adaptive and efficient offloading scheme to minimize the transmission energy consumption with queueing latency guarantee. Gao *et al.* [10] investigated efficient offloading and social-awareness-aided network resource allocation for device-to-device-enabled (D2D-enabled) IoT users. Zhang *et al.* [11] designed an online rewards-optimal scheme for the computation offloading of energy harvesting-enabled IoT devices based on Lyapunov optimization and Vickrey-Clarke-Groves auction. Different from such works that focus on fog computing systems with flat or two-tiered architectures, our solution is applicable to general multi-tiered fog computing systems with time-varying wireless channel states and unknown traffic statistics. Moreover, to the best of our knowledge, our solution is also the first to proactively leverage the predicted traffic information to optimize the system performance with theoretical guarantee. We are also the first to investigate the fundamental benefits of predictive offloading in fog computing systems. We compare our work with the above mentioned works in TABLE I.

III. MOTIVATING EXAMPLE

In this section, we provide a motivating example to show the potential power-latency tradeoff in multi-tiered fog computing systems. The objective is to achieve low power consumptions and short average workload latency (in the unit of packets).

Figure 1 shows an instance of time-slotted fog computing system with two fog tiers, *i.e.*, edge fog tier and central fog tier. Within each fog tier resides one fog node, *i.e.*, an edge fog node (EFN) in edge fog tier and a central fog node (CFN) in central fog tier. The EFN connects to the CFN via a wireless link, while the CFN connects to the cloud data center over wired links. Each fog node maintains one queue to store packets. Figure 1(a) shows that during time slot t_0 , both the EFN and the CFN store 8 packets in their queues.

We assume that each fog node sticks to one policy all the time to handle packets, *i.e.*, either *processing packets locally* or *offloading them to its next tier*. The local processing capacities of EFN and CFN are 1 and 8 packets per time slot, respectively. The transmission capacities from EFN to CFN and from CFN to cloud are 4 and 5 packets per time slot, respectively. The power consumption is assumed linearly proportional to the number of processed/transmitted packets. In particular, processing one packet locally consumes 1 mW power, while transmitting one packet over wireless link

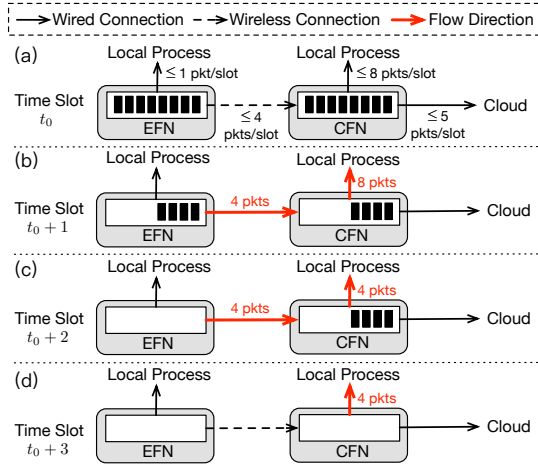


Fig. 1: Motivating example of dynamic offloading and resource consumption in multi-tiered fog computing systems.

consumes 0.5 mW. We ignore the processing latency in the cloud due to its powerful processing capacity.

TABLE II lists the total power consumptions and average packet latencies under all four possible settings. Figures 1(b)-1(d) show the case when EFN sticks to offloading and CFN sticks to local processing. In time slot ($t_0 + 1$), EFN offloads four packets to CFN at its full transmission capacity, while CFN processes all the eight packets locally. In time slot ($t_0 + 2$), EFN offloads the rest four packets to CFN; meanwhile, CFN locally processes the four packets that arrive in previous time slot. In time slot ($t_0 + 3$), CFN finishes processing the rest four packets. In this case, the system consumes 16 mW power in local processing and 4 mW power in transmission, with an average packet latency of 1.75 time slots.

TABLE II: Performance under different offloading policies

Policy of EFN	Policy of CFN	Total Power Consumptions (mW)	Average Packet Latency (time slot)
Local	Local	16	2.75
Local	Offload	8	2.9375
Offload	Local	20	1.75
Offload	Offload	4	2.125

From TABLE II, we conclude that: First, when EFN sticks to offloading and CFN sticks to local processing, the system achieves the lowest average packet latency of 1.75 slots but the maximum power consumption of 20mW. Second, with the same offloading policy on EFN, there is a tradeoff between the total power consumptions and the average packet latency when CFN sticks to different policies. The reason is that offloading to the cloud can not only reduce power consumptions but also prolong latency as well. Third, when CFN sticks to local processing, there is a power-latency tradeoff with different policies at EFN, in that offloading to CFN can induce lower processing latency but at the cost of even higher power consumption for wireless transmissions.

IV. MODEL AND PROBLEM FORMULATION

We consider a multi-tiered fog computing system, as shown in Figure 2. The system evolves over time slots indexed by

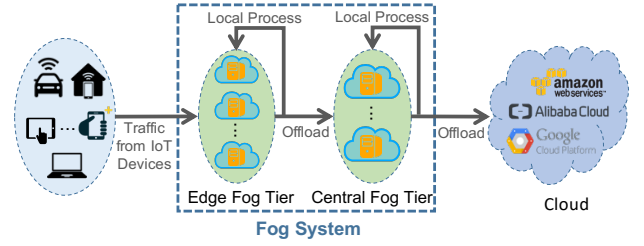


Fig. 2: An example of fog computing systems with two fog tiers.

$t \in \{0, 1, 2, \dots\}$. Each time slot has a length of τ_0 . Inside the edge fog tier (EFT) are a set of edge fog nodes (EFNs) that offer low-latency access to IoT devices. On the other hand, the central fog tier (CFT) comprises of central fog nodes (CFNs) with greater processing capacities than EFNs. We assume that the workload on each EFN can be offloaded to and processed by any of its accessible CFNs, and that each CFN can offload its workload to the cloud. In our model, we do not consider the power consumptions and latencies within the cloud. We mainly focus on the power consumptions and latencies within fog tiers, as shown in TABLE III. First, the power consumptions we consider include two parts: processing power and transmit power. The processing power consumption is induced by the workload processing on both EFT and CFT. The transmit power is induced by the transmissions from EFT to CFT. We do not consider the transmit power consumption from CFT to cloud because we assume that the CFT communicates with the cloud through wireline connections. Second, the latencies we consider include three parts: queuing latency, processing latency and transmit latency. We focus on the queuing latency on both EFT and CFT. We assume that the workload processing in each time slot can be completed by the end of the same time slot, and then we can ignore the processing latency. Since the EFT communicates with the CFT through high-speed wireless connections and the CFT communicates with the cloud through high-speed wireline connections, we assume that transmission latencies from both EFT to CFT and CFT to Cloud are negligible.

TABLE III: Performance Metrics in Our Model

	Power Consumption		Latency		
	Processing	Transmit	Queuing	Processing	Transmit
EFT	✓		✓		
EFT2CFT		✓			
CFT	✓		✓		
CFT2Cloud					

In the following, we first introduce the basic settings in Section IV-A, then elaborate the queuing models in Section IV-D. Next, we define the optimization objective in Section IV-E while proposing the problem formulation in Section IV-F. We summarize the key notations in TABLE IV.

A. Basic Settings

The fog computing system consists of N EFNs in EFT and M CFNs in CFT. Let \mathcal{N} and \mathcal{M} be the sets of EFNs and CFNs. Each EFN i has access to a subset of CFNs in their proximities. We denote the subset by $\mathcal{M}_i \subset \mathcal{M}$. For each CFN j , $\mathcal{N}_j \subset \mathcal{N}$

TABLE IV: Key notations

Notation	Description
τ_0	Length of each time slot
\mathcal{N}	\mathcal{N} is the set of EFNs with $ \mathcal{N} \triangleq N$
\mathcal{M}	\mathcal{M} is the set of CFNs with $ \mathcal{M} \triangleq M$
\mathcal{N}_j	Set of accessible EFNs from CFN j
\mathcal{M}_i	Set of accessible CFNs from EFN i
$A_i(t)$	Amount of workload arriving to EFN i in time slot t
λ_i	Average workload arriving rate on EFN i , $\lambda_i \triangleq \mathbb{E}\{A_i(t)\}$
W_i	Prediction window size of EFN i
$A_{i,-1}(t)$	Arrival queue backlog of EFN i in time slot t
$A_{i,w}(t)$	Prediction queue backlog of EFN i in time slot t , such that $0 \leq w \leq W_i - 1$
$Q_i^{(e,a)}(t)$	Integrate queue backlog of EFN i in time slot t
$Q_i^{(e,l)}(t)$	Local processing queue backlog of EFN i in time slot t
$Q_i^{(e,o)}(t)$	Offloading queue backlog of EFN i in time slot t
$b_i^{(e,l)}(t)$	Amount of workload to be sent to $Q_i^{(e,l)}(t)$ in time slot t
$b_i^{(e,o)}(t)$	Amount of workload to be sent to $Q_i^{(e,o)}(t)$ in time slot t
$f_i^{(e)}(t)$	CPU frequency of EFN i in time slot t
$H_{i,j}(t)$	Wireless channel gain between EFN i and CFN j
$p_{i,j}(t)$	Transmit power from EFN i to CFN j in time slot t
$R_{i,j}(t)$	Transmit rate from EFN i to CFN j in time slot t
$Q_j^{(c,a)}(t)$	Arrival queue backlog of CFN j in time slot t
$Q_j^{(c,l)}(t)$	Local processing queue backlog of CFN j in time slot t
$Q_j^{(c,o)}(t)$	Offloading queue backlog of CFN j in time slot t
$b_j^{(c,l)}(t)$	Amount of workload to be sent to $Q_j^{(c,l)}(t)$ in time slot t
$b_j^{(c,o)}(t)$	Amount of workload to be sent to $Q_j^{(c,o)}(t)$ in time slot t
$f_j^{(c)}(t)$	CPU frequency of CFN j in time slot t
$P(t)$	Total power consumptions in time slot t

denotes the set of its accessible EFNs. Accordingly, for any $i \in \mathcal{N}_j$ we have $j \in \mathcal{M}_i$.

B. Queueing Model for Edge Fog Node

During time slot t , there is an amount $A_i(t)$ ($\leq A_{\max}$ for some constant A_{\max}) of workload generated from IoT devices arrive to be processed on EFN i such that $\mathbb{E}\{A_i(t)\} = \lambda_i$. We assume that such arrivals are independent over time slots and different EFNs. Each EFN i is equipped with a learning module¹ that can predict the future workload within a *prediction window* of size W_i , i.e. workload will arrive in the next W_i time slots. The predicted arrivals are pre-generated and recorded, then arrive to EFN i for pre-serving. Once the predicted arrivals actually arrive after pre-serving, they will be considered finished.

On each EFN, as Figure 3 shows, there are four types of queues: prediction queues with the backlogs as $A_{i,0}(t)$, ..., $A_{i,W_i-1}(t)$, arrival queue $A_{i,-1}(t)$, local processing queue $Q_i^{(e,l)}(t)$, and offloading queue $Q_i^{(e,o)}(t)$. In time slot t , prediction queue $A_{i,w}(t)$ ($0 \leq w \leq W_i - 1$) stores untreated workload that will arrive in time slot $(t+w)$. Workload that actually arrives at EFN i is stored in the arrival queue $A_{i,-1}(t)$, awaiting being forwarded to the local processing

¹We do not specify any particular learning method in this paper, since our work aims to explore the *fundamental* benefits of predictive offloading. In practice, one can leverage machine learning techniques such as time-series prediction methods [15] for workload arrival prediction.

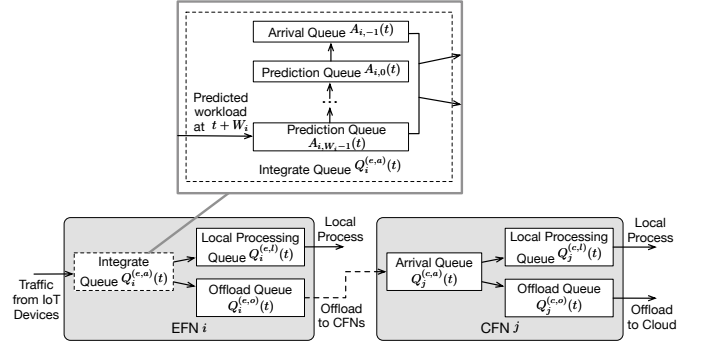


Fig. 3: Queueing model of the system.

queue $Q_i^{(e,l)}(t)$ or the offloading queue $Q_i^{(e,o)}(t)$. Workload in $Q_i^{(e,l)}(t)$ will be processed locally by EFN i , while workload in $Q_i^{(e,o)}(t)$ will be offloaded to CFNs in set \mathcal{M}_i .

1) *Prediction Queues and Arrival Queues in EFNs*: Within each time slot t , in addition to the current arrivals in the arrival queue, EFN i can also forward future arrivals in the prediction queues. We define $\mu_{i,w}(t)$ as the amount of output workload from $A_{i,w}(t)$, for $w \in \{-1, 0, \dots, W_i - 1\}$. Such workload should be distributed to the local processing queue and offloading queue. We denote the amounts of workload to be distributed to the local processing queue and offloading queue as $b_i^{(e,l)}(t)$ and $b_i^{(e,o)}(t)$, respectively, such that

$$0 \leq b_i^{(e,\beta)}(t) \leq b_{i,\max}^{(e,\beta)}, \quad \forall \beta \in \{l, o\} \quad (1)$$

where each $b_{i,\max}^{(e,\beta)}$ is a positive constant. As a result, we have

$$\sum_{w=-1}^{W_i-1} \mu_{i,w}(t) = b_i^{(e,l)}(t) + b_i^{(e,o)}(t). \quad (2)$$

Next, we consider the queueing dynamics for different types of queues in EFN, respectively.

Regarding $A_{i,w}(t)$, it is updated whenever pre-service is finished and the lookahead window moves one slot ahead at the end of each time slot. Therefore, we have

(i) If $w = W_i - 1$, then

$$A_{i,W_i-1}(t+1) = A_i(t + W_i). \quad (3)$$

(ii) If $0 \leq w \leq W_i - 2$, then

$$A_{i,w}(t+1) = [A_{i,w+1}(t) - \mu_{i,w+1}(t)]^+, \quad (4)$$

where $[x]^+ \triangleq \max\{x, 0\}$ for $x \in \mathbb{R}$. In time slot $(t+1)$, the amount of workload that will arrive after $(W_i - 1)$ time slots is $A_i(t + W_i)$ and it remains unknown until time slot $(t+1)$.

Regarding the arrival queue $A_{i,-1}(t)$, it records the actual backlog of EFN i with the update equation as follows:

$$A_{i,-1}(t+1) = [A_{i,-1}(t) - \mu_{i,-1}(t)]^+ + [A_{i,0}(t) - \mu_{i,0}(t)]^+. \quad (5)$$

Note that $\mu_{i,-1}(t)$ denotes the amount of distributed workload that have already being in $A_{i,-1}(t)$.

Next, we introduce an integrate queue with a backlog size as the sum of all prediction queues and the arrival queue on EFN i , denoted by $Q_i^{(e,a)}(t) \triangleq \sum_{w=-1}^{W_i-1} A_{i,w}(t)$. Under *fully-efficient* [16] service policy, $Q_i^{(e,a)}(t)$ is updated as

$$Q_i^{(e,a)}(t+1) = [Q_i^{(e,a)}(t) - (b_i^{(e,l)}(t) + b_i^{(e,o)}(t))]^+ + A_i(t + W_i). \quad (6)$$

The input of integrate queue $Q_i^{(e,a)}(t)$ consists of the predicted workload that will arrive at EFN i in time slot $(t + W_i)$, while its output consists of workloads being forwarded to the local processing queue and the offloading queue. Note that $b_i^{(e,l)}(t) + b_i^{(e,o)}(t)$ is the output capacity of integrate queue $Q_i^{(e,a)}(t)$ in time slot t . If the capacity is larger than the queue backlog size, the true output amount will be smaller than $b_i^{(e,l)}(t) + b_i^{(e,o)}(t)$.

2) *Offloading Queues in EFNs*: In time slot t , workload in queue $Q_i^{(e,o)}(t)$ will be offloaded to CFNs in set \mathcal{M}_i . The transmission capacities are determined by the transmit power decisions $(p_{i,j}(t))_{j \in \mathcal{M}_i}$, where $p_{i,j}(t)$ is the transmit power from EFN i to CFN j . The transmit power is nonnegative and the total transmit power of each EFN is upper bounded, *i.e.*,

$$p_{i,j}(t) \geq 0, \quad \forall i \in \mathcal{N}, j \in \mathcal{M}_i \text{ and } t, \quad (7)$$

$$\sum_{j \in \mathcal{M}_i} p_{i,j}(t) \leq p_{i,\max}, \quad \forall i \in \mathcal{N} \text{ and } t. \quad (8)$$

According to Shannon's capacity formula [17], the transmission capacity from EFN i to CFN j is

$$R_{i,j}(t) \triangleq \hat{R}_{i,j}(p_{i,j}(t)) = \tau_0 B \log_2 \left(1 + \frac{p_{i,j}(t) H_{i,j}(t)}{N_0 B} \right), \quad (9)$$

where τ_0 is the length of each time slot, B is the channel bandwidth, $H_{i,j}(t)$ is the wireless channel gain between EFN i and CFN j , and N_0 is the system power spectral density of the additive white Gaussian noise. Note that $H_{i,j}(t)$ is an uncontrollable environment state with positive upper bound H_{\max} . We do not consider the interferences among fog nodes and tiers. By adjusting the transmit power $p_{i,j}(t)$, we can offload different amounts of workload from EFN i to CFN j in time slot t . Accordingly, the update equation of offloading queue $Q_i^{(e,o)}(t)$ is

$$Q_i^{(e,o)}(t+1) \leq [Q_i^{(e,o)}(t) - \sum_{j \in \mathcal{M}_i} R_{i,j}(t)]^+ + b_i^{(e,o)}(t), \quad (10)$$

where $\sum_{j \in \mathcal{M}_i} R_{i,j}(t)$ is the total transmission capacity to EFN i in time slot t . The inequality here means that the actual arrival of $Q_i^{(e,o)}(t)$ may be less than $b_i^{(e,o)}(t)$, because $b_i^{(e,o)}(t)$ is the transmission capacity from integrate queue $Q_i^{(e,a)}(t)$ to offloading queue $Q_i^{(e,o)}(t)$ instead of the amount of truly transmitted workload. Recall that we assume the transmission latency from EFT to CFT is negligible compared to the length of each time slot, the workload transmission in each time slot can be accomplished by the end of that time slot.

C. Queueing Model for Central Fog Node

Figure 3 also shows the queueing model on CFN. Each CFN $j \in \mathcal{M}$ maintains three queues: an arrival queue $Q_j^{(c,a)}(t)$, a local processing queue $Q_j^{(c,l)}(t)$, and an offloading queue $Q_j^{(c,o)}(t)$. Similar to EFNs, workload offloaded from the EFT will be firstly stored in the arrival queue, then distributed to $Q_j^{(c,l)}(t)$ for local processing and to $Q_j^{(c,o)}(t)$ for further offloading.

1) *Arrival Queues in CFNs*: The arrivals on CFN j consist of workloads offloaded from EFNs in the set \mathcal{N}_j . We denote the amounts of workloads distributed to the local processing queue and offloading queue in time slot t as $b_j^{(c,l)}(t)$ and $b_j^{(c,o)}(t)$, respectively, such that

$$0 \leq b_j^{(c,\beta)}(t) \leq b_{j,\max}^{(c,\beta)}, \quad \forall \beta \in \{l, o\}, \quad (11)$$

where each $b_{j,\max}^{(c,\beta)}$ is a positive constant. Accordingly, $Q_j^{(c,a)}(t)$ is updated as follows:

$$Q_j^{(c,a)}(t+1) \leq [Q_j^{(c,a)}(t) - (b_j^{(c,l)}(t) + b_j^{(c,o)}(t))]^+ + \sum_{i \in \mathcal{N}_j} R_{i,j}(t). \quad (12)$$

2) *Offloading Queues in CFNs*: For each CFN $j \in \mathcal{M}$, its offloading queue $Q_j^{(c,o)}(t)$ stores the workload to be offloaded to the cloud. We define $D_j(t)$ as the transmission capacity of the wired link from CFN j to the cloud during time slot t , which depends on the network state and is upper bounded by some constant D_{\max} for all j and t . Then we have the following update function for $Q_j^{(c,o)}(t)$:

$$Q_j^{(c,o)}(t+1) \leq [Q_j^{(c,o)}(t) - D_j(t)]^+ + b_j^{(c,o)}(t). \quad (13)$$

Note that the amount of actually offloaded workload to the cloud is $\min\{Q_j^{(c,o)}(t), D_j(t)\}$.

D. Local Processing Queues on EFNs and CFNs

We assume that all fog nodes are able to adjust their CPU frequencies in each time slot, by applying *dynamic voltage and frequency scaling* (DVFS) techniques [18]. Next, we define $L_k^{(\alpha)}$ as the number of CPU cycles that fog node $k \in \mathcal{N} \cup \mathcal{M}$ requires to process one bit of workload, where α is an indicator of fog node k 's type ($\alpha = e$ if k is an EFN, and $\alpha = c$ if k is a CFN). $L_k^{(\alpha)}$ is assumed constant and can be measured offline [19]. Therefore, the local processing capacity of fog node k is $f_k^{(\alpha)}(t)/L_k^{(\alpha)}$. The local processing queue on fog node k evolves as follows:

$$Q_k^{(\alpha,l)}(t+1) \leq [Q_k^{(\alpha,l)}(t) - \tau_0 f_k^{(\alpha)}(t)/L_k^{(\alpha)}]^+ + b_k^{(\alpha,l)}(t). \quad (14)$$

All CPU frequencies are nonnegative and finite:

$$0 \leq f_k^{(\alpha)}(t) \leq f_{k,\max}^{(\alpha)}, \quad \forall k \in \mathcal{N} \cup \mathcal{M} \text{ and } t, \quad (15)$$

where each $f_{k,\max}^{(\alpha)}$ is a positive constant.

E. Power Consumptions

The total power consumptions $P(t)$ of fog tiers in time slot t consist of the processing power consumption and wireless transmit power consumption. Given a local CPU with frequency f , its power consumption per time slot is $\tau_0 \varsigma f^3$, where ς is a parameter depending on the deployed hardware and is measurable in practice [20]. Thus $P(t)$ is defined as follows:

$$P(t) \triangleq \hat{P}(\mathbf{f}(t), \mathbf{p}(t)) = \sum_{i \in \mathcal{N}} \tau_0 \varsigma (f_i^{(e)}(t))^3 + \sum_{j \in \mathcal{M}} \tau_0 \varsigma (f_j^{(c)}(t))^3 + \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}_i} \tau_0 p_{i,j}(t), \quad (16)$$

where $\mathbf{f}(t) \triangleq ((f_i^{(e)}(t))_{i \in \mathcal{N}}, (f_j^{(c)}(t))_{j \in \mathcal{M}})$ is the vector of all CPU frequencies, and $\mathbf{p}(t) \triangleq (\mathbf{p}_i(t))_{i \in \mathcal{N}}$ in which $\mathbf{p}_i(t) = (p_{i,j}(t))_{j \in \mathcal{M}_i}$ is the transmit power allocation of EFN i .

F. Problem Formulation

We define the long-term time-average expectation of total power consumptions \bar{P} and total queue backlog \bar{Q} as follows:

$$\bar{P} \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\{P(t)\}, \quad (17)$$

$$\bar{Q} \triangleq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{\beta \in \{a, l, o\}} \left(\sum_{i \in \mathcal{N}} \mathbb{E}\{Q_i^{(e, \beta)}(t)\} + \sum_{j \in \mathcal{M}} \mathbb{E}\{Q_j^{(c, \beta)}(t)\} \right). \quad (18)$$

In this paper, we aim to minimize the long-term time-average expectation of total power consumptions \bar{P} , while ensuring the stability of all queues in the system, *i.e.*, $\bar{Q} < \infty$. The problem formulation is given by

$$\begin{aligned} & \underset{\{\mathbf{b}(t), \mathbf{f}(t), \mathbf{p}(t)\}_t}{\text{Minimize}} && \bar{P} \\ & \text{Subject to} && (1)(7)(8)(11)(15), \\ & && \bar{Q} < \infty. \end{aligned} \quad (19)$$

V. ALGORITHM DESIGN

A. Predictive Algorithm

To solve problem (19), we adopt Lyapunov optimization techniques [16] [21] to decouple the problem into a series of subproblems over time slots. We show the detail of this process in Appendix A. By solving each of these subproblems during each time slot, we propose PORA, an efficient and predictive scheme conducts workload offloading in an online and distributed manner. We show the pseudocode of PORA in Algorithm 1. Note that symbol $\alpha \in \{e, c\}$ indicates the type of fog node. Specifically, for each fog node k , $\alpha = e$ if k is an EFN and CFN otherwise. Next, we introduce PORA in detail.

1) *Offloading Decision*: In each time slot t , under PORA, each fog node $k \in \mathcal{N} \cup \mathcal{M}$ decides the amounts of workload scheduled to the offloading queue and the local processing queue, denoted by $b_k^{(\alpha, l)}(t)$ and $b_k^{(\alpha, o)}(t)$, respectively. Such decisions are obtained by solving the following problem:

$$\underset{0 \leq b_k^{(\alpha, \beta)} \leq b_{k, \max}^{(\alpha, \beta)}}{\text{Minimize}} \left(Q_k^{(\alpha, \beta)}(t) - Q_k^{(\alpha, a)}(t) \right) b_k^{(\alpha, \beta)}, \quad (20)$$

where $\beta \in \{l, o\}$. Accordingly, the optimal solution to (20) is

$$b_k^{(\alpha, \beta)}(t) = \begin{cases} b_{k, \max}^{(\alpha, \beta)}, & \text{if } Q_k^{(\alpha, \beta)}(t) < Q_k^{(\alpha, a)}(t), \\ 0, & \text{otherwise.} \end{cases} \quad (21)$$

From (21), we see that, to determine the optimal solutions $b_i^{(e, l)}(t)$ and $b_i^{(e, o)}(t)$, each EFN i would compare its integrate queue backlog size $Q_i^{(e, a)}(t)$ with its local processing queue backlog size $Q_i^{(e, l)}(t)$ and offloading queue backlog size $Q_i^{(e, o)}(t)$, respectively. Particularly, if there is too much workload in its integrate queue compared to its local queue ($Q_i^{(e, l)}(t) < Q_i^{(e, a)}(t)$), then it will offload as much workload (up to $b_{i, \max}^{(e, l)}$) as possible to its local queue. Likewise, if

Algorithm 1 Predictive Offloading and Resource Allocation (PORA) in One Time Slot

```

1: Initialize  $\mathbf{b}(t) \leftarrow \mathbf{0}$ ,  $\mathbf{f}(t) \leftarrow \mathbf{0}$ ,  $\mathbf{p}(t) \leftarrow \mathbf{0}$ .
2: for each fog node  $k \in \mathcal{N} \cup \mathcal{M}$  do
3:   %%Make Offloading Decisions
4:   if  $Q_k^{(\alpha, a)}(t) > Q_k^{(\alpha, l)}(t)$  then
5:     Set  $b_k^{(\alpha, l)}(t) \leftarrow b_{k, \max}^{(\alpha, l)}$ .
6:   end if
7:   if  $Q_k^{(\alpha, a)}(t) > Q_k^{(\alpha, o)}(t)$  then
8:     Set  $b_k^{(\alpha, o)}(t) \leftarrow b_{k, \max}^{(\alpha, o)}$ .
9:   end if
10:  %%Local CPU Resource Allocation
11:  Set  $f_k^{(\alpha)}(t) \leftarrow \min\{\sqrt{Q_k^{(\alpha, l)}(t)}/3V\zeta L_k^{(\alpha)}, f_{k, \max}^{(\alpha)}\}$ .
12: end for
13: %%Transmit Power Allocation
14: for each EFN  $i \in \mathcal{N}$  do
15:   Set  $\lambda_{\min} \leftarrow 0$ .
16:   Set  $\lambda_{\max} \leftarrow \max_{j \in \mathcal{M}_i} \frac{(Q_i^{(e, o)} - Q_j^{(c, a)})H_{i, j}(t)}{N_0} - V$ .
17:   while  $\lambda_{\max} - \lambda_{\min} > \varepsilon$  do
18:     %%Water Filling with Bisection Method
19:     Set  $\lambda^* \leftarrow (\lambda_{\min} + \lambda_{\max})/2$ .
20:     Set  $p_{i, j}(t) \leftarrow B \left[ \frac{Q_i^{(e, o)}(t) - Q_j^{(c, a)}(t)}{V + \lambda^*} - \frac{N_0}{H_{i, j}(t)} \right]^+$ .
21:     if  $\sum_{j \in \mathcal{M}_i} p_{i, j}(t) > p_{i, \max}$  then
22:       Set  $\lambda_{\max} \leftarrow \lambda^*$ .
23:     else
24:       Set  $\lambda_{\min} \leftarrow \lambda^*$ .
25:     end if
26:   end while
27: end for
28: Enforce scheduling decisions  $\mathbf{b}(t)$ ,  $\mathbf{f}(t)$ , and  $\mathbf{p}(t)$ .

```

its integrate queue is loaded with more workload than its offloading queue ($Q_i^{(e, o)}(t) < Q_i^{(e, a)}(t)$), it will offload up to an amount of $b_{i, \max}^{(e, o)}$ workload to its offloading queue.

Notably, if the backlog size of the EFN i 's integrate queue is larger than both its local queue and offloading queue, then the EFN will transmit the workload one by one unit (*e.g.* packets); each unit of workload is either sent to the EFN i 's local queue or its offloading queue, such that the amounts of workload distributed to such two queues are no greater than $b_{i, \max}^{(e, l)}$ and $b_{i, \max}^{(e, o)}$, respectively. In practice, the workload distributing strategy is left as a degree of freedom to be specified in the implementation of PORA. In our simulation, we adopt the following distributing strategy. When an EFN i 's integrate queue backlog size is greater than both its local queue and its offloading backlog size, then it will transmit workload to its local queue until the amount of transmitted workload reaches $b_{i, \max}^{(e, l)}$. Then the rest workload in the integrate queue is transmitted to the offloading queue until the amount of distributed workload reaches $b_{i, \max}^{(e, o)}$. Such a process terminates whenever the integrate queue becomes empty.

The decision making process is similar for CFNs. Specifically, each CFN j determines $b_j^{(c, l)}(t)$ and $b_j^{(c, o)}(t)$ by comparing its arrival queue backlog size $Q_j^{(c, a)}(t)$ with its local processing queue backlog size $Q_j^{(c, l)}(t)$ and offloading queue backlog size $Q_j^{(c, o)}(t)$, respectively.

Remark: For each EFN, we can view the difference between the backlog sizes of its integrate queue and its local processing/offloading queue as its willingness of workload transmission. If such willingness is positive, then the EFN will transmit as much workload as possible from its integrate queue; otherwise, the EFN will leave the workload not distributed in the current time slot. In such a way, PORA always endeavors to balance the integrate queue backlog and the local/offloading queue backlog. Likewise, under PORA, each CFN determines its offloading decisions upon the difference between the backlog sizes of its arrival queue and its local processing/offloading queue to ensure the queue stability.

2) *Local CPU Frequency Allocation:* Under PORA, in each time slot t , each fog node $k \in \mathcal{N} \cup \mathcal{M}$ sets its local CPU frequency $f_k^{(\alpha)}(t)$ by solving the following subproblem:

$$\text{Minimize } V\zeta(f_k^{(\alpha)})^3 - Q_k^{(\alpha,l)}(t) f_k^{(\alpha)} / L_k^{(\alpha)}. \quad (22)$$

$$0 \leq f_k^{(\alpha)} \leq f_{k,\max}^{(\alpha)}$$

By setting the second derivative of the objective function in (22) to zero, we can obtain the optimal CPU frequency $f_k^{(\alpha)}(t)$ to be set by fog node k as

$$f_k^{(\alpha)}(t) = \min \left\{ \sqrt{Q_k^{(\alpha,l)}(t) / 3V\zeta L_k^{(\alpha)}}, f_{k,\max}^{(\alpha)} \right\}. \quad (23)$$

We prove the optimality of (23) in Appendix B.

Remark: When $f_k^{(\alpha)}(t) < f_{k,\max}^{(\alpha)}$, the allocated CPU frequency $f_k^{(\alpha)}(t)$ is proportional to the square root of the backlog size of local processing queue $Q_k^{(\alpha,l)}(t)$ and the inverse of the value of parameter V . This shows that, on the one hand, PORA would allocate as much CPU frequency as possible to process the workload in the queues. On the other hand, the value of parameter V determines the tradeoff between power consumption and the backlog sizes of queues: a small value of V will encourage the fog node to allocate more CPU frequency to process the workload and hence a small queue backlog size; in contrast, a large value of V will make the fog node more conservative to allocate resources, leading to less power consumptions but a large queue backlog size as well. In practice, the choice of the value of V is dependent on the system design objective.

3) *Power Allocations for EFNs:* In each time slot t , under PORA, each EFN $i \in \mathcal{N}$ determines its allocated transmit power $p_{i,j}(t)$ by solving the following optimization problem.

$$\begin{aligned} & \text{Minimize } \sum_{j \in \mathcal{M}_i} \left[V p_{i,j} - m_{i,j}(t) \log_2(1 + l_{i,j}(t) p_{i,j}) \right] \\ & \text{Subject to } \sum_{j \in \mathcal{M}_i} p_{i,j} \leq p_{i,\max}, \\ & p_{i,j} \geq 0, \quad \forall j \in \mathcal{M}_i, \end{aligned} \quad (24)$$

where $m_{i,j}(t) \triangleq (Q_i^{(e,o)}(t) - Q_j^{(c,a)}(t))B$ and $l_{i,j}(t) \triangleq \frac{H_{i,j}(t)}{N_0 B}$. By applying *water-filling* algorithm [22], we obtain the optimal solution to problem (24) as

$$p_{i,j}(t) = [m_{i,j}(t) / (V + \lambda^*) - 1 / l_{i,j}(t)]^+, \quad \forall j \in \mathcal{M}_i, \quad (25)$$

where λ^* is the optimal Lagrangian variable that satisfies

$$\sum_{j \in \mathcal{M}_i} [m_{i,j}(t) / (V + \lambda^*) - 1 / l_{i,j}(t)]^+ = p_{i,\max}. \quad (26)$$

The optimality of such solutions is proven in Appendix C. We adopt bisection method (line 15-25 in Algorithm 1) to obtain the value of λ^* with its lower and upper bounds as λ_{\min} and λ_{\max} , respectively. Note that the value of λ^* converges asymptotically to the optimum λ^{opt} as the tolerance parameter ε approaches zero, such that $|\lambda^* - \lambda^{\text{opt}}| \leq \varepsilon/2$.

Remark: PORA tends to allocate more transmit power to the CFN with smaller arrival queue backlog size $Q_j^{(c,a)}(t)$ for load balancing. When $Q_j^{(c,a)}(t) \geq Q_i^{(e,o)}(t)$, we have $m_{i,j}(t) \leq 0$ and $p_{i,j}(t) = 0$, i.e., EFN i allocates no transmit power to CFN j unless the backlog size of the arrival queue on CFN j is greater than that of the offloading queue on EFN i . By increasing the value of V , transmit power consumption will be reduced but the backlog size will increase as well.

B. Computational Complexity of PORA

During each time slot, part of the computational complexity concentrates on the calculation for CPU frequency settings and offloading decision makings. Since the calculation (line 3-11) requires only constant time for each fog node, the total complexity of these steps is $O(N + M)$. Next, each EFN i applies the bisection method (line 15-26) to calculate the optimal dual variable, with a complexity of $O(\log_2((\lambda_{\max} - \lambda_{\min})/\varepsilon) + |\mathcal{M}_i|)$. After that, EFN i determines the transmit power to each CFN in the set \mathcal{M}_i . In the worst case, each EFN is potentially connected to all CFNs, thus the total complexity of PORA algorithm is $O(M \times N)$.

C. Performance Analysis

We conduct theoretical analysis on the relationship between the average power consumption \bar{P} and queue backlog \bar{Q} under PORA scheme in the non-predictive case ($W_i = 0, \forall i \in \mathcal{N}$), and then analyze the benefits of predictive offloading in terms of latency reduction.

1) *Time-average Power Consumption and Queue Backlog:* Let P^* be the achievable minimum of \bar{P} over all feasible non-predictive policies. We have the following theorem.

Theorem 1: Assume the system arrival lies in the interior of the capacity region and $\mathbf{Q}(0) < \infty$. Under PORA, without prediction, there exist constants $\theta > 0$ and $\epsilon > 0$ such that

$$\bar{P} \leq \theta/V + P^*, \quad \bar{Q} \leq (\theta + V P_{\max})/\epsilon,$$

where \bar{P} and \bar{Q} are defined in (17) and (18), respectively.

The proof is quite standard and hence omitted here.

Remark: By *Little's* theorem [23], the average queue backlog size is proportional to the average queuing latency. Therefore, Theorem 1 implies that by adjusting parameter V , PORA can achieve an $[O(1/V), O(V)]$ power-latency tradeoff in the non-predictive case. Furthermore, the average power consumption \bar{P} approaches the optimum P^* asymptotically as the value of V increases to infinity.

2) *Latency Reduction:* We analyze the latency reduction induced by PORA under perfect prediction compared to the non-predictive case. In particular, we denote the prediction window vector $(W_i)_{i \in \mathcal{N}}$ by \mathbf{W} and the corresponding delay reduction by $\eta(\mathbf{W})$. For each unit of workload on EFN i , let $\pi_{i,w}$ denote the steady-state probability that it experiences a latency of w time slots in $A_{i,-1}(t)$. Without prediction, the average latency on its arrival queues is $d = \sum_{i \in \mathcal{N}} \lambda_i \sum_{w \geq 1} w \pi_{i,w} / \sum_{i \in \mathcal{N}} \lambda_i$. Then we have the following theorem.

Theorem 2: Suppose the system steady-state behavior depends only on the statistical behaviors of the arrivals and service processes. Then the latency reduction $\eta(\mathbf{W})$ is

$$\eta(\mathbf{W}) = \frac{\sum_{i \in \mathcal{N}} \lambda_i \left(\sum_{1 \leq w \leq W_i} w \pi_{i,w} + W_i \sum_{w \geq 1} \pi_{i,w+W_i} \right)}{\sum_{i \in \mathcal{N}} \lambda_i}. \quad (27)$$

Furthermore, if $d < \infty$, as $\mathbf{W} \rightarrow \infty$, i.e., with infinite predictive information, we have

$$\lim_{\mathbf{W} \rightarrow \infty} \eta(\mathbf{W}) = d. \quad (28)$$

We relegate the proof of Theorem 2 to Appendix D.

Remark: Theorem 2 implies that predictive offloading conduces to a shorter workload latency; in other words, with predicted information, PORA can break the barrier of $[O(1/V), O(V)]$ power-latency tradeoff. Furthermore, the latency reduction induced by PORA is proportional to the inverse of the prediction window size, and approaches zero as prediction window sizes go to infinity. In our simulations, we see that PORA can effectively shorten the average arrival queue latency with only mild-value of future information.

D. Impact of Network Topology

Fog computing systems generally proceed in wireless environments, thus the network topology of such systems is usually dynamic and may change over time slots. However, at the beginning of each time slot, the network topology is observed and deemed fixed by the end of the time slot. Therefore, in the following, we put the focus of our discussion on the impact of network topology within each time slot.

Recall that in our settings, each EFN has access to only a subset of CFNs in its vicinity. For each EFN i , the subset of its accessible EFNs is denoted by \mathcal{M}_i with a size of $|\mathcal{M}_i|$. From the perspective of graph theory, we can view the interconnection among fog nodes of different tiers as a directed graph, in which each vertex corresponds to a fog node and each edge indicates a directed connection between nodes. Hence, the value of $|\mathcal{M}_i|$ can be regarded as the out-degree of EFN i , which is an important parameter of network topology that measures the number of directed connections originating from EFN i . Due to time-varying wireless dynamics, the out-degree of each fog node may vary over time slots; consequentially, the resulting topology would significantly affect the system performance. In the following, we discuss such impacts under two channel conditions, respectively.

On the one hand, within each time slot, poor channel conditions (e.g. in terms of low SINR) would often lead to unreliable or even unavailable connections among fog nodes and hence a network topology with a relatively smaller out-degree of nodes. In this case, each fog node may have a very limited freedom to choose the best target node to offload its workloads, further leading to backlog imbalance among fog nodes or even overloading in its upper tier with a large cumulative queue backlog size. Besides, poor channel conditions may also require more power consumptions to ensure reliable communication between successive fog nodes.

On the other hand, within each time slot, good channel conditions allow each fog node to have a broader access to the fog nodes in its upper tier, resulting a network topology

TABLE V: Simulation Settings

Parameter	Value
B	2 MHz
$H_{i,j}(t), \forall i \in \mathcal{N}, j \in \mathcal{M}$	$24 \log_{10} d_{i,j} + 20 \log_{10} 5.8+60^a$
N_0	-174 dBm/Hz
$P_{i,\max}, \forall i \in \mathcal{N}$	500 mW
$L_i^{(e)} \forall i \in \mathcal{N}, L_j^{(c)} \forall j \in \mathcal{M}$	297.62 cycles/bit
$f_{i,\max}^{(e)}, \forall i \in \mathcal{N}$	4 G cycles/s
$f_{j,\max}^{(c)}, \forall j \in \mathcal{M}$	8 G cycles/s
ς	$10^{-27} \text{ W} \cdot \text{s}^3 / \text{cycle}^3$
$b_{i,\max}^{(e,t)}, b_{i,\max}^{(e,o)}, \forall i \in \mathcal{N}$	6 Mb/s
$b_{j,\max}^{(c,t)}, b_{j,\max}^{(c,o)}, \forall j \in \mathcal{M}$	12 Mb/s
$D_j(t), \forall j \in \mathcal{M}, t$	6 Mb/s

^a $d_{i,j}$ is the distance between EFN i and CFN j .

with a relatively larger out-degree of nodes. In this case, each fog node is able to conduct better decision-making with more freedom in choosing the fog nodes in its upper fog tier, thereby achieving a better tradeoff between power consumptions and backlog sizes.

E. Use Cases

In practice, PORA can be applied as a theoretical framework to design the offloading schemes for fog computing systems under various use cases, such as public safety systems, intelligent transportation, and smart healthcare systems. For example, in a public safety system, each street is usually deployed with multiple smart cameras (IoT devices). At runtime, such smart cameras would upload real-time vision data to one of their accessible EFNs. Each EFN aggregates such data to extract or even analyze the instant road conditions within multiple streets. Such EFNs can upload some of the workload to their upper-layered CFNs (each taking charge of one community consisting of several streets) with greater computing capacities. Each CFN can further offload the workload to the cloud via optical fiber links. For latency-sensitive applications, the real-time vision data will be processed locally on EFNs or offloaded to CFNs. For latency-insensitive applications with intensive computation demand, the data will be offloaded to the cloud through the fog nodes. PORA conduces to the design of dynamic and online offloading and resource allocation schemes to support such fog systems with various applications.

VI. NUMERICAL RESULTS

We conduct extensive simulations to evaluate PORA and its variants. The parameter settings in our simulation are based on the commonly adopted wireless environment settings that have been used in [24], [25]. The simulation is conducted on a MacBook Pro with 2.3 GHz Intel Core i5 processor and 8GB 2133 MHz LPDDR3 memory, and the simulation program is implemented using Python 3.7. This section firstly presents the basic settings of our simulations, and then provides the key results under perfect and imperfect prediction, respectively.

A. Basic Settings

We simulate a hierarchal fog computing system with 80 EFNs and 20 CFNs. All EFNs have a uniform prediction window size W , which varies from 0 to 30. Note that $W = 0$

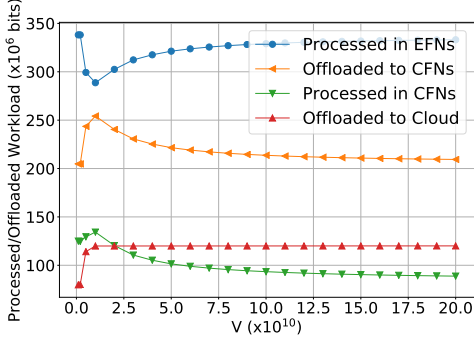
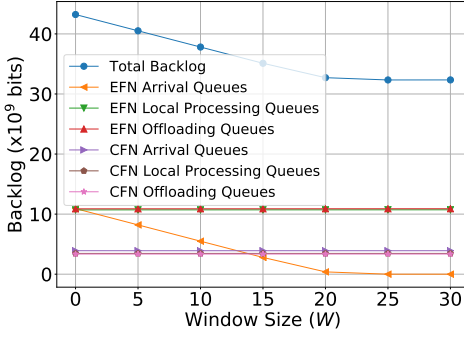
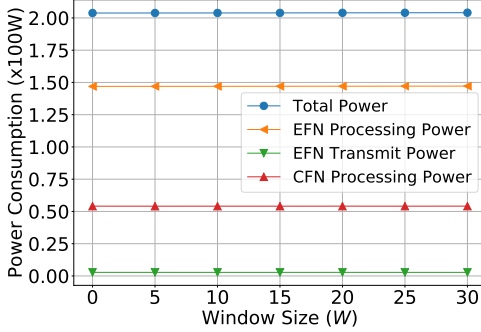


Fig. 4: Offloading decisions when $W = 10$.



(a) Queue backlogs.



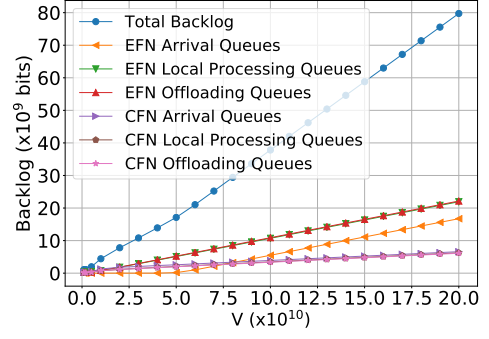
(b) Power consumptions.

Fig. 5: Performance of PORA vs. W when $V = 10^{11}$.

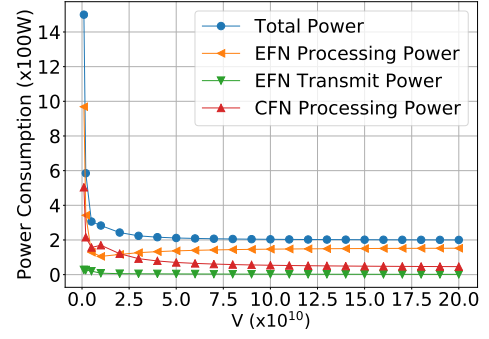
refers to the case without prediction. For each EFN i , its accessible CFN set \mathcal{M}_i is chosen uniformly randomly from the power set of the CFN set with size $|\mathcal{M}_i| = 5$. We set the time slot length $\tau_0 = 1$ second. During each time slot, workload arrives to the system in the unit of packets, each with a fixed size of 4096 bits. The packet arrivals are drawn from previous measurements [26], where the average flow arrival rate is 538 flows/s, and the distribution of flow size has a mean of 13 Kb. Given these settings, the average arrival rate is about 7 Mbps. All results are averaged over 50000 time slots. We list all other parameter settings in TABLE V.

B. Evaluation with Perfect Prediction

Under the perfect prediction settings, we evaluate how the values of parameter V and prediction window size W



(a) Queue backlogs.



(b) Power consumptions.

Fig. 6: Performance of PORA when $W = 10$.

influence the performance of PORA, respectively.

System Performance under Different Values of V : Figure 4 shows the impact of parameter V on the offloading decisions of PORA: When the value of V is around 10^{10} , the time-average amount of locally processed workload on EFNs reaches the bottom of the curve, while other offloading decisions induce the peak workload. The reason is that the offloading decisions are not only determined by the value of V , but also influenced by the queue backlog sizes.

Figure 6 presents the impact of the value of V on different types of queues and power consumptions in the system, respectively. As the value of V increases, we see a rising trend in the sizes of all types of queue backlogs, and a roughly falling trend in all types of power consumptions.

System Performance with Different Values of Prediction Window Size W : Figures 5(a) and 5(b) show the system performance with the prediction window size W varying from 0 to 30. With perfect prediction, PORA effectively shortens the average queuing latencies on EFN arrival queues – eventually close to zero with no extra power consumption and only a mild-value of prediction window size ($W = 20$ in this case).

PORA vs. PORA- d (Low-Sampling Variant): In practice, since PORA requires to sample system dynamics across various fog nodes, it may incur considerable sampling overheads. By adopting the idea of randomized load balancing techniques [27], we propose PORA- d , a variant of PORA that reduces the sampling overheads by probing d ($d \in \{1, 2, 3, 4\}$)² CFNs and conducting resource allocation on which are uniformly chosen for each EFN from its accessible CFN set.

²When $d = 1$, the scheme degenerates to uniform random sampling.

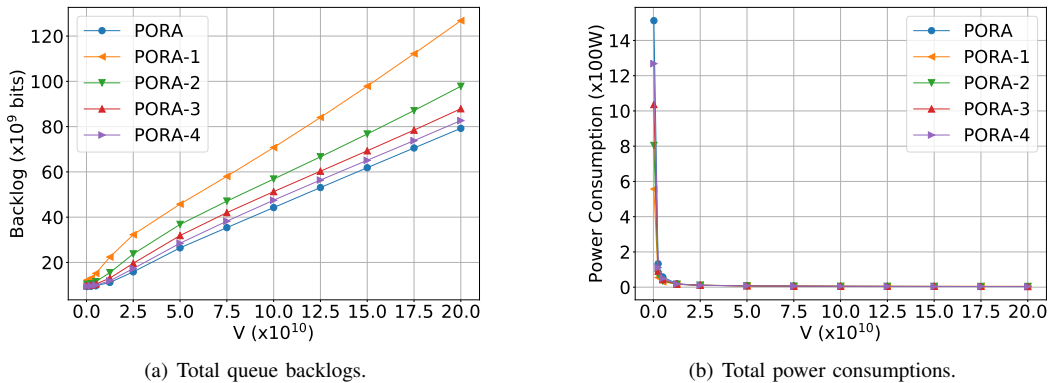


Fig. 7: Performance of variants of PORA.

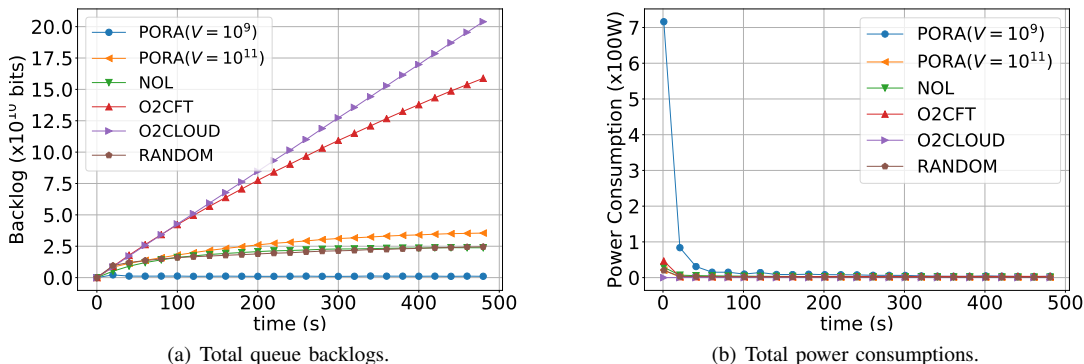


Fig. 8: Comparison between PORA and baselines.

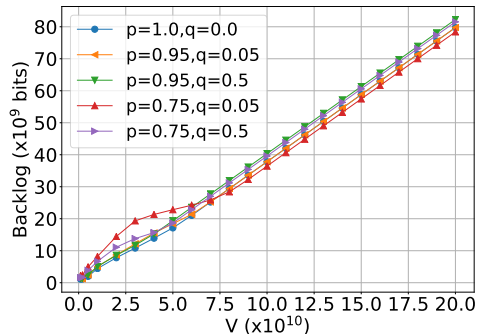
Figure 7 compares the performance of PORA with PORA- d . We observe that PORA achieves the smallest queue backlog size. The result is reasonable since each EFN has access to 5 CFNs under PORA, more than the $d \leq 4$ CFNs under PORA- d . As a result, each EFN has more chance to access to the CFNs with better wireless channel condition and processing capacity under PORA when compared with PORA- d . The observation that the queue backlog size increases as d decreases further verifies our analysis. In fact, we can view d as the degree of each EFN in the network topology. As d decreases, the system performance degrades. However, when the value of V is sufficiently large, PORA- d achieves the similar power consumptions as PORA and the ratio of increment in the backlog size is small. For example, when $V = 2 \times 10^{11}$, PORA-4 achieves 4.3% larger backlog size than PORA, and PORA-3 achieves 10.9% larger backlog size than PORA. In summary, PORA- d (when $d = 2, 3, 4$) can reduce the sampling overheads by trading off only a little performance degradation under large V .

Comparison of PORA and Baselines: We introduce four baselines to evaluate the performance of PORA: (1) NOL (No Offloading): All nodes in the EFT process packets locally. (2) O2CFT (Offload to CFT): All packets are offloaded to the CFT and processed therein. (3) O2CLOUD (Offload to Cloud): All packets are offloaded to the cloud. (4) RANDOM: Each fog node randomly chooses to offload each packet or process

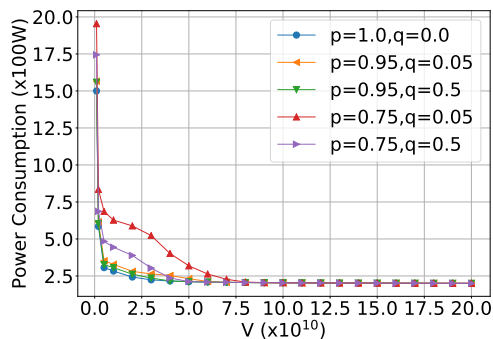
it locally with equal chance. Note that all above baselines are also assumed capable of pre-serving future workloads in the prediction window. Figure 8 compares the instant total queue backlog sizes and power consumptions over time slots under the five schemes (PORA, NOL, O2CFT, O2CLOUD, RANDOM), where $W = 10$ and $V \in \{10^9, 10^{11}\}$.

We observe that scheme O2CLOUD achieves the minimum power consumptions, but incurs constantly increasing queue backlog sizes over time. The reasons are shown as follows. On one hand, in our settings, the mean power consumption for transmitting workload from EFT to CFT is smaller than the mean power consumption of processing the same amount of workload on fog nodes; under scheme O2CLOUD, only wireless transmit power is consumed and hence the minimum is achieved. On the other hand, all the workload must travel through all fog tiers before being offloaded to the cloud, which results in network congestion within fog tiers and thus workload accumulation with increasing queue backlogs.

As Figure 8 illustrates, PORA achieves the maximum power consumptions but the smallest backlog size when $V = 10^9$. Upon convergence of PORA, the power consumptions under all these schemes reach the same level, but the differences between their queue backlog sizes become more obvious: PORA ($V = 10^9$) reduces 96% of the queue backlog when compared with NOL and RANDOM. The results demonstrate that with the appropriate choice of the value of V , PORA can



(a) Total queue backlogs.



(b) Total power consumptions.

Fig. 9: Performance of PORA under imperfect prediction.

achieve less latency than the four baselines under the same power consumptions.

C. Evaluation with Imperfect Prediction

In practice, prediction errors are inevitable. Hence, we investigate the performance of PORA in the presence of prediction errors [28]. Particularly, we consider two kinds of prediction errors: false alarm and missed detection. A packet is falsely alarmed if it is predicted to arrive but it does not arrive actually. A packet is missed to be detected if it will arrive but is not predicted. We assume that all EFNs have the uniform false-alarm rate p_1 and missed-detection rate p_2 . In our simulation, we consider different pairs of values of (p_1, p_2) : $(0.0, 0.0)$, $(0.05, 0.05)$, $(0.5, 0.05)$, $(0.05, 0.25)$, and $(0.5, 0.25)$. Note that $(p_1, p_2) = (0.0, 0.0)$ corresponds to the case when the prediction is perfect.

Figure 9 presents the results under prediction window size $W = 10$. We observe when $V \leq 7.5 \times 10^{10}$, both the total queue backlog sizes and power consumptions under imperfect prediction are larger than that under perfect prediction. The reason for this performance degradation is twofold: First, arrivals that are missed to be detected cannot be pre-served, thus leading to larger queue backlog sizes. Second, PORA allocates redundant resources to handle the falsely predicted arrivals, thus causing more power consumptions. As the value of V increases, this performance degradation becomes negligible. Taking the total queue backlog under $(p_1, p_2) = (0.25, 0.5)$ as an example, when compared with the case under perfect prediction, it increases by 4.72% at $V = 10^{11}$, and increases by 2.24% at $V = 2 \times 10^{11}$. Moreover, there is no extra power

consumption under imperfect prediction when $V \geq 7.5 \times 10^{10}$ since PORA tends to reserve resources to reduce power consumptions under large V .

In summary, there will be performance degradation in both total queue backlog sizes and power consumptions in the presence of prediction errors. However, as the value of V increases, this degradation decreases and becomes negligible. Though a large value of V can improve the robustness of PORA and achieve small power consumptions, it brings long workload latencies. In practice, the choice of the value of V depends on how the system designer trades off all these criteria.

VII. CONCLUSION

In this paper, we studied the problem of dynamic offloading and resource allocation with prediction in a fog computing system with multiple tiers. By formulating it as a stochastic network optimization problem, we proposed PORA, an efficient online scheme that exploits predictive offloading to minimize power consumption with queue stability guarantee. Our theoretical analysis and trace-driven simulations showed that PORA achieves a tunable power-latency tradeoff, while effectively shortening latency with only mild-value of future information, even in the presence of prediction errors. As for future work, our model can be further extended to more general settings such that the instant wireless channel states may be unknown by the moment of decision making or the underlying system dynamics is non-stationary.

REFERENCES

- [1] Y. Xiao and M. Krunz, "QoE and power efficiency tradeoff for fog computing networks with fog node cooperation," in *Proceedings of IEEE INFOCOM*, 2017.
- [2] S. Yi, Z. Hao, Z. Qin, and Q. Li, "Fog computing: Platform and applications," in *Proceedings of IEEE HotWeb*, 2015.
- [3] J. Broughton, "Netflix adds download functionality," <https://technology.ihc.com/586280/netflix-adds-download-support>, 2016.
- [4] Y. Wang, X. Tao, X. Zhang, P. Zhang, and Y. T. Hou, "Cooperative task offloading in three-tier mobile computing networks: An ADMM framework," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 3, pp. 2763–2776, 2019.
- [5] L. Liu, Z. Chang, and X. Guo, "Socially-aware dynamic computation offloading scheme for fog computing system with energy harvesting devices," *IEEE Internet of Things Journal*, vol. 5, no. 3, pp. 1869–1879, 2018.
- [6] S. Misra and N. Saha, "Detour: Dynamic task offloading in software-defined fog for IoT applications," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 5, pp. 1159–1166, 2019.
- [7] L. Lei, H. Xu, X. Xiong, K. Zheng, and W. Xiang, "Joint computation offloading and multi-user scheduling using approximate dynamic programming in NB-IoT edge computing system," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 5345–5362, 2019.
- [8] Y. Mao, J. Zhang, S. Song, and K. B. Letaief, "Power-delay tradeoff in multi-user mobile-edge computing systems," in *Proceedings of IEEE GLOBECOM*, 2016.
- [9] Y. Chen, N. Zhang, Y. Zhang, X. Chen, W. Wu, and X. S. Shen, "Energy efficient dynamic offloading in mobile edge computing for Internet of Things," *IEEE Transactions on Cloud Computing*, 2019, doi: 10.1109/TCC.2019.2898657.
- [10] Y. Gao, W. Tang, M. Wu, P. Yang, and L. Dan, "Dynamic social-aware computation offloading for low-latency communications in IoT," *IEEE Internet of Things Journal*, 2019, doi: 10.1109/JIOT.2019.2909299.
- [11] D. Zhang, L. Tan, J. Ren, M. K. Awad, S. Zhang, Y. Zhang, and P.-J. Wan, "Near-optimal and truthful online auction for computation offloading in green edge-computing systems," *IEEE Transactions on Mobile Computing*, 2019, doi: 10.1109/TMC.2019.2901474.
- [12] M. Taneja and A. Davy, "Resource aware placement of IoT application modules in fog-cloud computing paradigm," in *Proceedings of IFIP/IEEE IM*, 2017.

- [13] M. Chen, W. Li, G. Fortino, Y. Hao, L. Hu, and I. Humar, "A dynamic service migration mechanism in edge cognitive computing," *ACM Transactions on Internet Technology*, vol. 19, no. 2, p. 30, 2019.
- [14] D. Zhang, Z. Chen, L. X. Cai, H. Zhou, S. Duan, J. Ren, X. Shen, and Y. Zhang, "Resource allocation for green cloud radio access networks with hybrid energy supplies," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 2, pp. 1684–1697, 2017.
- [15] N. K. Ahmed, A. F. Atiya, N. E. Gayar, and H. El-Shishiny, "An empirical comparison of machine learning models for time series forecasting," *Econometric Reviews*, vol. 29, no. 5-6, pp. 594–621, 2010.
- [16] L. Huang, S. Zhang, M. Chen, and X. Liu, "When backpressure meets predictive scheduling," *IEEE/ACM Transactions on Networking*, vol. 24, no. 4, pp. 2237–2250, 2016.
- [17] R. G. Gallager, *Principles of Digital Communication*. Cambridge University Press, 2008.
- [18] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.
- [19] A. P. Miettinen and J. K. Nurminen, "Energy efficiency of mobile clients in cloud computing," in *Proceedings of ACM HotCloud*, 2010.
- [20] Y. Kim, J. Kwak, and S. Chong, "Dual-side optimization for cost-delay tradeoff in mobile edge computing," *IEEE Transactions on Vehicular Technology*, vol. 67, no. 2, pp. 1765–1781, 2018.
- [21] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synthesis Lectures on Communication Networks*, vol. 3, no. 1, pp. 1–211, 2010.
- [22] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [23] A. Leon-Garcia, *Probability, Statistics, and Random Processes for Electrical Engineering*, 3rd ed. Pearson Education, 2017.
- [24] C.-F. Liu, M. Bennis, and H. V. Poor, "Latency and reliability-aware task offloading and resource allocation for mobile edge computing," in *Proceedings of IEEE GLOBECOM*, 2017.
- [25] J. Du, L. Zhao, J. Feng, and X. Chu, "Computation offloading and resource allocation in mixed fog/cloud computing systems with min-max fairness guarantee," *IEEE Transactions on Communications*, vol. 66, no. 4, pp. 1594–1608, 2017.
- [26] T. Benson, A. Akella, and D. A. Maltz, "Network traffic characteristics of data centers in the wild," in *Proceedings of ACM IMC*, 2010.
- [27] M. Mitzenmacher, "The power of two choices in randomized load balancing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 12, no. 10, pp. 1094–1104, 2001.
- [28] K. Chen and L. Huang, "Timely-throughput optimal scheduling with prediction," in *Proceedings of IEEE INFOCOM*, 2018.

APPENDIX A

DESIGN OF SCHEME PORA

First, we define Lyapunov function [21] $L(\mathbf{Q}(t))$ as

$$L(\mathbf{Q}(t)) \triangleq \frac{1}{2} \sum_{i \in \mathcal{N}} \sum_{\beta \in \{a, l, o\}} (Q_i^{(e, \beta)}(t))^2 + \frac{1}{2} \sum_{j \in \mathcal{M}} \sum_{\beta \in \{a, l, o\}} (Q_j^{(c, \beta)}(t))^2. \quad (29)$$

Next, we define the drift-plus-penalty $\Delta_V L(\mathbf{Q}(t))$ as

$$\Delta_V L(\mathbf{Q}(t)) \triangleq \mathbb{E}[L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) | \mathbf{Q}(t)] + V \mathbb{E}\{P(t) | \mathbf{Q}(t)\}, \quad (30)$$

where V is a positive parameter. According to definition (29), the update functions (6), (10), (12), (13), and (14), there exists a positive constant $\theta > 0$ such that

$$\begin{aligned} & L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) \\ & \leq \theta + \sum_{i \in \mathcal{N}} Q_i^{(e, a)}(t) (A_i(t + W_i) - b_i^{(e, l)}(t) - b_i^{(e, o)}(t)) \\ & \quad + \sum_{i \in \mathcal{N}} Q_i^{(e, l)}(t) (b_i^{(e, l)}(t) - \tau_0 f_i^{(e)}(t) / L_i^{(e)}) \\ & \quad + \sum_{i \in \mathcal{N}} Q_i^{(e, o)}(t) (b_i^{(e, o)}(t) - \sum_{j \in \mathcal{M}_i} R_{i, j}(t)) \end{aligned}$$

$$\begin{aligned} & + \sum_{j \in \mathcal{M}} Q_j^{(c, a)}(t) \left[\sum_{i \in \mathcal{N}_j} R_{i, j}(t) - b_j^{(c, l)}(t) - b_j^{(c, o)}(t) \right] \\ & + \sum_{j \in \mathcal{M}} Q_j^{(c, l)}(t) (b_j^{(c, l)}(t) - \tau_0 f_j^{(c)}(t) / L_j^{(c)}) \\ & + \sum_{j \in \mathcal{M}} Q_j^{(c, o)}(t) (b_j^{(c, o)}(t) - D_j(t)). \quad (31) \end{aligned}$$

Substituting (31) into the definition of drift-plus-penalty shown in (30) and by $\mathbb{E}[A_i(t + W_i)] = \lambda_i$, we obtain

$$\begin{aligned} & \Delta_V L(\mathbf{Q}(t)) \\ & \leq \theta + V \mathbb{E}\{P(t) | \mathbf{Q}(t)\} \\ & \quad + \sum_{i \in \mathcal{N}} Q_i^{(e, a)}(t) \mathbb{E}\left\{ \lambda_i - (b_i^{(e, l)}(t) + b_i^{(e, o)}(t)) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{i \in \mathcal{N}} Q_i^{(e, l)}(t) \mathbb{E}\left\{ b_i^{(e, l)}(t) - \tau_0 f_i^{(e)}(t) / L_i^{(e)} | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{i \in \mathcal{N}} Q_i^{(e, o)}(t) \mathbb{E}\left\{ b_i^{(e, o)}(t) - \sum_{j \in \mathcal{M}_i} R_{i, j}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} Q_j^{(c, a)}(t) \mathbb{E}\left\{ \sum_{i \in \mathcal{N}_j} R_{i, j}(t) \right. \\ & \quad \left. - (b_j^{(c, l)}(t) + b_j^{(c, o)}(t)) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} Q_j^{(c, l)}(t) \mathbb{E}\left\{ b_j^{(c, l)}(t) - \tau_0 f_j^{(c)}(t) / L_j^{(c)} | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} Q_j^{(c, o)}(t) \mathbb{E}\left\{ b_j^{(c, o)}(t) - D_j(t) | \mathbf{Q}(t) \right\}. \quad (32) \end{aligned}$$

Then by the expression of transmission capacity from EFN i to CFN j shown in (9) and the expression of total power consumptions shown in (16), we have

$$\begin{aligned} & \Delta_V L(\mathbf{Q}(t)) \\ & \leq \theta + \sum_{i \in \mathcal{N}} Q_i^{(e, a)}(t) \mathbb{E}\{A_i(t + W_i) | \mathbf{Q}(t)\} \\ & \quad + \sum_{i \in \mathcal{N}} \mathbb{E}\left\{ (Q_i^{(e, l)}(t) - Q_i^{(e, a)}(t)) b_i^{(e, l)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{i \in \mathcal{N}} \mathbb{E}\left\{ (Q_i^{(e, o)}(t) - Q_i^{(e, a)}(t)) b_i^{(e, o)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{i \in \mathcal{N}} \mathbb{E}\left\{ V \tau_0 \varsigma (f_i^{(e)}(t))^3 - \frac{\tau_0 Q_i^{(e, l)}(t)}{L_i^{(e)}} f_i^{(e)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} \mathbb{E}\left\{ (Q_j^{(c, l)}(t) - Q_j^{(c, a)}(t)) b_j^{(c, l)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} \mathbb{E}\left\{ (Q_j^{(c, o)}(t) - Q_j^{(c, a)}(t)) b_j^{(c, o)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{j \in \mathcal{M}} \mathbb{E}\left\{ V \tau_0 \varsigma (f_j^{(c)}(t))^3 - \frac{\tau_0 Q_j^{(c, l)}(t)}{L_j^{(c)}} f_j^{(c)}(t) | \mathbf{Q}(t) \right\} \\ & \quad + \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}_i} \mathbb{E}\left\{ V \tau_0 p_{i, j}(t) \right. \\ & \quad \left. - \tau_0 m_{i, j}(t) \log_2(1 + l_{i, j}(t) p_{i, j}(t)) | \mathbf{Q}(t) \right\} \\ & \quad - \sum_{j \in \mathcal{M}} Q_j^{(c, o)}(t) \mathbb{E}\{D_j(t) | \mathbf{Q}(t)\} \quad (33) \end{aligned}$$

where $m_{i,j}(t) \triangleq (Q_i^{(e,o)}(t) - Q_j^{(c,a)}(t))B$ and $l_{i,j}(t) \triangleq \frac{H_{i,j}(t)}{N_0B}$ for all $i \in \mathcal{N}, j \in \mathcal{M}_i$.

To solve problem (19), we should minimize the upper bound of $\Delta_V L(\mathbf{Q}(t))$ in every time slot. However, it is hard to solve a minimization problem with expectation. Thus we approximately solve the problem by considering the following deterministic problem in every time slot t :

$$\begin{aligned}
& \text{Minimize} \sum_{\mathbf{b}, \mathbf{f}, \mathbf{p}} \sum_{i \in \mathcal{N}} \left(Q_i^{(e,l)}(t) - Q_i^{(e,a)}(t) \right) b_i^{(e,l)} \\
& + \sum_{i \in \mathcal{N}} \left(Q_i^{(e,o)}(t) - Q_i^{(e,a)}(t) \right) b_i^{(e,o)} \\
& + \sum_{i \in \mathcal{N}} \left(V\tau_0 \zeta \left(f_i^{(e)} \right)^3 - \frac{\tau_0 Q_i^{(e,l)}(t)}{L_i^{(e)}} f_i^{(e)} \right) \\
& + \sum_{j \in \mathcal{M}} \left(Q_j^{(c,l)}(t) - Q_j^{(c,a)}(t) \right) b_j^{(c,l)} \\
& + \sum_{j \in \mathcal{M}} \left(Q_j^{(c,o)}(t) - Q_j^{(c,a)}(t) \right) b_j^{(c,o)} \quad (34) \\
& + \sum_{j \in \mathcal{M}} \left(V\tau_0 \zeta \left(f_j^{(c)} \right)^3 - \frac{\tau_0 Q_j^{(c,l)}(t)}{L_j^{(c)}} f_j^{(c)} \right) \\
& + \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \left[V\tau_0 p_{i,j} - \tau_0 m_{i,j}(t) \log_2(1 + l_{i,j}(t) p_{i,j}) \right. \\
& \left. + l_{i,j}(t) p_{i,j} \right]
\end{aligned}$$

Subject to (1)(7)(8)(11)(15).

Problem (34) can be decomposed into subproblems shown in Section V. By solving these subproblems, we develop PORA, an online scheme that independently makes predictive offloading decisions $\mathbf{b}(t)$, sets CPU frequencies $\mathbf{f}(t)$, and allocates transmit powers $\mathbf{p}(t)$ in every time slot t . ■

APPENDIX B

PROOF OF OPTIMAL LOCAL CPU FREQUENCY

To solve the optimal solution to subproblem (22), we denote its objective function by

$$F_k^{(\alpha,t)}(f_k^{(\alpha)}) \triangleq V\zeta \left(f_k^{(\alpha)} \right)^3 - \frac{Q_k^{(\alpha,l)}(t)}{L_k^{(\alpha)}} f_k^{(\alpha)}. \quad (35)$$

Its first- and second-order derivatives are shown as follows:

$$\frac{dF_k^{(\alpha,t)}(f_k^{(\alpha)})}{df_k^{(\alpha)}} = 3V\zeta \left(f_k^{(\alpha)} \right)^2 - \frac{Q_k^{(\alpha,l)}(t)}{L_k^{(\alpha)}}, \quad (36)$$

$$\frac{d^2 F_k^{(\alpha,t)}(f_k^{(\alpha)})}{(df_k^{(\alpha)})^2} = 6V\zeta f_k^{(\alpha)}. \quad (37)$$

From the above two derivatives, we conclude that function $F_k^{(\alpha,t)}(\cdot)$ is convex in interval $[0, f_{k,\max}^{(\alpha)}]$ since its second order derivative satisfies $d^2 F_k^{(\alpha,t)}(\cdot)/(df_k^{(\alpha)})^2 \geq 0$ for $f_k^{(\alpha)} \geq 0$. On the other hand, its first order derivative satisfies $dF_k^{(\alpha,t)}(\cdot)/df_k^{(\alpha)} = 0$ when $f_k^{(\alpha)} = \sqrt{Q_k^{(\alpha,l)}(t)/3V\zeta L_k^{(\alpha)}}$. Thus the minimum point of $F_i^{(\alpha,t)}(\cdot)$ over interval $[0, f_{k,\max}^{(\alpha)}]$ is $\min \left\{ \sqrt{Q_k^{(\alpha,l)}(t)/3V\zeta L_k^{(\alpha)}}, f_{k,\max}^{(\alpha)} \right\}$. ■

APPENDIX C

PROOF OF OPTIMAL TRANSMIT POWER ALLOCATION

We denote the optimal solution to subproblem (24) by $\mathbf{p}_i^*(t)$ and the objective function in subproblem (24) by $G_i^{(t)}(\mathbf{p}_i)$. Moreover, we define the following function

$$G_{i,j}^{(t)}(p_{i,j}) \triangleq Vp_{i,j} - m_{i,j}(t) \log_2(1 + l_{i,j}(t) p_{i,j}) \quad (38)$$

for each $j \in \mathcal{M}_i$. Then $G_i^{(t)}(\mathbf{p}_i)$ can be expressed as

$$G_i^{(t)}(\mathbf{p}_i) = \sum_{j \in \mathcal{M}_i} G_{i,j}^{(t)}(p_{i,j}). \quad (39)$$

We denote the minimizer of function $G_{i,j}^{(t)}(\cdot)$ in interval $[0, \infty)$ by $\tilde{p}_{i,j}^{(t)}$, i.e.,

$$\tilde{p}_{i,j}^{(t)} \triangleq \arg \min_{p_{i,j} \geq 0} G_{i,j}^{(t)}(p_{i,j}). \quad (40)$$

When $m_{i,j}(t) \leq 0$, $G_{i,j}^{(t)}(\cdot)$ is increasing over interval $[0, \infty)$ and $\tilde{p}_{i,j}^{(t)} = 0$. In this case, we have $\mathbf{p}_i^*(t) = \tilde{\mathbf{p}}_i^{(t)}$. When $m_{i,j}(t) > 0$, $G_{i,j}^{(t)}(\cdot)$ is convex in interval $[0, \infty)$ since its second-order derivative satisfies

$$\frac{d^2 G_{i,j}^{(t)}(p_{i,j})}{dp_{i,j}^2} = \frac{m_{i,j}(t) (l_{i,j}(t))^2}{(1 + l_{i,j}(t) p_{i,j})^2} > 0. \quad (41)$$

Thus we obtain $\tilde{p}_{i,j}^{(t)}$ by letting its first-order derivative to be zero:

$$\frac{dG_{i,j}^{(t)}(p_{i,j})}{dp_{i,j}} \Big|_{p_{i,j} = \tilde{p}_{i,j}^{(t)}} = V - \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) \tilde{p}_{i,j}^{(t)}} = 0. \quad (42)$$

It follows that when $m_{i,j}(t) > 0$,

$$\tilde{p}_{i,j}^{(t)} = \left[\frac{m_{i,j}(t)}{V} - \frac{1}{l_{i,j}(t)} \right]^+. \quad (43)$$

If $\sum_{j \in \mathcal{M}_i} \tilde{p}_{i,j}^{(t)} \leq p_{i,\max}$, we have $\mathbf{p}_i^*(t) = \tilde{\mathbf{p}}_i^{(t)}$ as the constraints in (24) are satisfied. Otherwise, we have the following lemma.

Lemma 1: If $\sum_{j \in \mathcal{M}_i} \tilde{p}_{i,j}^{(t)} > p_{i,\max}$, then $\mathbf{p}_i^*(t)$ must satisfy $\sum_{j \in \mathcal{M}_i} p_{i,j}^*(t) = p_{i,\max}$.

Proof: We prove Lemma 1 by contradiction. Suppose that there exists $\theta_1 > 0$ such that $\sum_{j \in \mathcal{M}_i} p_{i,j}^*(t) + \theta_1 = p_{i,\max}$. Since $\sum_{j \in \mathcal{M}_i} \tilde{p}_{i,j}^{(t)} > p_{i,\max}$, there exist $j' \in \mathcal{M}_i$ and $\theta_2 > 0$ such that $p_{i,j'}^*(t) < \tilde{p}_{i,j'}^{(t)} - \theta_2$. Note that $m_{i,j'}(t) > 0$ must hold for j' since $\tilde{p}_{i,j'}^{(t)} > 0$. Now we consider a solution $\mathbf{p}_i^0(t)$ to subproblem (24) which satisfies

$$\begin{aligned}
p_{i,j'}^0(t) &= p_{i,j'}^*(t) + \theta_3, \\
p_{i,j}^0(t) &= p_{i,j}^*(t), \quad \forall j \in \mathcal{M}_i/j',
\end{aligned} \quad (44)$$

where $\theta_3 \in (0, \min(\theta_1, \theta_2)]$. Then $\mathbf{p}_i^0(t)$ is a feasible solution since

$$\begin{aligned}
\sum_{j \in \mathcal{M}_i} p_{i,j}^0(t) &= \sum_{j \in \mathcal{M}_i} p_{i,j}^*(t) + \theta_3 \\
&\leq \sum_{j \in \mathcal{M}_i} p_{i,j}^*(t) + \theta_1 = p_{i,\max}.
\end{aligned} \quad (45)$$

By the definition of $p_{i,j'}^0(t)$ in (44), we have

$$p_{i,j'}^*(t) < p_{i,j'}^0(t) < \tilde{p}_{i,j'}^{(t)}. \quad (46)$$

Since $G_{i,j'}^{(t)}(\cdot)$ is convex and $\tilde{p}_{i,j'}^{(t)}$ is its unique minimizer, we have

$$G_{i,j'}^{(t)}(p_{i,j'}^*(t)) > G_{i,j'}^{(t)}(p_{i,j'}^0(t)) > G_{i,j'}^{(t)}(\tilde{p}_{i,j'}^{(t)}). \quad (47)$$

It follows that

$$G_i^{(t)}(\mathbf{p}_i^*(t)) > G_i^{(t)}(\mathbf{p}_i^0(t)), \quad (48)$$

which contradicts the fact that $\mathbf{p}_i^*(t)$ is the optimal solution to (24). Thus θ_1 must equal zero and $\mathbf{p}_i^*(t)$ satisfies $\sum_{j \in \mathcal{M}_i} p_{i,j}^*(t) = p_{i,\max}$. ■

When $\sum_{j \in \mathcal{M}_i} \tilde{p}_{i,j}^{(t)} > p_{i,\max}$, to find the optimal solution to problem (24), we need the following lemma as well.

Lemma 2: For any $j \in \mathcal{M}_i$, if $m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$, then $p_{i,j}^*(t) = \tilde{p}_{i,j}^{(t)} = 0$.

Proof: By (43), $\tilde{p}_{i,j}^{(t)} = 0$ if and only if $m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$. Next, we show that if $m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$, then the optimal $p_{i,j}^*(t)$ must be zero. Particularly, we prove it by contradiction.

Assume the optimal $p_{i,j}^*(t) > 0$, then there must exist a feasible solution $\mathbf{p}_i^1(t)$ such that $p_{i,j'}^1(t) = p_{i,j'}^*(t)$ for all $j' \in \mathcal{M}_i/j$ and $p_{i,j}^1(t) = 0 < p_{i,j}^*(t)$. Then we have

$$G_i^{(t)}(\mathbf{p}_i^*(t)) - G_i^{(t)}(\mathbf{p}_i^1(t)) = V p_{i,j}^*(t) - m_{i,j}(t) \log_2(1 + l_{i,j}(t) p_{i,j}^*(t)). \quad (49)$$

If $m_{i,j}(t) \leq 0$, according to $p_{i,j}^*(t) > 0$, we have

$$G_i^{(t)}(\mathbf{p}_i^*(t)) - G_i^{(t)}(\mathbf{p}_i^1(t)) > 0. \quad (50)$$

If $0 < m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$, since $\tilde{p}_{i,j}^{(t)} = 0 < p_{i,j}^*(t)$ is the unique minimizer of $G_{i,j}^{(h)}(\cdot)$ over $[0, \infty)$, we have

$$G_i^{(t)}(\mathbf{p}_i^*(t)) - G_i^{(t)}(\mathbf{p}_i^1(t)) = G_{i,j}^{(t)}(p_{i,j}^*(t)) > G_{i,j}^{(t)}(\tilde{p}_{i,j}^{(t)}), \quad (51)$$

which contradicts the fact that $\mathbf{p}_i^*(t)$ is the optimal solution of problem (24). Thus for any j with $m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$, the optimal $p_{i,j}^*(t)$ must be zero. ■

We define $\mathcal{M}_i^+ \triangleq \{j | j \in \mathcal{M}_i, m_{i,j}(t) > \frac{V}{l_{i,j}(t)}\}$. By applying Lemma 1 and Lemma 2, when $\sum_{j \in \mathcal{M}_i} \tilde{p}_{i,j}^{(t)} > p_{i,\max}$, we just need to solve the following problem:

$$\begin{aligned} & \text{Minimize} \sum_{(p_{i,j})_{j \in \mathcal{M}_i^+}} \sum_{j \in \mathcal{M}_i^+} \left[V p_{i,j} - m_{i,j}(t) \log_2(1 + l_{i,j}(t) p_{i,j}) \right] \\ & \text{Subject to} \sum_{j \in \mathcal{M}_i^+} p_{i,j} = P_{i,\max}, \\ & p_{i,j} \geq 0, \forall j \in \mathcal{M}_i^+. \end{aligned} \quad (52)$$

Note that $(p_{i,j}^*(t))_{j \in \mathcal{M}_i^+}$ is the optimal solution to problem (52) and it satisfies the following KKT conditions:

$$V - \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} + \lambda^* - \mu_j^* = 0, \forall j \in \mathcal{M}_i^+, \quad (53)$$

$$\mu_j^* p_{i,j}^*(t) = 0, \forall j \in \mathcal{M}_i^+, \quad (54)$$

$$\lambda^*, \mu_j^* \geq 0, \forall j \in \mathcal{M}_i^+, \quad (55)$$

$$\sum_{j \in \mathcal{M}_i^+} p_{i,j}^*(t) = p_{i,\max}, \quad (56)$$

$$p_{i,j}^*(t) \geq 0, \quad (57)$$

where λ^* and $(\mu_j^*)_{j \in \mathcal{M}_i^+}$ are the corresponding optimal dual variables. Multiplying both sides of (53) by $p_{i,j}^*(t)$, we have

$$\left(V - \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} + \lambda^* \right) p_{i,j}^*(t) - \mu_j^* p_{i,j}^*(t) = 0. \quad (58)$$

It follows by (54) that

$$\left(V - \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} + \lambda^* \right) p_{i,j}^*(t) = 0. \quad (59)$$

On the other hand, according to (53) and (55), we have

$$\begin{aligned} \lambda^* &= \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} - V + \mu_j^* \\ &\geq \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} - V \end{aligned} \quad (60)$$

for every $j \in \mathcal{M}_i^+$. Now we consider two cases:

1) If $\lambda^* < m_{i,j}(t) l_{i,j}(t) - V$, then (60) holds only if $p_{i,j}^*(t) > 0$. It follows by (59) that

$$\lambda^* = \frac{m_{i,j}(t) l_{i,j}(t)}{1 + l_{i,j}(t) p_{i,j}^*(t)} - V, \quad (61)$$

which yields $p_{i,j}^*(t) = \frac{m_{i,j}(t)}{V + \lambda^*} - \frac{1}{l_{i,j}(t)}$.

2) If $\lambda^* \geq m_{i,j}(t) l_{i,j}(t) - V$, then condition (59) holds if and only if $p_{i,j}^*(t) = 0$.

In conclusion, we have

$$p_{i,j}^*(t) = \begin{cases} \frac{m_{i,j}(t)}{V + \lambda^*} - \frac{1}{l_{i,j}(t)}, & \text{if } \lambda^* < m_{i,j}(t) l_{i,j}(t) - V, \\ 0, & \text{if } \lambda^* \geq m_{i,j}(t) l_{i,j}(t) - V, \end{cases} \quad (62)$$

or equivalently,

$$p_{i,j}^*(t) = \left[\frac{m_{i,j}(t)}{V + \lambda^*} - \frac{1}{l_{i,j}(t)} \right]^+. \quad (63)$$

Note that the above expression also applies to the case when $m_{i,j}(t) \leq \frac{V}{l_{i,j}(t)}$. Then by substituting (63) into (56), we obtain

$$\sum_{j \in \mathcal{M}_i} \left[\frac{m_{i,j}(t)}{V + \lambda^*} - \frac{1}{l_{i,j}(t)} \right]^+ = p_{i,\max}. \quad (64)$$

The left-hand side is a piecewise-linear decreasing function of λ^* , with the breakpoint at $(m_{i,j}(t) l_{i,j}(t) - V)$. Therefore, the equation has a unique solution. ■

APPENDIX D
PROOF OF THEOREM 2

Applying the Corollary 1 in [16], given prediction window size W_i , the average latency of workload in arrival queue $A_{i,-1}(t)$ of EFN i under PORA is

$$d_i^p = \sum_{w \geq 1} w \pi_{i,w+W_i}. \quad (65)$$

According to Little's theorem [23], the average arrival queue backlog size of EFN i under prediction is

$$\psi_i^p = \lambda_i d_i^p = \lambda_i \sum_{w \geq 1} w \pi_{i,w+W_i}. \quad (66)$$

Therefore, the total average arrival queue backlog sizes of all EFNs is

$$\psi^p = \sum_{i \in \mathcal{N}} \psi_i^p = \sum_{i \in \mathcal{N}} \lambda_i \sum_{w \geq 1} w \pi_{i,w+W_i}. \quad (67)$$

When the prediction window size is zero, *i.e.*, when there is no prediction, the corresponding total average arrival queue backlog size of all EFNs is

$$\psi = \sum_{i \in \mathcal{N}} \psi_i = \sum_{i \in \mathcal{N}} \lambda_i \sum_{w \geq 1} w \pi_{i,w}. \quad (68)$$

Using (67) and (68), we conclude that

$$\begin{aligned} \psi - \psi^p &= \sum_{i \in \mathcal{N}} \lambda_i \left(\sum_{w \geq 1} w \pi_{i,w} - \sum_{w \geq 1} w \pi_{i,w+W_i} \right) \\ &= \sum_{i \in \mathcal{N}} \lambda_i \left(\sum_{w \geq 1} w \pi_{i,w} - \sum_{w \geq 1} (w + W_i) \pi_{i,w+W_i} \right. \\ &\quad \left. + \sum_{w \geq 1} W_i \pi_{i,w+W_i} \right) \\ &= \sum_{i \in \mathcal{N}} \lambda_i \left(\sum_{w \geq 1} w \pi_{i,w} - \sum_{w \geq W_i+1} w \pi_{i,w} \right. \\ &\quad \left. + \sum_{w \geq 1} W_i \pi_{i,w+W_i} \right) \\ &= \sum_{i \in \mathcal{N}} \lambda_i \left(\sum_{1 \leq w \leq W_i} w \pi_{i,w} + W_i \sum_{w \geq 1} \pi_{i,w+W_i} \right). \quad (69) \end{aligned}$$

Dividing both sides by $\sum_{i \in \mathcal{N}} \lambda_i$ and using Little's theorem, we obtain (27).

Next, we prove (28). Taking the limit of \mathbf{W} ($\mathbf{W} \rightarrow \infty$), we obtain

$$\lim_{\mathbf{W} \rightarrow \infty} \sum_{i \in \mathcal{N}} \lambda_i \sum_{1 \leq w \leq W_i} w \pi_{i,w} = \psi. \quad (70)$$

It follows that

$$\lim_{\mathbf{W} \rightarrow \infty} \eta(\mathbf{W}) = d + \lim_{\mathbf{W} \rightarrow \infty} \frac{\sum_{i \in \mathcal{N}} \lambda_i W_i \sum_{w \geq 1} \pi_{i,w+W_i}}{\sum_{i \in \mathcal{N}} \lambda_i}. \quad (71)$$

On the other hand, we have

$$\begin{aligned} \lim_{\mathbf{W} \rightarrow \infty} \eta(\mathbf{W}) &= \frac{\psi}{\sum_{i \in \mathcal{N}} \lambda_i} - \lim_{\mathbf{W} \rightarrow \infty} \frac{\psi^p}{\sum_{i \in \mathcal{N}} \lambda_i} \\ &\leq \frac{\psi}{\sum_{i \in \mathcal{N}} \lambda_i}. \quad (72) \end{aligned}$$

Combining (71) and (72), we have

$$\lim_{\mathbf{W} \rightarrow \infty} \frac{\sum_{i \in \mathcal{N}} \lambda_i W_i \sum_{w \geq 1} \pi_{i,w+W_i}}{\sum_{i \in \mathcal{N}} \lambda_i} = 0, \quad (73)$$

since it cannot be negative. Substituting (73) into (71), we obtain

$$\lim_{\mathbf{W} \rightarrow \infty} \eta(\mathbf{W}) = d. \quad (74)$$

■