

When Backpressure Meets Predictive Scheduling

Longbo Huang^{*}, Shaoquan Zhang[#], Minghua Chen[#], Xin Liu[†]
^{*}longbohuang@tsinghua.edu.cn, IIS, Tsinghua University
[#]{zsq008, minghua}@ie.cuhk.edu.hk, IE, Chinese University of Hong Kong
[†]liuxin@microsoft.com, Microsoft Research Asia

ABSTRACT

Motivated by the increasing popularity of learning and predicting human user behavior in communication and computing systems, in this paper, we investigate the fundamental benefit of predictive scheduling, i.e., predicting and pre-serving arrivals, in controlled queueing systems. Based on a lookahead-window prediction model, we first establish a novel queue-equivalence between the predictive queueing system with a *fully-efficient* scheduling scheme and an equivalent queueing system without prediction. This result allows us to analytically demonstrate that predictive scheduling necessarily improves system delay performance and drives it to zero with increasing prediction power. It also enables us to exactly determine the required prediction power for different systems and study its impact on tail delay.

We then propose the **Predictive Backpressure (PBP)** algorithm for achieving optimal utility performance in such predictive systems. PBP efficiently incorporates prediction into stochastic system control and avoids the great complication due to the exponential state space growth in the prediction window size. We show that PBP achieves a utility performance that is within $O(\epsilon)$ of the optimal, for any $\epsilon > 0$, while guaranteeing that the system delay distribution is a *shifted-to-the-left* version of that under the original Backpressure algorithm. Hence, the average delay under PBP is strictly better than that under Backpressure, and vanishes with increasing prediction window size. This implies that the resulting utility-delay tradeoff with predictive scheduling can beat the known optimal $[O(\epsilon), O(\log(1/\epsilon))]$ tradeoff for systems without prediction.

1. INTRODUCTION

Due to the rapid development of powerful handheld devices, e.g., smartphones or tablet computers, human users now interact much more easily and frequently with the communication and computing infrastructures, e.g., E-commerce websites, cellular networks, and crowdsourcing platforms. Thus, in order to provide high level quality-of-service, it is

important to understand human behavior features and to utilize such information in guiding system control algorithm design. Therefore, various studies have been conducted to learn and predict human behavior patterns, e.g., online social networking [1], online searching behavior [2], and online browsing [3].

In this paper, we take one step further and ask the following important question: *What is the fundamental system benefit of having such user-behavior information?* Our objective is to obtain a theoretical quantification of this gain. To mathematically carry out our investigation, we consider a multi-user single-server queueing system. At every time, user workload arriving at the system will first be queued at corresponding buffer space. Then, the server allocates resources and decides the scheduling for serving the jobs. These operations allow the server to serve certain amount of workload for each user, but also result in a system cost due to resource utilization. Different from most existing work in multi-queue systems, here we assume that the server can *predict and serve future arrivals before they arrive at the system*. Hence, at every time, the server updates his prediction of future arrivals and adapts his control action. The objective is to serve all user workload with minimum cost, and to ensure small job latency for each user.

This is an important problem and can be used to model many practical systems where traffic prediction and pre-serving can be performed. The first example is scheduling in cellular networks. In this case, the base station handles users' data demand. Instead of waiting for the users to submit their requests, and suffering from a potentially big burst of traffic, which can lead to a large service latency, the base station can "push" the information to the users beforehand, e.g., push the news information at 7am in the morning. The second example scenario is prefetching in computing systems, e.g., [4] [5]. Here, data or instructions are preloaded into memory before they are actually requested. Doing so enables faster access or execution of the commands and enhances system performance. Another example is computing management, e.g., in computers. In this case, each user represents a software application and the server represents a workload management unit. Then, according to the needs of the applications, the managing unit pre-computes some information in case some later applications request them, e.g., branch prediction in computer architecture [6] [7].

There have been many previous works studying multi-queue system scheduling with utility optimization. [8] studies the fundamental tradeoff between energy consumption and packet delay for a single-queue system. [9] extends

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.
MobiHoc '14, August 11–14, 2014, Philadelphia, PA, USA.
Copyright 2014 ACM 978-1-4503-2620-9/14/08 ...\$15.00.
<http://dx.doi.org/10.1145/2632951.2632983>.

the results to a downlink system and designs algorithms to achieve the optimal tradeoff. [10] designs algorithms for minimizing energy consumption of a stochastic network. [11] designs energy optimal scheme for satellites. [12] looks at the problem of quality-of-service guaranteed energy efficient transmission using a calculus approach. [13] studies the tradeoff between energy and robustness for downlink systems. [14] and [15] develop algorithms for achieving the optimal utility-delay tradeoff in multihop networks.

However, we note that all the aforementioned works assume that the system only takes *causal* scheduling actions, i.e., the server will start serving packets only after they enter the system. While this is necessary in many systems, preserving future traffic can actually be done in systems that have highly predictable traffic. While predictive scheduling approaches have been investigated, e.g., [7], not much analytical study has been conducted. Closest to our work are [16], [17], which study the benefit of proactive scheduling, and [18], which studies the impact of future arrival information on queueing delay in $M/M/1$ queues. However, we note that [16] and [17] do not consider the effect of queueing, which very commonly appears in communication and computing systems, whereas [18] does not consider controllable rates and scheduling. Indeed, due to the joint existence of prediction and controlled queueing, the problem considered here is very complicated. Delay problems in stochastic controlled queueing systems are known to be hard. Moreover, arrival prediction advances in a sliding-window pattern over time, i.e., at every time, the system can predict slightly further into the future. Designing control algorithms for such systems often involves dynamic programming (DP). However, since the state space size grows exponentially with the prediction window size, the DP approach may not be computationally practical even for small systems. Even without prediction, the complexity of DP can still be very high due to the large queue state space. Moreover, since the system prediction evolves according to a sliding-window pattern, it is also not possible to apply the frame-based Lyapunov technique as in [19] and [20].

To resolve the above difficulties, we first establish a novel equivalence between the queueing system under prediction and a class of *fully-efficient* scheduling scheme and a queueing system without prediction but with a different initial condition and an equivalent scheduling policy. This connection is made by carrying out a sample-path queueing argument and enables us to analytically quantify the delay gain due to predictive scheduling for general multi-queue single-server systems. Our result shows that for such systems, the packet delay distribution is *shifted-to-the-left* under predictive scheduling. Hence, the average delay necessarily decreases and *approaches zero* as the prediction window size increases. Based on this result, we further propose a low-complexity **Predictive Backpressure** (PBP) scheduling policy for utility maximization in such predictive systems. PBP retains all desired features of the original Backpressure algorithm [21], e.g., greedy, does not require statistical information of the system dynamics, and has strong theoretical performance guarantee. We prove that the PBP algorithm can achieve an average cost that is $O(\epsilon)$ of the minimum cost for any $\epsilon > 0$, while guaranteeing an average delay that is strictly smaller than that under the original Backpressure. Hence, the resulting utility-delay tradeoff with predictive scheduling can beat the known optimal $[O(\epsilon), O(\log(1/\epsilon))]$

tradeoff for systems without prediction. We also demonstrate analytically and numerically with real data trace that when the first-in-first-out (FIFO) queueing discipline is used, PBP achieves an average packet delay reduction that is linear in the prediction window size, and that when the last-in-first-out (LIFO) discipline is used, the average packet delay under PBP decreases *exponentially* in the window size. These results demonstrate the power of predictive scheduling and provide explicit quantification of the benefits, which also provides useful guidelines for predictive algorithm design.

The rest of the paper is organized as follows. In Section 2, we present our system model and problem formulation. We develop the **Predictive Backpressure** (PBP) algorithm in Section 3. The analysis of delay performance under general predictive scheduling and PBP is given in Section 4. We extend the results to predictable-only systems in Section 5. Simulation results are presented in Section 6, followed by the conclusions in Section 7.

2. SYSTEM MODEL

We consider a general multi-queue single-server system shown in Fig. 1. In this system, a server serves N queues, one for each user that utilizes the service of the server.¹ This multi-queue system has many applications. For instance, it can be used to model downlink transmission in cellular networks, where the server represents the base station and the users are mobile users. Another example is a task management system of smartphones, where each user represents an application and the server represents the operating system that manages all computing workloads. We assume that the system operates in slotted time, i.e., $t \in \{0, 1, \dots\}$.

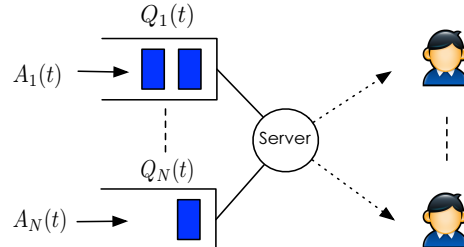


Figure 1: A multi-queue system where a server is serving workloads for different users/applications.

2.1 The Traffic Model

We use $A_n(t)$ to denote the amount of new workload arriving at the system at time t (called packets below). Here the workload can represent newly arrived data units that need to be delivered to their destinations, or new computing tasks that the server must fulfill eventually. We use $\mathbf{A}(t) = (A_1(t), \dots, A_N(t))$ to denote the vector of arrivals at time t . We assume that $\mathbf{A}(t)$ is i.i.d. with $\mathbb{E}\{A_n(t)\} = \lambda_n$.² We also assume that for each n , $0 \leq A_n(t) \leq A_{\max}$.

2.2 The Service Rate Model

¹Though we only consider a multi-queue system here, our results can likely be extended to general multihop queueing systems using the queue-equivalent results in Section 4.

²The arrivals can be arbitrarily correlated among different users.

Every time slot, the server allocates power for serving the pending packets.³ However, due to the potential system dynamics, e.g., channel fading coefficient changes, serving different users at different time may result in different resource consumption and generate different service rates. We model this fact by assuming that the server connects to each user n with a time-varying channel, whose state is denoted by $S_n(t)$. We then denote $\mathbf{S}(t) = (S_1(t), \dots, S_N(t))$ the system link state. We assume that $\mathbf{S}(t)$ is i.i.d. and takes values in $\{s_1, \dots, s_K\}$.⁴ We use π_{s_i} to denote the probability that $\mathbf{S}(t) = s_i$.

The server's power allocation over link n at time t is denoted by $P_n(t)$. We denote the aggregate system power allocation vector by $\mathbf{P}(t) = (P_1(t), \dots, P_N(t))$. Under a system link state s_i , we assume that the power allocation vector $\mathbf{P}(t)$ must be chosen from some feasible power allocation set $\mathcal{P}^{(s_i)}$, which is compact and contains the constraint $0 \leq P_n(t) \leq P_{\max}$. Then, under the given link state $\mathbf{S}(t)$ and the power allocation vector $\mathbf{P}(t)$, the amount of backlog that can be served for user n is determined by $\mu_n(t) = \mu_n(\mathbf{S}(t), \mathbf{P}(t))$. We assume that $\mu_n(\mathbf{S}(t), \mathbf{P}(t))$ is a continuous function of $\mathbf{P}(t)$ for all $\mathbf{S}(t)$. Also, we assume that there exists μ_{\max} such that $\mu_n(\mathbf{S}(t), \mathbf{P}(t)) \leq \mu_{\max}$ for all n , all time t , and under any $\mathbf{S}(t)$ and $\mathbf{P}(t)$.

2.3 The Predictive Service Model

Different from most previous works, we assume that the server can *predict and serve* future packet arrivals. Specifically, we first parameterize our prediction model by a vector $\mathbf{D} = (D_1, \dots, D_N)$, where $D_n \geq 1$ is the prediction window size of user n . That is, at each time t , the server has access to the arrival information in the *lookahead window* $\{A_n(t), \dots, A_n(t + D_n - 1)\}$, and can allocate rates to serve the future arrivals in the current time slot.⁵ Such a lookahead window model approximates practical scenarios and was also used in [18] and [23].

We then use $\{\mu_n^{(d)}(t)\}_{d=0}^{D_n-1}$ to denote the rate allocated to serving the arriving packets in time slot $t + d$ and let $\mu_n^{(-1)}(t)$ denote the rate allocated for serving the packets that are already in the system. Note that we always have $\sum_{d=-1}^{D_n-1} \mu_n^{(d)}(t) \leq \mu_n(t)$. Fig. 2 shows the slot structure and the predictive service model.

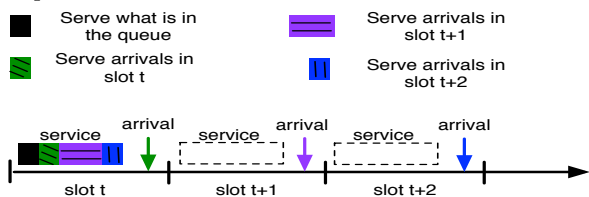


Figure 2: The figure shows what happens in a single time slot in the case of $D_n = 3$. The server predicts and serves the arrivals in timeslots t , $t + 1$ and $t + 2$, respectively.

³Our results can be extended to the case where the server consumes multiple types of resources, e.g., power and CPU cycles.

⁴Our results can easily be generalized to the case when both arrivals and channel conditions are Markovian, using the variable-size drift analysis developed in [22].

⁵Since we assume that the arrivals in a time slot can only be served in the next slot, we also consider $A_n(t)$ to be future arrivals.

2.4 Queueing

Denote $Q_n(t)$ the number of packets queued at the server for user n . We assume the following queueing dynamics:

$$Q_n(t+1) = \left[Q_n(t) - \mu_n^{(-1)}(t) \right]^+ + A_n^{(-1)}(t). \quad (1)$$

Here $A_n^{(-1)}(t)$ denotes the number of packets that *actually* enter the queue after going through a series of predictive service phases, i.e., for all $-1 \leq d \leq D_n - 2$,

$$A_n^{(d)}(t) = [A_n^{(d+1)}(t) - \mu_n^{(d+1)}(t - d - 1)]^+, \quad (2)$$

and $A_n^{(D_n-1)}(t) = A_n(t)$. In this paper, we say that the system is *stable* if the following condition holds:

$$Q_{\text{av}} \triangleq \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \sum_n \mathbb{E}\{Q_n(\tau)\} < \infty. \quad (3)$$

2.5 System Objective

In every time slot, the server spends certain cost due to power expenditure. We denote this cost by $f(\mathbf{S}(t), \mathbf{P}(t))$. One simple example is $f(\mathbf{S}(t), \mathbf{P}(t)) = \sum_n P_n(t)$, which denotes the total power consumption. We assume that under any state $\mathbf{S}(t)$, there exists a constant f^{\max} such that $f(\mathbf{S}(t), \mathbf{P}(t)) \leq f^{\max}$. The special case when $f(\mathbf{S}(t), \mathbf{P}(t))$ is independent of $\mathbf{P}(t)$ corresponds to the stability scheduling problem [21].

The system's objective is to find a power allocation and scheduling scheme for minimizing the time average cost, defined as:

$$f_{\text{av}} \triangleq \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \mathbb{E}\{f(\tau)\}, \quad (4)$$

subject to the constraint that the queues in the system must be stable, i.e., (3) holds. We use $f_{\text{av}}^{D_n*}$ to denote the minimum average cost under any feasible predictive scheduling algorithm with prediction vector \mathbf{D} , i.e., those that predict the arrivals for D_n slots and allocates service rates to serving the arrivals within the window $[t, t + D_n - 1]$ for each user n . We then use f_{av}^* to denote the minimum average power consumption incurred under any non-predictive scheduling policy.

2.6 Discussion of the Model

Note that the lookahead-window model is an idealized model which assumes that the system can perfectly predict future arrivals. Because of this, our results can be viewed as upper bounds of the fundamental benefit of predictive scheduling, which provide important criteria for evaluating predictive control algorithms. We also investigate the impact of prediction error in Section 6.

We also note that our model is very different from previous controlled queueing system works, which almost all assume that the system operates in a *causal* manner, i.e., only serves packets after they arrive at the system. Our model is motivated by pre-fetching techniques used in memory management [4], branch prediction in computer architecture [6], as well as recent advancement in data mining for learning user behavior patterns [3].

Our model is most relevant for modeling problems where future workload can be predicted and served before they enter the system. One such application scenario is in mobile networks, where the base station handles users' demand.

Since each user typically requests certain news information at specific times, e.g., 7am in the morning. Instead of waiting for the all the users to submit their requests at the same time, which can lead to a large service latency and high power consumption, one can “push” some information to the users beforehand at times when the link condition is good.

Without such predictive control, the cost minimization problem has been extensively studied and algorithms have been proposed, e.g., [10]. However, very little is known about the fundamental impact of prediction on system performance, let alone finding optimal control policies for such predictive queueing systems. Moreover, due to the existence of prediction windows and the fact that arrival processes are stochastic, the system naturally evolves according to a Markov chain whose state space size grows exponentially in the prediction window size. Thus, this problem is very challenging to solve.

3. PREDICTIVE BACKPRESSURE

In this section, we present our algorithm, which is designed by efficiently incorporating prediction information into the Backpressure technique [21]. Note that since future arrival information is made available in a sliding-window form, prediction couples the current action with future arrivals in every time slot. This prohibits the use of frame-based Lyapunov technique [22], and makes the problem complicated. Fortunately, as we will see, with the development of a novel *queue-equivalence* result, one can incorporate prediction into system control cleanly and significantly reduce the complexity in both algorithm design and analysis.

3.1 Prediction Queues

For our algorithm development and analysis, we now introduce the notion of a prediction queue, which records the number of residual arrivals in every slot in time window $[t, t + D_n - 1]$. Specifically, we denote $Q_n^{(d)}(t)$ the number of remaining arrivals currently in future slot $t + d$, i.e., d slots into the future, and denote $Q_n^{(-1)}(t)$ the number of packets already in the system. We see then the queues evolve according to the following dynamics:

1. If $d = D_n - 1$, then:

$$Q_n^{(d)}(t+1) = A_n(t+D_n). \quad (5)$$

2. If $0 \leq d \leq D_n - 2$, then:

$$Q_n^{(d)}(t+1) = [Q_n^{(d+1)}(t) - \mu_n^{(d+1)}(t)]^+. \quad (6)$$

3. For $Q_n^{(-1)}(t)$, we have:

$$Q_n^{(-1)}(t+1) = \left[Q_n^{(-1)}(t) - \mu_n^{(-1)}(t) \right]^+ + \left[Q_n^{(0)}(t) - \mu_n^{(0)}(t) \right]^+, \quad (7)$$

with $Q_n^{(-1)}(0) = 0$.

Fig. 3 shows the definition of the prediction queues. One can see that $\{Q_n^{(d)}(t)\}_{d=0}^{D_n-1}$ are not really queues. They simply record the residual arrivals going through the timeline, whereas $Q_n^{(-1)}(t)$ records the true backlog in the system. Notice that $Q_n^{(-1)}(t)$ is exactly the same as $Q_n(t)$ in (1). Since $Q_n^{(-1)}(t)$ is the only actual queue, the system is stable if and only if $Q_n^{(-1)}(t)$ is stable.

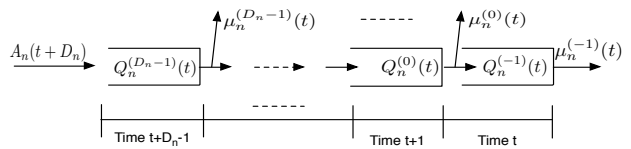


Figure 3: The prediction queues that describe the system evolution.

3.2 Predictive Backpressure

Here we construct our algorithm based on the above prediction queues and Backpressure. Our main idea is to use the sum of all the queues $Q_n^{\text{sum}} \triangleq \sum_{d=-1}^{D_n-1} Q_n^{(d)}(t)$ for decision making.

To describe the algorithm in details, we define the notion of queueing discipline for the predictive system, i.e., how to select packets to serve from queues $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$. Specifically, we order the packets in $Q_n^{(d)}(t)$ with labels $p_{\sum_{d' < d} Q_n^{(d')}(t)+1}, \dots, p_{\sum_{d' \leq d} Q_n^{(d')}(t)}$. Then, all the packets in $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$ are ordered from p_1 to $p_{Q_n^{\text{sum}}(t)}$. When a particular queueing discipline is applied in the predictive system, we select packets to serve according to the discipline using the order of the packets. For instance, if FIFO is used, then the server will serve the packets $p_1, \dots, p_{\min[\mu_n(t), Q_n^{\text{sum}}(t)]}$ from the queues every time. We now also define the notion of a *fully-efficient* predictive scheduling policy.

Definition 1. A predictive scheduling policy is called fully-efficient if for every user n , we have: (i) $\sum_d \mu_n^{(d)}(t) = \mu_n(t)$, and (ii) whenever there exists any $-1 \leq d \leq D_n - 1$ such that $\mu_n^{(d)}(t) > Q_n^{(d)}(t)$, $\mu_n^{(d')}(t) \geq Q_n^{(d')}(t)$, $\forall d' \neq d$. \diamond

In other words, if a policy is fully-efficient, it will always try to utilize all service opportunities and not allocate more service rate to serve any queue unless all other queues are already fully served.⁶ Hence, it will not waste any service opportunity unless there are more. With this definition, we now present our algorithm, in which $V \geq 1$ is a control parameter used to tradeoff utility performance and system delay (See Theorem 5).

Predictive Backpressure (PBP): In every time slot, compute $Q_n^{\text{sum}}(t) = \sum_{d=-1}^{D_n-1} Q_n^{(d)}(t)$ for all n . Then, observe the current channel state vector $\mathbf{S}(t)$ and perform:

- Choose the power allocation vector $\mathbf{P}(t)$ to solve the following problem:

$$\text{min : } \quad Vf(\mathbf{S}(t), \mathbf{P}(t)) - \sum_n Q_n^{\text{sum}}(t) \mu_n(\mathbf{S}(t), \mathbf{P}(t)) \quad (8)$$

$$\text{s.t. } \quad \mathbf{P}(t) \in \mathcal{P}(\mathbf{S}(t)). \quad (9)$$

Then, allocate the service rates $\{\mu_n^{(d)}(t)\}_{d=-1}^{D_n-1}$ to the queues $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$ in a fully-efficient manner according to any pre-specified queueing discipline.

- Update the queues according to (5), (6) and (7). \diamond

REMARK 1. Notice that the PBP algorithm has a very clean format. Indeed, PBP can be viewed as weighting the predicted

⁶It is equivalent to work-conserving in queue scheduling.

future arrivals for different users into current system control. It is worth emphasizing that such a low-complexity algorithm is not possible without the use of prediction queues and the Backpressure technique.

4. PERFORMANCE ANALYSIS

In this section, we first present an important theorem which states that if a predictive scheduling policy is fully-efficient, then the queueing system under the scheme evolves in the exact same way as a non-predictive queueing system with delayed arrivals and a different initial queue state. Using this queue-equivalence result, we obtain an interesting delay distribution shifting theorem. After that, we present our delay analysis for the PBP algorithm.

4.1 Performance of Fully-Efficient Scheduling Policies

We start by presenting the theorem regarding the equivalence between predictive and non-predictive systems.

Theorem 1. Let $\hat{Q}_n(t)$ be the queue size of a single queue system that (i) has $\hat{Q}_n(0) = \sum_{t=0}^{D_n-1} A_n(t)$, (ii) has arrival $\hat{A}_n(t) = A_n(t+D_n)$, (iii) has service $\hat{\mu}_n(t) = \sum_{d=-1}^{D_n-1} \mu_n^{(d)}(t)$, and (iv) evolves according to:

$$\hat{Q}_n(t+1) = \left[\hat{Q}_n(t) - \hat{\mu}_n(t) \right]^+ + \hat{A}_n(t). \quad (10)$$

Then, if the predictive system uses a fully-efficient predictive scheduling policy (with any queueing discipline) with $Q_n^{(-1)}(0) = 0$ for all n , we have for all t that:

$$Q_n^{sum}(t) = \hat{Q}_n(t), \quad \forall n. \quad \diamond \quad (11)$$

PROOF. See Appendix A. \square

Theorem 1 provides an important connection between the predictive system and the system without prediction. Using this result, we derive the following theorem, which relates the delay distribution of the predictive system to the equivalent system without prediction.

Theorem 2. (Delay Distribution Shifting) Denote $\pi_{n,k}^{(D_n)}$ the steady-state probability that a user n packet experiences a delay of k slots under a fully-efficient predictive scheduling policy in the predictive system, and let $\hat{\pi}_{n,k}$ denote the steady-state probability that a user n packet experiences a k -slot delay in $\hat{Q}_n(t)$. Suppose the set of queues $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$ and $\hat{Q}_n(t)$ use the same queueing discipline. Then, we have for each queue n that:

$$\pi_{n,0}^{(D_n)} = \sum_{k=0}^{D_n} \hat{\pi}_{n,k} \quad \text{and} \quad \pi_{n,k}^{(D_n)} = \hat{\pi}_{n,k+D_n}, \quad k \geq 1. \quad (12)$$

That is, the distribution of the original queue can be viewed as “shifted to the left by D_n slots” under predictive scheduling with D_n -slot prediction. \diamond

PROOF. See Appendix B. \square

Note that Theorem 2 is important for the general framework of predictive scheduling. It allows us to compare scheduling with prediction with the original queueing system, and enables us to leverage existing results in queueing theory for

analyzing predictive systems. To formalize this idea, first notice that if we start with $\hat{Q}_n(0) = 0$ and have $\hat{A}_n(t) = A_n(t)$, then $\hat{Q}_n(t)$ becomes exactly the same as the queueing process in the original system without prediction. Thus, if the steady-state behavior of $\hat{Q}_n(t)$ does not depend on the initial condition and the shift of the arrival process, e.g., a $G/G/1$ queue [24], then the delay performance of the predictive system can be understood by studying the delay distribution of the original system without prediction.

Corollary 1. Suppose $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$ and $\hat{Q}_n(t)$ use the same queueing discipline. For any arrival and service processes under which the delay distribution of $\hat{Q}_n(t)$ does not depend on $\hat{Q}_n(0)$ and the shift in the arrival process, we have:

$$\pi_{n,0}^{(D_n)} = \sum_{k=0}^{D_n} \pi_{n,k} \quad \text{and} \quad \pi_{n,k}^{(D_n)} = \pi_{n,k+D_n}, \quad k \geq 1. \quad (13)$$

Here $\pi_{n,k}$ is the steady-state probability that a user n packet experiences a delay of k slots in the system without prediction. \diamond

Note that Corollary 1 applies to general multi-queue single-server systems where the steady-state behavior depends only on the statistical behavior of the arrival and service processes. Note that (13) also quantifies how *tail delay* changes with predictive scheduling. Specifically, we now have in the predictive system that:

$$\Pr\{\text{Delay}_n > K\} = \sum_{k>K} \pi_{n,k+D_n}. \quad (14)$$

This is often exponentially smaller compared to that under the non-predictive system (see the simulation section).

With Theorem 2, we can now quantify how much delay improvement one can obtain via predictive scheduling. This is summarized in the following theorem, in which we use W_{tot} to denote the average delay of the original system without prediction, i.e.,

$$W_{tot} = \frac{\sum_{n=1}^N \lambda_n \sum_{k \geq 0} k \pi_{n,k}}{\sum_{n=1}^N \lambda_n}. \quad (15)$$

Theorem 3. Suppose the conditions in Corollary 1 hold. The delay reduction offered by predictive scheduling with prediction window vector \mathbf{D} , denoted by $R(\mathbf{D})$, is given by:

$$R(\mathbf{D}) = \frac{\sum_{n=1}^N \lambda_n \left(\sum_{1 \leq k \leq D_n} k \pi_{n,k} + D_n \sum_{k \geq 1} \pi_{n,k+D_n} \right)}{\sum_{n=1}^N \lambda_n}. \quad (16)$$

In particular, if $W_{tot} < \infty$, the average system delay goes to zero as D_n goes to infinity for all queue n , i.e.,

$$\lim_{D \rightarrow \infty} R(\mathbf{D}) = W_{tot}. \quad (17)$$

Here $\mathbf{D} \rightarrow \infty$ means $D_n \rightarrow \infty$ for all n . \diamond

PROOF. See Appendix C. \square

Note that Theorems 2 and 3, and Corollary 1 show that systems with predictive scheduling can be analyzed by studying the original system without prediction. Also note that the above results hold under any queueing discipline. The resulting delay distribution, of course, changes under different disciplines. Hence, the results also provide an efficient way for deciding how much prediction power is needed for

different systems under different control policies, e.g., if a system has a delay distribution under which most packets experience a delay of no more than D slots, then using a prediction power of D slots suffices to reap most of the benefit of predictive scheduling, and further investment on improving prediction power can readily be saved.

4.2 Performance of PBP

In this section, we analyze the performance of PBP. We will assume the following *system slack* condition: there exist a set of power vectors and probabilities $\{\mathbf{P}_m^{(s_i)}, a_m^{(s_i)}\}$, and a constant $\eta > 0$, such that:

$$\lambda_n - \sum_{s_i} \pi_{s_i} \sum_{m=0}^{\infty} a_m^{(s_i)} \mu_n(s_i, \mathbf{P}_m^{(s_i)}) \leq -\eta, \quad \forall n. \quad (18)$$

Note that (18) is commonly assumed in stochastic queueing system works and $\eta \geq 0$ is necessary for system stability [21]. The following theorem states that allowing predictive scheduling does not change the optimal average cost.

Theorem 4. *For any vector $\mathbf{0} \preceq \mathbf{D} \prec \infty$, we have:*

$$f_{av}^{D*} = f_{av}^*. \quad \diamond \quad (19)$$

PROOF. Omitted due to space limitation. Please see [25] for proof details. \square

Theorem 4 is interesting and shows that predictive scheduling does not reduce the minimum cost needed for system stability. Instead, the theorem, together with Theorem 5 below, deliver an important message that predictive scheduling *it improves the system delay given the same utility performance.*

We now have the following theorem, which shows that PBP achieves an average power consumption that is within $O(1/V)$ of the minimum and guarantees an average congestion bound.

Theorem 5. *The PBP algorithm achieves the following:*

$$f_{av}^{PBP} \leq f_{av}^{D*} + \frac{B}{V}, \quad Q_{av}^{sum} = Q_{av}^{BP} = \frac{B + V f_{av}^{max}}{\eta}. \quad (20)$$

Here $B = \frac{N}{2}(\mu_{max}^2 + A_{max}^2)$ is a constant independent of V , Q_{av}^{sum} denotes the average expected queue size of $\sum_n Q_n^{sum}(t)$, and Q_{av}^{BP} denotes the average expected queue size of the non-predictive system under Backpressure. \diamond

PROOF. Omitted due to space limitation. Please see [25] for proof details. \square

Theorem 5 is similar to the results in previous literature of Backpressure, e.g., [21]. It states that the average size of $\sum_n Q_n^{sum}(t)$ is the same as Q_{av}^{BP} under Backpressure without prediction. Since $Q_n^{sum}(t)$ is the total size of the actual queue and the prediction queues, we see that the actual queue size is strictly smaller than that under Backpressure. Since the average queue size under PBP is finite, we can apply Theorem 3 to obtain the following immediate corollary.

Corollary 2. *Suppose there exists a steady-state distribution of the queue vector under PBP. Then, the average delay under PBP goes to zero as $\mathbf{D} \rightarrow \infty$. \diamond*

Corollary 2 shows that with predictive scheduling, it is possible to achieve an $O(1/V)$ performance with an average

delay that is smaller than $O(\log(V))$. This is fundamentally different from the non-predictive case, in which the best utility-delay tradeoff is $[O(1/V), O(\log(V))]$ [9]. It is also tempting to analyze the exact delay reduction offered by PBP. However, due to the complex queueing dynamics under Backpressure, it is challenging to compute the exact distributions $\pi_{n,k}$ even without prediction. Thus, in the following, we consider a general class of cost-minimization problems, and study the delay reduction due to prediction in this case. For stating the results, we define the following optimization problem:

$$\max : \quad g(\boldsymbol{\gamma}), \quad \text{s.t. } \boldsymbol{\gamma} \succeq \mathbf{0}, \quad (21)$$

where $g(\boldsymbol{\gamma})$ is defined as:

$$g(\boldsymbol{\gamma}) = \sum_{s_i} \pi_{s_i} \inf_{\mathbf{P}_m^{(s_i)} \in \mathcal{P}^{(s_i)}} \left\{ V f(s_i, \mathbf{P}_m^{(s_i)}) + \sum_n \gamma_n [\lambda_n - \mu_n(s_i, \mathbf{P}_m^{(s_i)})] \right\}. \quad (22)$$

We now state our theorem regarding the average backlog reduction due to predictive scheduling. In the theorem, we use $\boldsymbol{\gamma}^*$ to denote an optimal solution of (21).

Theorem 6. *Suppose (i) $\boldsymbol{\gamma}^* = \Theta(V) > 0$ is unique, (ii) the η -slack condition (18) is satisfied with $\eta > 0$, (iii) the dual function $g(\boldsymbol{\gamma})$ satisfies:*

$$g(\boldsymbol{\gamma}^*) \geq g(\boldsymbol{\gamma}) + L \|\boldsymbol{\gamma}^* - \boldsymbol{\gamma}\|, \quad \forall \boldsymbol{\gamma} \succeq \mathbf{0}, \quad (23)$$

for some constant $L > 0$ independent of V , (iv) there exists a steady-state distribution of $\mathbf{Q}^{sum}(t)$ under PBP, (v) $D_n = O(\frac{1}{A_{max}}[\boldsymbol{\gamma}_n^* - G - K(\log(V))^2 - \mu_{max}]^+)$ for all n , and (vi) FIFO is used. Then, under PBP with a sufficiently large V , we have:

$$Q_{av}^{(-1)} \leq Q_{av}^{BP} - \sum_n D_n (\lambda_n - O(\frac{1}{V \log(V)}))^+. \quad \diamond \quad (24)$$

PROOF. See Appendix D. \square

As shown in [15], conditions (i)-(iii) in Theorem 6 are satisfied in many practical network optimization problems, especially when the power allocation sets $\{\mathcal{P}^{(s_i)}\}_{i=1}^M$ are finite. In this case, queue vector $\mathbf{Q}(t) = (Q_1(t), \dots, Q_N(t))$ mostly stays close to the fixed point $\boldsymbol{\gamma}^*$ [15]. Using Little's theorem, Theorem 6 implies that the system delay is reduced roughly linearly in the prediction window size D_n . Note that the linear reduction in D_n is due to the use of the FIFO discipline and the fact that D_n is constrained. In the case when D_n is larger than $O(\frac{1}{A_{max}}[\boldsymbol{\gamma}_n^* - G - K(\log(V))^2 - \mu_{max}]^+)$, or when LIFO is used, we observe that packet delay decreases *exponentially* in D_n (see Section 6).

5. PREDICTABLE-ONLY ARRIVAL

Here we discuss how PBP can also be applied (with slight modification) to the case when arrivals can only be predicted but not pre-served. The idea is to first pretend that the predictable-only traffic can also be pre-served, and then construct the algorithm and show that pre-serving rarely happens.

Predictable – Only PBP (POPBP): In every time slot, run PBP. In addition, for each queue n , do:

- (Marking) Mark all packets in $\{Q_n^{(d)}(t)\}_{d=0}^{D_n-1}$ served in the current time slot as *mistaken packets*.

- (Dropping) Drop all mistaken packets when they enter $Q_n^{(-1)}(t)$. \diamond

Here “run PBP” means carrying out all the steps in PBP including choosing and implementing actions, and updating queue values. Note that $\{Q_n^{(d)}(t)\}_{d=0}^{D_n-1}$ do not exactly correspond to the number of remaining future arrivals, as they are not served under POPBP. Instead, they are equal to the number of future arrivals excluding the mistaken packets.

The performance of POPBP is summarized as follows.

Theorem 7. *Suppose the conditions in Theorem 6 hold. Then, under POPBP with a sufficiently large V , we have:*

$$f_{av}^{\text{POPBP}} \leq f_{av}^{D*} + \frac{B}{V}, \quad (25)$$

$$Q_{av}^{(-1)} \leq Q_{av}^{\text{BP}} - \sum_n D_n (\lambda_n - O(\frac{1}{V \log(V)}))^+. \quad (26)$$

Moreover, the average packet dropping rate is $O(\frac{1}{V \log(V)})$. \diamond

PROOF. Omitted due to space limitation. Please see [25] for proof details. \square

6. SIMULATION

We present simulation results of the PBP algorithm in this section in a 10-user system.

6.1 Parameters and Settings

Real Data Trace: We collected data from 10 different mobile users in a 12-day long period (over five-minute intervals). The data for each user represents the aggregate amount of mobile traffic (over all applications) that is delivered by the base station to the user in that slot. Fig. 4 shows the traffic of a single user and the aggregate traffic of all users in 12 days. In the simulations, we take the traffic data as the workload arriving at the system which needs delivered by the server to different users.

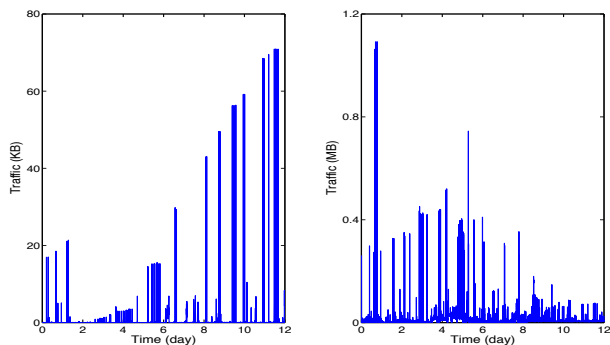


Figure 4: Left: The mobile traffic delivered by the base station to a single user. Right: The aggregate mobile traffic delivered by the base station to 10 different users.

Channel, Power and Prediction: For each user n , we assume that the channel condition $S_n(t)$ takes values $\{1, 2\}$ with equal probabilities, and $P_n(t) \in \mathcal{P}^{(S_n)} \triangleq \{0, 5, 10\}$. We assume that at any time, only one channel receives nonzero power allocation. The service rate is given by $\mu_n(t) = \lfloor 100 \log(1 + S_n(t) P_n(t)) \rfloor$ KB/s. The cost function $f(S(t), \mathbf{P}(t))$ is set to be $\sum_n P_n(t)$, which denotes the total power consumption. We set $D_n = \rho * 5$ for all n . Then, we simulate the cases $\rho \in \{1, 3, 5, 10\}$ to see the effect of the pre-

dictive scheduling. We simulate the algorithm with $V \in \{1, 5, 10, 20, 50, 100\}$.

6.2 Performance of PBP

Fig. 5 shows the performance of PBP with FIFO queueing policy. We see from the left plot that the average power consumption decreases as V increases. The right plot shows the average backlog under PBP. It is not hard to see that the average system backlog scales as $O(V)$. One also sees that as the prediction window sizes increase, the network delay decreases linearly in D .

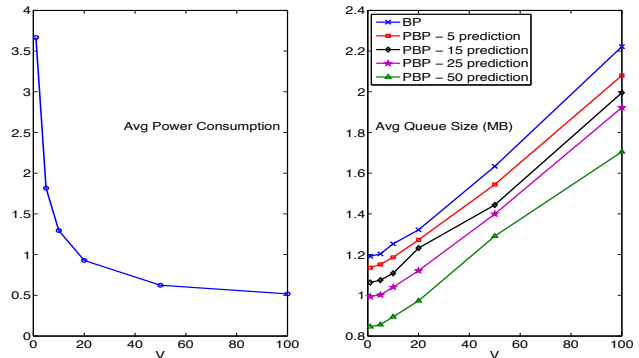


Figure 5: Performance of PBP. Left: Average power consumption under PBP. Right: Average queue size under PBP with different prediction window sizes.

Fig. 6 shows the delay distribution under PBP for the setting with $V = 100$ and $\rho = 1$. We see that the distributions of the latency for queue n are shifted to the left by D_n , as shown in Theorem 2. It can also be verified that Corollary 1 holds.

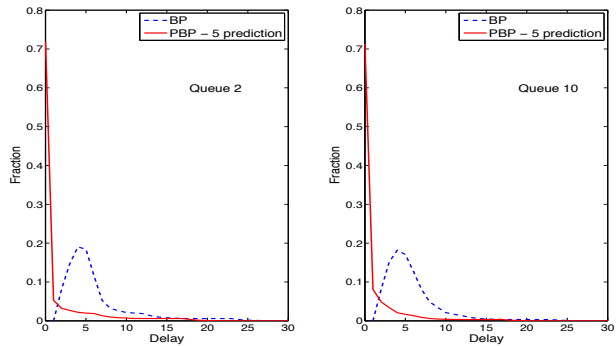


Figure 6: Packet delay distribution under PBP with FIFO scheduling with $V = 100$ and $D_n = 5$ for all n . We see that predictive scheduling effectively shifts the distribution to the left by 5 slots for both queues.

Fig. 7 then shows the delay distribution under PBP and the original Backpressure, under the LIFO discipline. It can be verified that the distribution for the predictive system is also a left-shifted version of the one under Backpressure. We see that large fractions of packets experience zero delay in both queues, i.e., they are served before they arrive at $Q_n^{(-1)}(t)$. This is so because under LIFO Backpressure (no prediction), most packets roughly experience $(\log(V))^2$ delay. Thus, with a moderate size prediction window size, the server can serve most packets before they enter the system. Since we use a log-scale for the x-axis, we do not plot the fraction for packets that have zero delay. Instead, we show the numbers

in the plot. We see that, 69% of the packets for queue 2 are served before they enter the system, whereas 73% of the packets are served for queue 10. These results demonstrate the power of predictive scheduling on delay reduction.

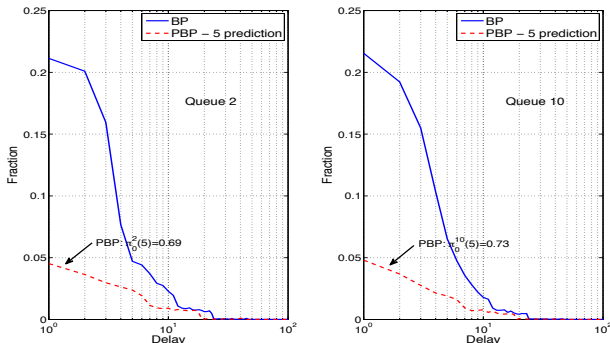


Figure 7: Packet Delay distribution under PBP with LIFO scheduling ($V = 100$ and $D_n = 5$ for all n). We see that a large fraction of the packets now experience zero delay! This is because with a moderate size prediction window, most packets are served before they arrive at $Q_n^{(-1)}(t)$. In the plots, $\pi_0^n(D_n)$ denotes the fraction of packets experience zero delay.

6.3 Impact of Imperfect Prediction

To investigate the impact of imperfect prediction, we consider two types of prediction errors. The first type is failing to predict actual arrivals, i.e., miss detection. When it happens, the arrivals will be out of the system's vision and thus will not appear in prediction queues. Therefore, they cannot be served predictively. The other type is false alarm, which happens when the system mistakenly predicts the existence of non-existing arrivals. Such false arrivals will appear in prediction queues, but will not enter the system. However, the system may incorrectly allocate resources to serve them, resulting in wasted service opportunities.

We model miss detections and false alarms as follows. Each time unit, q fraction of the request arrivals are miss detections, and $\frac{1-p}{p}$ fraction of the rest request arrivals are false alarms ($p \geq 0.5$). Therefore, fraction $\frac{1-q}{p}$ of all predicted requests are actual arrivals, which can be served beforehand. Larger q means more miss detections, and smaller p means more false alarms in the system. For perfect prediction, $q = 0$ and $p = 1$. We simulate three different settings: ($q = 0.5, p = 0.95$), ($q = 0.05, p = 0.75$) and ($q = 0.2, p = 0.85$). The first setting corresponds to the case when the prediction mechanism works very conservatively, which leads to very few false alarms but many miss detections. The second setting corresponds to the case where the prediction mechanism works very aggressively and results in few miss detections but many false alarms. The third setting is in-between.

As discussed above, $Q_n^{(d)}(t)$ ($0 \leq d \leq D_n - 1$) may now contain false alarms besides real arrivals, and the fraction of false alarms is $1 - p$ on average. Thus, the effective queue size of $Q_n^{(d)}(t)$ (the number of real arrivals) is $p \cdot Q_n^{(d)}(t)$. Therefore, in the PBP algorithm, we use $Q_n^{(-1)}(t) + p \sum_{d=0}^{D_n-1} Q_n^{(d)}(t)$ as the weight in (8) instead of $\sum_{d=-1}^{D_n-1} Q_n^{(d)}(t)$.

Fig. 8 shows the performance of PBP with FIFO queueing policy under imperfect prediction when $D_n = 5$ for all n . We see that PBP still improves the delay performance

while at the same time keeping a good power performance. This is because that the chance that PBP serves future arrivals is decreased when V increases. Therefore, the impact of prediction errors on the utility performance of PBP is reduced, showing that PBP is robust against prediction error. Compared to PBP with perfect prediction, prediction errors increase the average backlog of the system and thus the average delay. This is intuitive since miss detections cannot be served beforehand and false alarms waste service opportunities. We also observe that the delay performance of PBP is more sensitive to miss detections. This is because miss detections enter the system directly and may increase the system backlog.

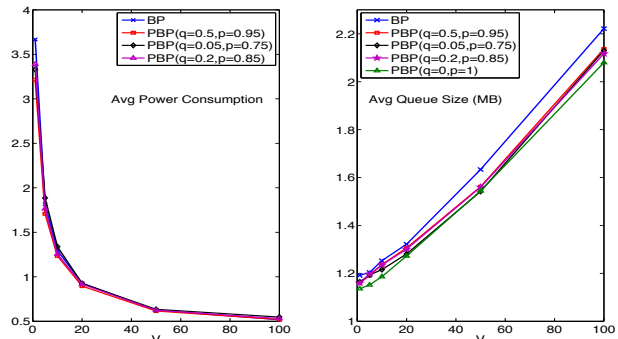


Figure 8: Performance of PBP under imperfect prediction when $D_n = 5$ for all n . Left: Average power consumption under different prediction errors. Right: Average queue size under different prediction errors.

7. CONCLUSION

In this paper, we investigate the fundamental benefit of predictive scheduling in controlled queueing systems. Based on a lookahead prediction window model, we establish a novel queue-equivalence result, which enables exact analysis of queueing systems under predictive scheduling using traditional queueing network control techniques. We then propose the **Predictive Backpressure** (PBP) algorithm, and show that PBP achieves a cost performance that is arbitrarily close to the optimal, while guaranteeing that the average system delay vanishes as the prediction window size increases. Our results can also be used to determine the required prediction power and analyze the tail delay reduction improvement.

8. ACKNOWLEDGEMENT

This work was supported in part by the National Basic Research Program of China Grant 2011CBA00300, 2011CBA00301, the National Natural Science Foundation of China Grant 61033001, 61361136003, 61303195, and the China youth 1000-talent grant. It was also partially supported by the University Grants Committee of the Hong Kong Special Administrative Region, China (Area of Excellence Grant Project No. AoE/E-02/08 and General Research Fund Project No. 411011), and Chancellor's Fellowship and NSF grant CNS1147930.

Appendix A – Proof of Theorem 1

Here we prove Theorem 1.

PROOF. (Theorem 1) We prove the result by induction with the aid of the following figure showing the evolution of $\hat{Q}_n(t)$.

First, we see that the result holds for $t = 0$: On one hand, $\hat{Q}_n(0) = \sum_{t=0}^{D_n-1} A_n(t)$. On the other hand, in the

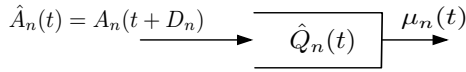


Figure 9: The original queue without prediction and with a delayed arrival process as well as a different initial queue state.

system under predictive scheduling, since $Q_n^{(-1)}(0) = 0$ and $Q_n^{(d)}(0) = A_n(d)$ for $d \in \{0, \dots, D_n - 1\}$, we have $Q_n^{\text{sum}}(0) = \sum_{t=0}^{D_n-1} A_n(t)$.

Now suppose the result holds for all $t = 0, \dots, k$, we show that it holds for $t = k+1$. Using the queueing dynamic equation (10), we know that in time slot k , $\tilde{\mu}_n(k) = \min[\mu_n(k), \hat{Q}_n(k)]$ packets will be served from $\hat{Q}_n(k)$. Now consider the queues $\{Q_n^{(d)}(k)\}_{d=-1}^{D_n-1}$. Since the scheduling policy is fully-efficient, we must have that the number of packets served from these queues is also $\tilde{\mu}_n(k) = \min[\mu_n(k), \hat{Q}_n(k)]$. To see this, note that if $\tilde{\mu}_n(k) = \mu_n(k)$, there are more packets in the queues than the number of packets that can be served. In this case, we must have $\mu_n^{(d)}(k) \leq Q_n^{(d)}(k)$ for all d . Also, because the policy is fully-efficient, we have $\sum_d \mu_n^{(d)}(k) = \mu_n(k)$. Hence, exactly $\mu_n(k)$ packets will be served from $\{Q_n^{(d)}(k)\}_{d=-1}^{D_n-1}$, resulting in $\hat{Q}_n(k+1) = Q_n^{\text{sum}}(k+1)$. On the other hand, suppose $\tilde{\mu}_n(k) = \hat{Q}_n(k)$. Then, there are enough service opportunities to clear all the awaiting packets. In this case, since the scheduling policy is fully-efficient, exactly $\hat{Q}_n(k)$ packets will be served. Thus, in both cases, we have $\hat{Q}_n(k+1) = Q_n^{\text{sum}}(k+1) = A_n(k+D_n)$. \square

Appendix B – Proof of Theorem 2

Here we prove Theorem 2.

PROOF. (Theorem 2) From Theorem 1, we see that $Q_n^{\text{sum}}(t) = \hat{Q}_n(t)$ for all time. Hence, if the two queueing systems use the same queueing discipline in choosing what packets to serve, then every packet will experience the exact same delay in both $\hat{Q}_n(t)$ and $\{Q_n^{(d)}(t)\}_{d=-1}^{D_n-1}$.

However, in $Q_n^{\text{sum}}(t)$, a packet will enter the actual system only after spending one unit of time in each of the queues in $\{Q_n^{(d)}(k)\}_{d=0}^{D_n-1}$, which takes exactly D_n slots in total. Thus, any packet experiencing a k -slot delay will experience $[k - D_n]^+$ delay in $Q_n^{(-1)}(t)$. This completes the proof of Theorem 2. \square

Appendix C – Proof of Theorem 3

We prove Theorem 3 here.

PROOF. (Theorem 3) Using Corollary 1, we see that in the predictive system, the average system backlog size is given by:

$$N_{\text{tot}}^{\text{P}} = \sum_{n=1}^N \lambda_n \sum_{k \geq 1} k \pi_{n,k+D_n}. \quad (27)$$

On the other hand, the average system backlog without prediction is given by:

$$N_{\text{tot}} = \sum_{n=1}^N \lambda_n \sum_{k \geq 1} k \pi_{n,k}. \quad (28)$$

Using (28) and (27), we conclude that:

$$N_{\text{tot}} - N_{\text{tot}}^{\text{P}} = \sum_{n=1}^N \lambda_n \left(\sum_{1 \leq k \leq D_n} k \pi_{n,k} + D_n \sum_{k \geq 1} \pi_{n,k+D_n} \right). \quad (29)$$

Using Little's theorem and dividing both sides by $\sum_n \lambda_n$, we see that (16) follows.

Now we prove (17). By taking a limit as $D \rightarrow \infty$, we first obtain:

$$\lim_{D \rightarrow \infty} \sum_n \lambda_n \sum_{1 \leq k \leq D_n} k \pi_{n,k} = N_{\text{tot}}. \quad (30)$$

Then, using the fact that $W_{\text{tot}} < \infty$, we have:

$$\lim_{D_n \rightarrow \infty} D_n \sum_{k \geq 1} \pi_{n,k+D_n} = 0, \quad \forall n. \quad (31)$$

Using the above in (16), we see that (17) follows. \square

Appendix D – Proof of Theorem 6

We prove Theorem 6. For our proof, we use the following theorem (which is Theorem 1 in [15]), in which γ^* denotes an optimal solution of (21). According to [15], γ^* is either $\Theta(V)$ or 0.

Theorem 8. Suppose (i) γ^* is unique, (ii) the η -slack condition (18) is satisfied with $\eta > 0$, (iii) the function $g(\gamma)$ satisfies:

$$g(\gamma^*) \geq g(\gamma) + L \|\gamma^* - \gamma\|, \quad \forall \gamma \succeq \mathbf{0}, \quad (32)$$

for some constant $L > 0$ independent of V . Then, under Backpressure, there exist constants $G, K, c = \Theta(1)$, i.e., all independent of V , such that for any $m \in \mathbb{R}_+$,

$$\mathcal{P}^{(r)}(G, Km) \leq ce^{-m}, \quad (33)$$

where $\mathcal{P}^{(r)}(G, Km)$ is defined:

$$\begin{aligned} \mathcal{P}^{(r)}(G, Km) & \quad (34) \\ & \triangleq \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=0}^{t-1} \Pr\{\exists n, |Q_n(\tau) - \gamma_n^*| > G + Km\}. \quad \diamond \end{aligned}$$

PROOF. See [15]. \square

We are now ready to present the proof of Theorem 6.

PROOF. (Theorem 6) We prove the results using Little's theorem. The main idea is to show that the average system queue length is roughly reduced by $\sum_n \lambda_n D_n$. To prove this, we show that the average total service rate allocated to the prediction queues is $O(\frac{1}{\sqrt{\log(V)}})$. Then, the average rate of the packets that go through $\{Q_n^{(d)}(t)\}_{d=0}^{D_n-1}$ will roughly be λ_n , and so the average queue size is reduced by roughly $\sum_n \lambda_n D_n$.

First, using (11) and (33), we see that in steady state,

$$\Pr\{|Q_n^{\text{sum}}(t) - \gamma_n^*| > G + Km\} \leq ce^{-m}.$$

Using the fact that $Q_n^{\text{sum}}(t) = \sum_{d=-1}^{D_n-1} Q_n^{(d)}(t)$, we have:

$$\Pr\{Q_n^{(-1)}(t) < \gamma_n^* - G - Km - \sum_{d=0}^{D_n-1} Q_n^{(d)}(t)\} \leq ce^{-m}.$$

Now let $m = (\log(V))^2$. Since $\gamma_n^* = \Theta(V)$, we see that when V is sufficiently large, we have:

$$\gamma_n^* - G - Km - \sum_d Q_n^{(d)}(t)$$

$$\begin{aligned}
&= \Theta(V) - G - K(\log(V))^2 - \sum_d Q_n^{(d)}(t) \\
&\stackrel{(a)}{\geq} \Theta(V) - G - K(\log(V))^2 - D_n A_{\max} \\
&\stackrel{(b)}{\geq} \mu_{\max}. \tag{35}
\end{aligned}$$

Here (a) follows from the fact that $Q_n^{(d)}(t) \leq A_{\max}$ for all $0 \leq d \leq D_n - 1$, and in (b) we use the fact that V is sufficiently large and $D_n = O(\frac{1}{A_{\max}}[\gamma_n^* - G - K(\log(V))^2 - \mu_{\max}]^+)$ for all n . This shows that the probability for $Q_n^{(-1)}(t)$ to go below μ_{\max} is at most $ce^{-(\log(V))^2} = \frac{c}{V^{\log(V)}}$.

Using the fact that under the FIFO queueing discipline, a prediction queue $Q_n^{(d)}(t)$ will be served only when $Q_n^{(-1)}(t) < \mu_{\max}$, we conclude that the average service rate allocated to the prediction queues is no more than $\frac{c\mu_{\max}}{V^{\log(V)}}$. Hence, the average traffic rate of the packets that traverse all prediction queues and eventually enter $Q_n^{(-1)}(t)$ is at least $[\lambda_n - \frac{c\mu_{\max}}{V^{\log(V)}}]^+$. Since every packet stays 1 slot in every prediction queue, using Little's theorem, we conclude that the average size of the prediction queues, denoted by $\sum_{d=0}^{D_n-1} \bar{Q}_n^{(d)}$ satisfies $\sum_{d=0}^{D_n-1} \bar{Q}_n^{(d)} \geq (\lambda_n - \frac{c\mu_{\max}}{V^{\log(V)}})^+ D_n$. Hence, (24) follows. \square

9. REFERENCES

- [1] M. Maia, J. Almeida, and V. Almeida. Identifying user behavior in online social networks. *Proceedings of the 1st Workshop on Social Network Systems*, pages 1-6, 2008.
- [2] I. Weber and A. Jaimes. Who uses web search for what? and how? *Web Search and Data Mining (WSDM)*, pages 21-30, 2011.
- [3] R. Kumar and A. Tomkins. A characterization of online browsing behavior. *Proceedings of the 19th international conference on World Wide Web*, pages 561-570, 2010.
- [4] V. N. Padmanabhan and J. C. Mogul. Using predictive prefetching to improve world wide web latency. *ACM SIGCOMM Computer Communication Review*, Volume 26, Issue 3, Pages 22-36, July 1996.
- [5] J. Lee, H. Kim, and R. Vuduc. When prefetching works, when it doesn't, and why. *ACM Transactions on Architecture and Code Optimization (TACO)*, Volume 9, Issue 1, March 2012.
- [6] T. Ball and J. R. Larus. Branch prediction for free. *Proceedings of the Conference on Programming Language Design and Implementation*, *ACM SIGPLAN Notices*, volume 28, pages 300-13, 1993.
- [7] M. U. Farooq, Khubaib, and L. K. John. Store-load-branch (slb) predictor: A compiler assisted branch prediction for data dependent branches. *Proceedings of the 19th IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, February 2013.
- [8] R. Berry and R. Gallager. Communication over fading channels with delay constraints. *IEEE Transactions on Information Theory*, vol. 48, no. 5, pp. 1135-1149, May 2002.
- [9] M. J. Neely. Optimal energy and delay tradeoffs for multi-user wireless downlinks. *IEEE Transactions on Information Theory* vol. 53, no. 9, pp. 3095-3113, Sept. 2007.
- [10] M. J. Neely. Energy optimal control for time-varying wireless networks. *IEEE Transactions on Information Theory* 52(7): 2915-2934, July 2006.
- [11] A. C. Fu, E. Modiano, and J. N. Tsitsiklis. Optimal energy allocation and admission control for communications satellites. *IEEE/ACM Transactions on Networking*, vol. 11, no. 3, pp. 488-500, 2003.
- [12] M. A. Zafer and E Modiano. A calculus approach to energy-efficient data transmission with quality-of-service constraints. *IEEE/ACM Transactions on Networking*, 17(3), 898-911, 2009.
- [13] C. W. Tan, D. P. Palomar, and M. Chiang. Energy-robustness tradeoff in cellular network power control. *IEEE/ACM Transactions on Networking*, Vol. 17, No. 3, pp. 912-925, 2009.
- [14] M. J. Neely. Super-fast delay tradeoffs for utility optimal fair scheduling in wireless networks. *IEEE Journal on Selected Areas in Communications (JSAC)*, Special Issue on Nonlinear Optimization of Communication Systems, vol. 24, no. 8, pp. 1489-1501, Aug. 2006.
- [15] L. Huang and M. J. Neely. Delay reduction via Lagrange multipliers in stochastic network optimization. *IEEE Trans. on Automatic Control*, Volume 56, Issue 4, pp. 842-857, April 2011.
- [16] J. Tadrous, A. Eryilmaz, and H. El Gamal. Proactive resource allocation: harnessing the diversity and multicast gains. *IEEE Transactions on Information Theory*, 2013.
- [17] J. Tadrous, A. Eryilmaz, and H. El Gamal. Pricing for demand shaping and proactive download in smart data networks. *The 2nd IEEE International Workshop on Smart Data Pricing (SDP)*, *INFOCOM*, 2013.
- [18] J. Spencer, M. Sudan, and K Xu. Queueing with future information. *ArXiv Technical Report arxiv:1211.0618*, 2012.
- [19] Y. Yao, L. Huang, A. Sharma, L. Golubchik, and M. J. Neely. Data centers power reduction: A two time scale approach for delay tolerant workloads. *USC CS Technical Report 11-9020*, 2011.
- [20] I. Hou and P.R. Kumar. Broadcasting delay-constrained traffic over unreliable wireless links with network coding. *Proceedings of MobiHoc*, 2011.
- [21] L. Georgiadis, M. J. Neely, and L. Tassiulas. *Resource Allocation and Cross-Layer Control in Wireless Networks*. Foundations and Trends in Networking Vol. 1, no. 1, pp. 1-144, 2006.
- [22] L. Huang and M. J. Neely. Max-weight achieves the exact $[O(1/V), O(V)]$ utility-delay tradeoff under Markov dynamics. *arXiv:1008.0200v1*, 2010.
- [23] A. Wierman M. Lin, Z. Liu and L. Andrew. Online algorithms for geographical load balancing. *International Green Computing Conference (IGCC)*, 2012.
- [24] R. G. Gallager. *Discrete Stochastic Processes*. Kluwer Academic Publishers, 1996.
- [25] L. Huang, S. Zhang, M. Chen, and X. Liu. When Backpressure meets Predictive Scheduling. *arXiv Technical Report*, arXiv:1309.1110v1.