

# Appendix to *Creating Work Breaks From Available Idleness*

Xu Sun and Ward Whitt

Department of Industrial Engineering and Operations Research, Columbia University  
New York, NY, 10027

April 8, 2017

## 1 Overview

This is an appendix to the main paper, *Creating Work Breaks From Available Idleness*. In §2 we provide additional simulation results for the work-conserving  $D_1$  server-assignment rule. In §3, we present additional simulation results for the  $D_2$  rule, in particular, for a large service system with  $n = 1000$ . In §4, we examine another way to announce work breaks from the available idleness. We refer to this rule as  $LISF - D_2$ .

## 2 Additional Results for the $D_1$ Assignment Rule

In this section we provide additional simulation results for the  $D_1$  assignment rule. In §2.1 we examine the impact of  $\theta$  (target length-of-break) on system performance. We present more results for the  $M/M/n$  queues in §2.2 and more results for the  $M/H_2/n$  queues in §2.3.

### 2.1 Impact of $\theta$ on the $D_1$ Rule

From Lemma 2.1 in the paper, we know that the choice of the parameter  $\theta$  directly influences the maximum possible rate at which breaks occur. Here we show via histograms that breaks occur less often as  $\theta$  increases. Besides, as  $\theta$  increases, the proportion of idleness on breaks also decreases.

The simulation results for the idle time distribution in the  $M/M/n$  model with rule  $D_1$  and model parameters  $\mu = 1$ ,  $\rho = 0.9$ ,  $n = 100$  and three different values of  $\theta$ :  $\theta = k/3$  for  $k = 4, 5, 6$  are displayed in Figure 1 and Table 1.

Panel (b) is for our base model with mean service times of 3 minutes, where the target duration was a 5-minute work break every one or two hours, so that  $\theta = 5/3$ . Evidently, the  $D_1$  rule is able to create work breaks from idleness. Figure 1 shows that the  $D_1$  rule creates a peak in the distribution at the target  $\theta$  and the rest of the distribution concentrates near the origin, decaying very rapidly. Overall, we obtain a probability density function which is bimodal.

We elaborate in Table 1 by showing the 95% confidence intervals (based on the data collected from the simulation experiments) for the mean and standard deviation of the idle time  $V_n$  and the long-run proportion of idle times that are work breaks and  $\pi_{\beta, I}$  to represent the long-run proportion of idle time that is made up of work break time.

Consistent with the conservation laws developed in the main paper,  $E[V_n]$  approximately equals  $1/9 = 0.1111$  for all  $\theta$ , because the long-run proportion of idleness is solely determined by the traffic intensity  $\rho$ , independent of

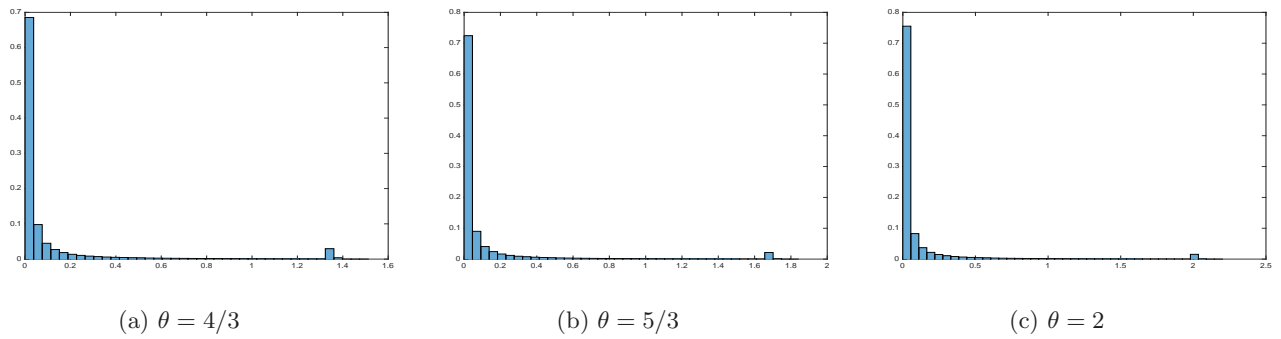


Figure 1: Histograms estimated from simulation of the idle-time distribution with rule  $D_1(\theta)$  for three candidate targets  $\theta$

$\theta$	$E[V_n]$	$std(V_n)$	$\pi_{\beta, I}$
6/6	$0.1112 \pm 5 \times 10^{-4}$	$0.2488 \pm 6 \times 10^{-4}$	$0.486 \pm 0.001$
7/6	$0.1112 \pm 5 \times 10^{-4}$	$0.2653 \pm 7 \times 10^{-4}$	$0.431 \pm 0.001$
8/6	$0.1113 \pm 4 \times 10^{-4}$	$0.2800 \pm 7 \times 10^{-4}$	$0.398 \pm 0.001$
9/6	$0.1113 \pm 4 \times 10^{-4}$	$0.2928 \pm 7 \times 10^{-4}$	$0.369 \pm 0.001$
10/6	$0.1108 \pm 5 \times 10^{-4}$	$0.3035 \pm 9 \times 10^{-4}$	$0.340 \pm 0.001$
11/6	$0.1108 \pm 5 \times 10^{-4}$	$0.3138 \pm 1 \times 10^{-3}$	$0.315 \pm 0.002$
12/6	$0.1112 \pm 5 \times 10^{-4}$	$0.3236 \pm 1 \times 10^{-3}$	$0.294 \pm 0.002$

Table 1: Estimated performance measures of  $D_1(\theta)$  as a function of  $\theta$

the server-assignment rule (provided that it is work-conserving). In addition,  $std(V_n)$  increases as  $\theta$  grows, but not significantly. Table 1 and Figure 1 also show that as  $\theta$  increases from 1 to 2, the proportion of idle time occupied by breaks decreases slowly, changing from 0.486 to 0.294.

## 2.2 Impact of System Size on the $D_1$ Assignment Rule

The simulation results for the idle time distribution in the  $M/M/n$  model with rule  $D_1$  and model parameters  $\mu = 1$ ,  $\rho = 0.9$ ,  $\theta = 5/3$  and 6 values of  $n$  ranging from 100 to 5000 are summarized in Table 2. See also the histograms and ECDFs in Figure 2 - 3 for rule  $D_1$  as functions of  $n$ . Consistent with the fluid limit derived in the main paper, these histograms have a tendency to converge to the suggested form, an extremal two-point distribution with mass  $p$  on  $\theta$  and mass  $1 - p$  on 0. Here we recall that  $p \equiv m/\theta = 0.0667$  with  $m \equiv (1 - \rho)/\rho = 0.1111$ . Moreover, from Table 2 we see that the probability values  $P(V_n \geq \theta)$  have a tendency to converge to the limit  $p = 0.0667$  as desired.

Simulation outputs for the period between successive work breaks are reported in Table 3 and the histograms and ECDFs in Figure 4 - 5 as functions of  $n$ . In line with our asymptotic analysis, these histograms converge slowly to the desired form, i.e., a shifted exponential distribution. Specifically, the limit  $T$  can be expressed as  $T \stackrel{d}{=} x^* + \theta + M$  where  $M$  denotes an exponential r.v. with unit rate. Using the formulas derived in the main paper, we get  $x^* = 1/p - 1 = 15 - 1 = 14$  and hence

$$\mathbb{E}[T] = x^* + 1 + \theta = 16.6667 \quad \text{and} \quad \sqrt{Var(T)} = \sqrt{Var(M)} = 1.$$

The first part of Table 3 shows strong evidence that both  $\mathbb{E}[T_n]$  and  $Var(T_n)$  converge to the correct limit as  $s \rightarrow \infty$ .

	$P(V_n \leq 0.001)$	$P(V_n \leq 0.01)$	$P(V_n \leq 0.1)$	$P(V_n \leq 1)$	$P(V_n \geq \theta)$
$n = 100$	$0.2539 \pm 0.0035$	$0.4629 \pm 0.0027$	$0.8240 \pm 0.0013$	$0.9657 \pm 4 \times 10^{-4}$	$0.0223 \pm 4 \times 10^{-4}$
$n = 250$	$0.1545 \pm 0.0027$	$0.5270 \pm 0.0017$	$0.8498 \pm 7 \times 10^{-4}$	$0.9589 \pm 3 \times 10^{-4}$	$0.0317 \pm 3 \times 10^{-4}$
$n = 500$	$0.1821 \pm 0.0012$	$0.6228 \pm 0.0010$	$0.8717 \pm 8 \times 10^{-4}$	$0.9531 \pm 5 \times 10^{-4}$	$0.0405 \pm 5 \times 10^{-4}$
$n = 1000$	$0.2813 \pm 6 \times 10^{-4}$	$0.7093 \pm 7 \times 10^{-4}$	$0.8896 \pm 8 \times 10^{-4}$	$0.9474 \pm 7 \times 10^{-4}$	$0.0492 \pm 7 \times 10^{-4}$
$n = 2500$	$0.4618 \pm 5 \times 10^{-4}$	$0.7921 \pm 4 \times 10^{-4}$	$0.9074 \pm 5 \times 10^{-4}$	$0.9424 \pm 5 \times 10^{-4}$	$0.0564 \pm 5 \times 10^{-4}$
$n = 5000$	$0.5893 \pm 3 \times 10^{-4}$	$0.8333 \pm 3 \times 10^{-4}$	$0.9155 \pm 2 \times 10^{-4}$	$0.9395 \pm 2 \times 10^{-4}$	$0.0601 \pm 2 \times 10^{-4}$
$n = \infty$	0.9333	0.9333	0.9333	0.9333	0.0667

Table 2: Statistics for the idle-time distribution with rule  $D_1$

system	$D_1$		$SISF$	
	$E[T_n]$	$std(T_n)$	$E[T_n]$	$std(T_n)$
$n = 100$	$48.06 \pm 0.79$	$18.73 \pm 0.41$	$37.85 \pm 0.49$	$36.68 \pm 0.52$
$n = 250$	$33.45 \pm 0.33$	$9.47 \pm 0.35$	$28.62 \pm 0.21$	$27.01 \pm 0.28$
$n = 500$	$25.79 \pm 0.34$	$5.66 \pm 0.21$	$23.38 \pm 0.20$	$21.65 \pm 0.17$
$n = 1000$	$20.84 \pm 0.30$	$3.06 \pm 0.12$	$20.28 \pm 0.16$	$18.54 \pm 0.16$
$n = 2500$	$17.99 \pm 0.14$	$1.75 \pm 0.06$	$18.18 \pm 0.09$	$16.46 \pm 0.07$
$n = 5000$	$16.75 \pm 0.07$	$1.38 \pm 0.03$	$17.28 \pm 0.05$	$15.59 \pm 0.06$
$n = \infty$	16.67	1.00	16.67	15.00

Table 3: Statistics for  $T_n$  with rule  $D_1$

From the MHHT  $M/M$  fluid model with rule  $D_1$ , we expect that the age  $A_B$  of a busy server to follow a mixture distribution consisting of  $U[0, \tau^*]$  and  $\tau^* + M$  glued together on each side of the threshold  $x^* = 14$  and the age  $A_I$  of an idle server to follow a mixture distribution consisting of a truncated-exponential and an exponential distribution spliced together back-to-back, as  $n \rightarrow \infty$ . The summary statistics for  $A_B$  and  $A_I$  are reported in Table 4. Figure 6 - 7 together with Table 4 strongly support the heavy-traffic fluid limits for  $A_B$  and  $A_I$ . See also the histograms as shown in Figure 6 - 7.

	Busy		Idle	
	$E[A_B]$	$std(A_B)$	$E[A_I]$	$std(A_I)$
$n = 100$	$26.510 \pm 0.055$	$19.146 \pm 0.076$	$41.725 \pm 0.073$	$19.725 \pm 0.088$
$n = 250$	$17.602 \pm 0.023$	$11.788 \pm 0.045$	$31.131 \pm 0.038$	$10.796 \pm 0.041$
$n = 500$	$13.178 \pm 0.016$	$8.395 \pm 0.037$	$24.858 \pm 0.022$	$6.565 \pm 0.028$
$n = 1000$	$10.518 \pm 0.010$	$6.380 \pm 0.021$	$20.865 \pm 0.014$	$3.828 \pm 0.019$
$n = 2500$	$9.012 \pm 0.006$	$5.349 \pm 0.018$	$18.408 \pm 0.009$	$2.406 \pm 0.012$
$n = 5000$	$8.399 \pm 0.004$	$4.935 \pm 0.011$	$17.378 \pm 0.005$	$1.797 \pm 0.009$

Table 4: Statistics for  $A_B$  and  $A_I$  ( $M/M/n$  model)

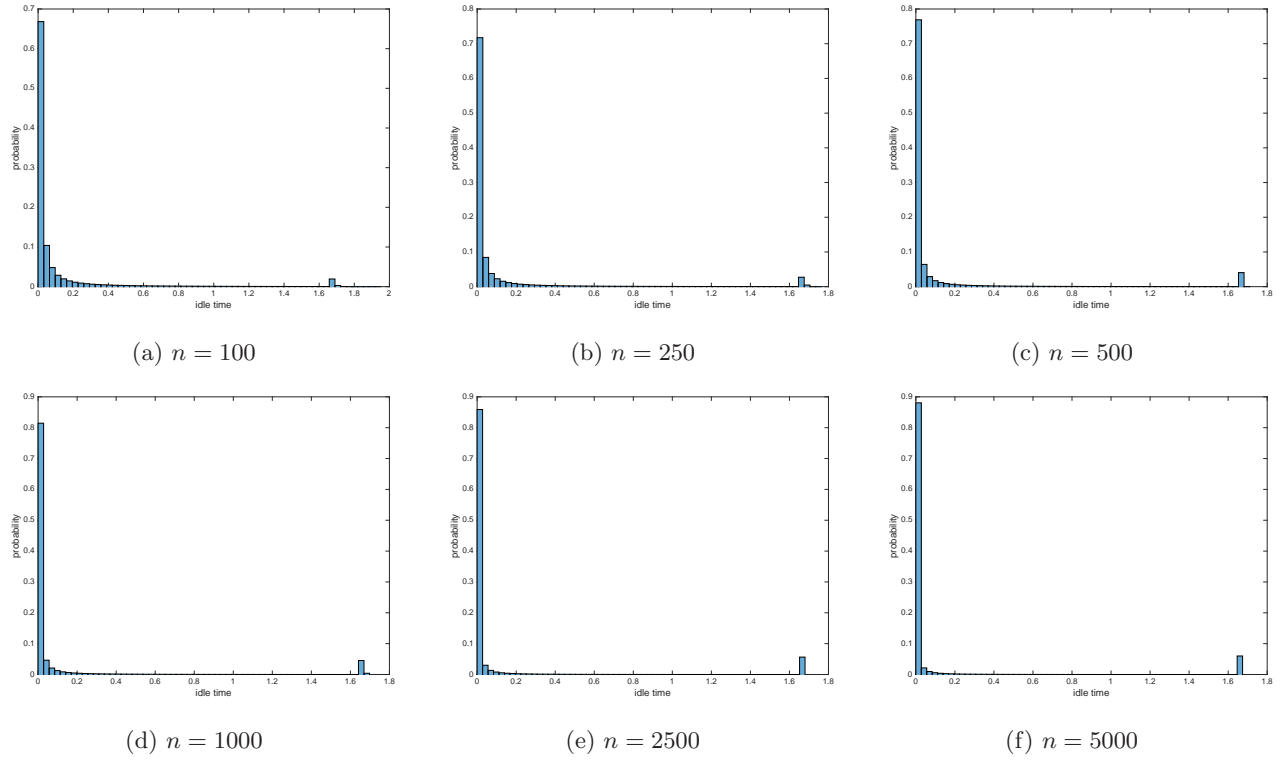


Figure 2: Histogram of idle periods estimated from computer simulation with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$

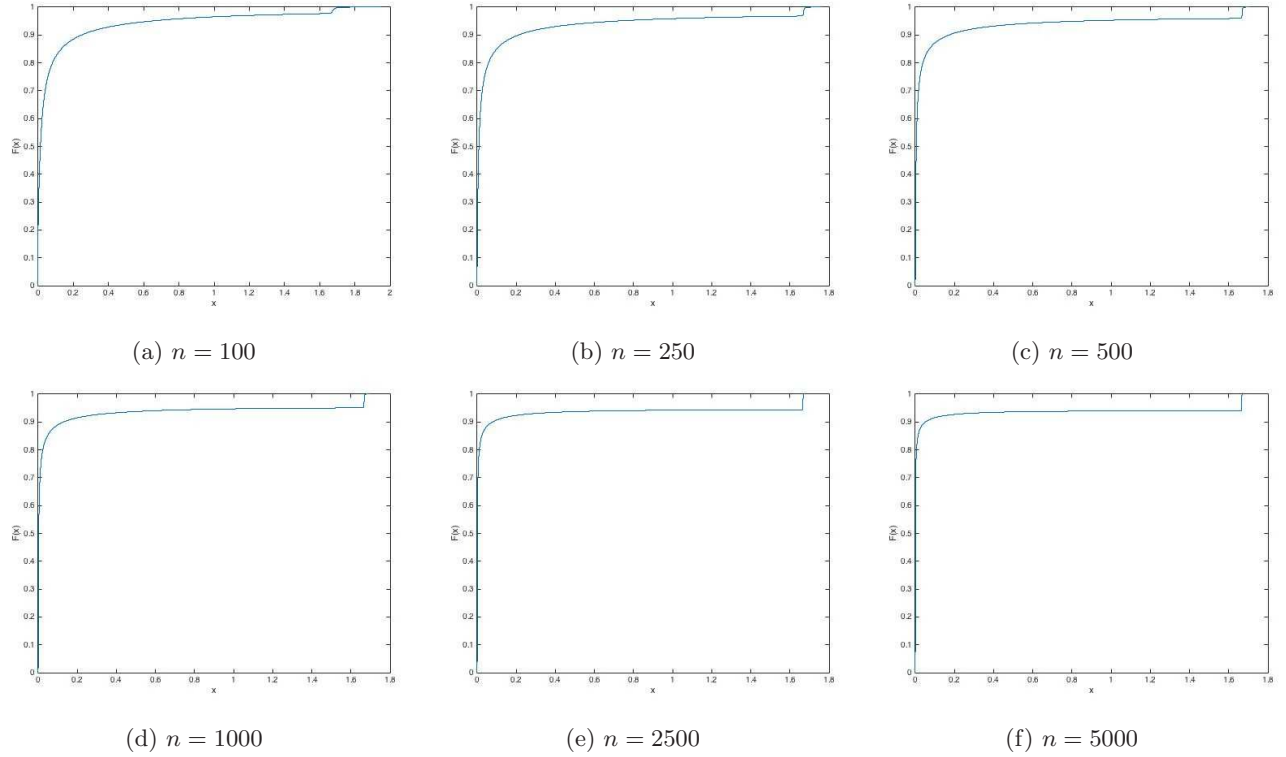
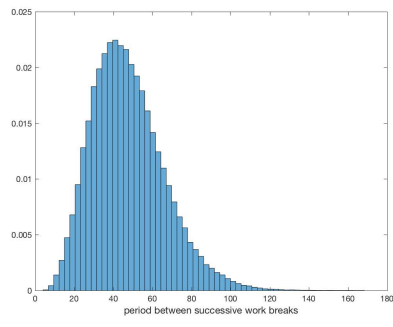
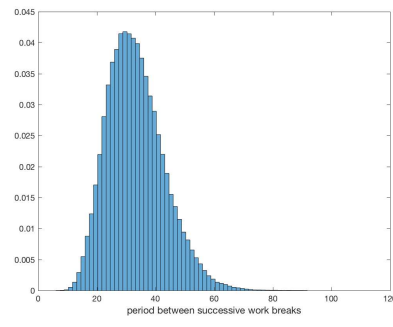


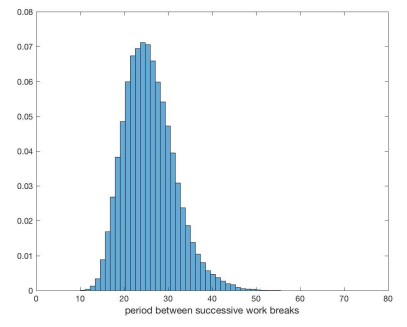
Figure 3: ECDF of idle periods estimated from computer simulation with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$



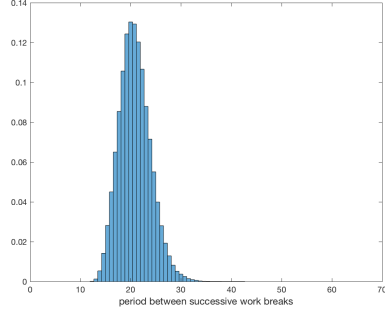
(a)  $n = 100$



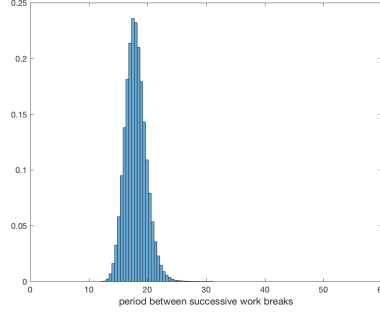
(b)  $n = 250$



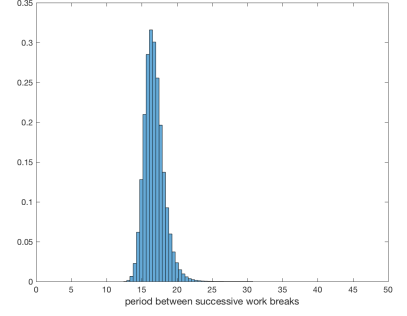
(c)  $n = 500$



(d)  $n = 1000$

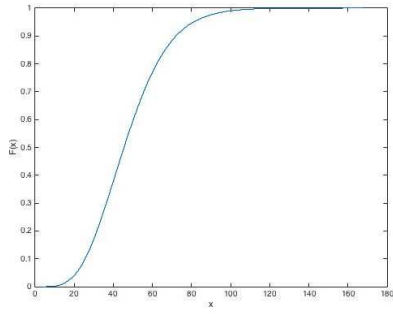


(e)  $n = 2500$

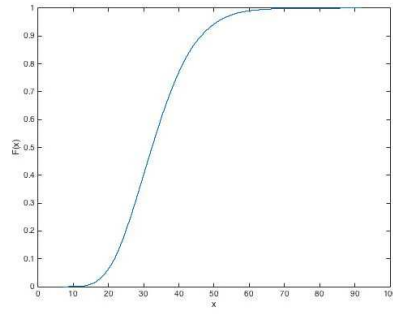


(f)  $n = 5000$

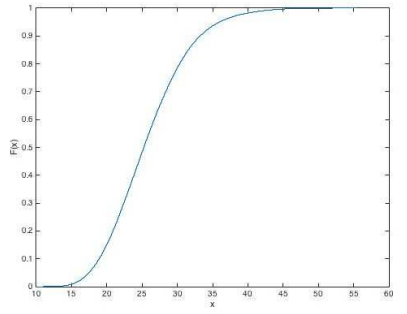
Figure 4: Histogram of periods between successive breaks estimated from computer simulation with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$



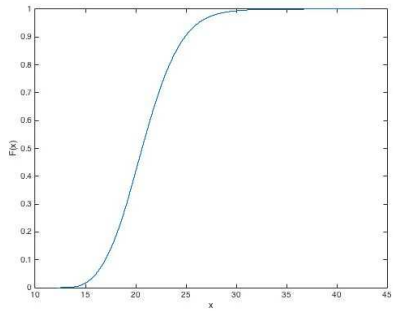
(a)  $n = 100$



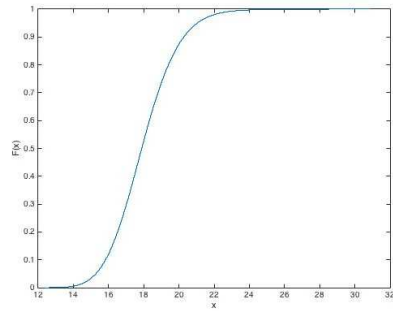
(b)  $n = 250$



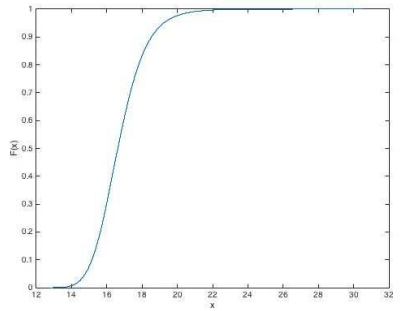
(c)  $n = 500$



(d)  $n = 1000$

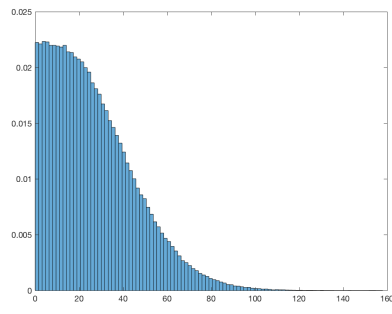


(e)  $n = 2500$

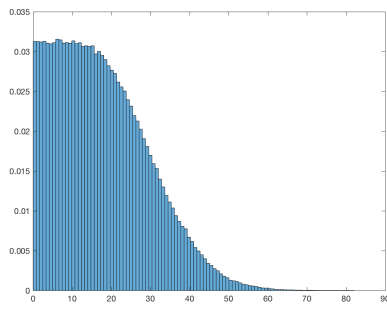


(f)  $n = 5000$

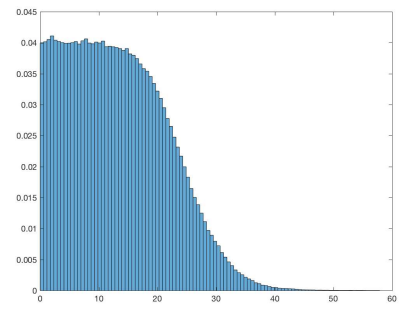
Figure 5: ECDF of periods between successive breaks estimated from computer simulation with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$



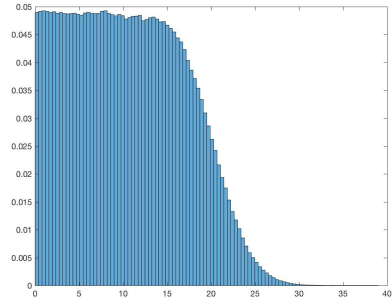
(a)  $n = 100$



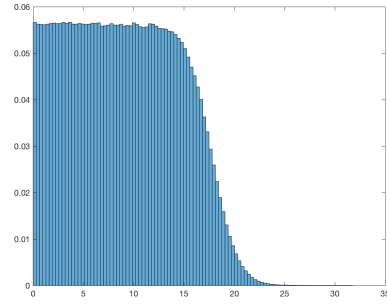
(b)  $n = 250$



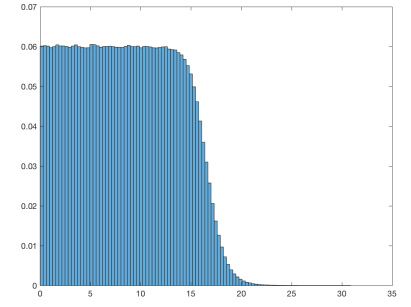
(c)  $n = 500$



(d)  $n = 1000$

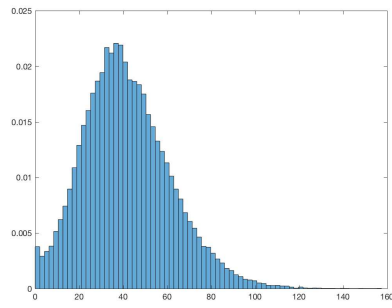


(e)  $n = 2500$

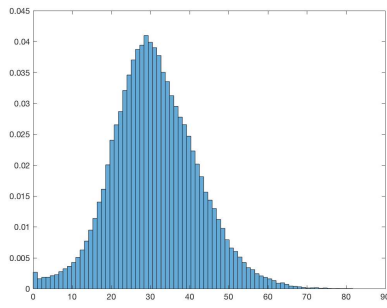


(f)  $n = 5000$

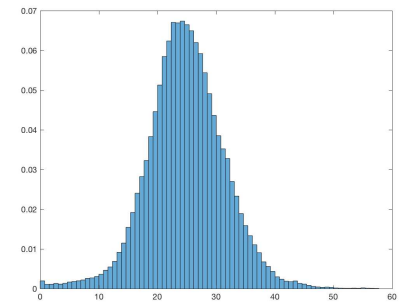
Figure 6: Histogram of age of a busy server estimated from computer simulation for  $M/M/n$  model with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$



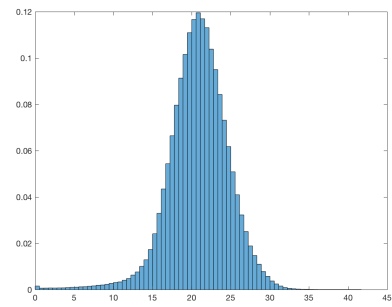
(a)  $n = 100$



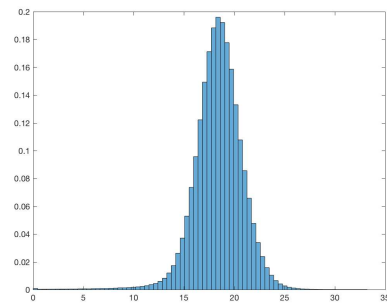
(b)  $n = 250$



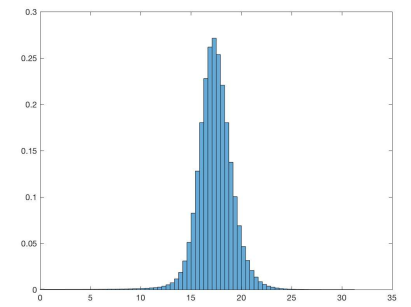
(c)  $n = 500$



(d)  $n = 1000$



(e)  $n = 2500$



(f)  $n = 5000$

Figure 7: Histogram of age of an idle server estimated from computer simulation for  $M/M/n$  model with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$

### 2.3 Impact of Non-Markovian Service-Times

In Section 4.5 of the main paper, we briefly described how non-exponential service time distribution affects the period between breaks,  $T_n$ . The present section expands the discussion and serves two purposes: (a) validate the fluid limit of time between successive breaks, i.e.,

$$T \stackrel{d}{=} x^* + R(x^*) + \theta = N(x^*) + 1 + \theta, \quad \mathbb{E}[T] = m(x^*) + 1 + \theta = 1/\beta, \quad \text{and} \quad \text{Var}(T) = \text{Var}(R(x^*)),$$

where  $N(\cdot), m(\cdot), R(\cdot), x^*$  and  $\beta$  are given in the paper when service times do not follow an exponential distribution; (b) expose the impact of non-exponential service times on the performance of rule  $D_1$ .

Here we assume the service time  $S$  to follow a hyper-exponential distribution or a mixture of exponential distributions; i.e., with probability  $\pi_j$ , the random variable  $S$  will take on the form of the exponential distribution with rate parameter  $\mu_j$ . A hyper-exponential r.v. does not enjoy memoryless property and possesses higher variability (see e.g., Whitt (1982)). Besides, it is particularly convenient to work because the expression of all moments is very explicit. Indeed, the moment-generating function

$$E[e^{tS}] = \int_0^\infty f_S(x)dx = \sum_{j=1}^K \pi_j \int_0^\infty e^{tx} \mu_j e^{-\mu_j x} dx = \sum_{j=1}^K \pi_j \mu_j / (\mu_j - t) \quad (2.1)$$

from which we can easily compute the first three moments

$$E[S] = \sum_{j \leq K} \pi_j / \mu_j, \quad E[S^2] = \sum_{j \leq K} 2\pi_j / \mu_j^2 \quad \text{and} \quad E[S^3] = \sum_{j \leq K} 6\pi_j / \mu_j^3.$$

Let  $S^e$  be its associated stationary-excess distribution. From the simple relationship between  $S$  and  $S^e$ , we have

$$E[S^e] = \frac{E[S^2]}{2E[S]} = \frac{\sum_{j \leq K} \pi_j / \mu_j^2}{\sum_{j \leq K} \pi_j / \mu_j} \quad \text{and} \quad E[(S^e)^2] = \frac{E[S^3]}{3E[S]} = \frac{\sum_{j \leq K} 2\pi_j / \mu_j^3}{\sum_{j \leq K} \pi_j / \mu_j}.$$

As a result,

$$\text{Var}[S^e] = E[(S^e)^2] - (E[S^e])^2 = \frac{\sum_{j \leq K} 2\pi_j / \mu_j^3}{\sum_{j \leq K} \pi_j / \mu_j} - \left( \frac{\sum_{j \leq K} \pi_j / \mu_j^2}{\sum_{j \leq K} \pi_j / \mu_j} \right)^2$$

For our purpose, we'd like to construct a r.v. with desired mean 1 and coefficient of variation  $c^2 = 4$ . To match the first two moments, it suffices to consider a 2-phase hyper-exponential distribution. In particular we use  $S = H_2$  with balanced means (see, e.g., Whitt (1982)). Applying the formulas in Whitt (1982), we get

$$\begin{aligned} \pi_1 &= \left( 1 + \sqrt{(c^2 - 1)/(c^2 + 1)} \right) / 2 = 0.8873, & \mu_1 &= 2\pi_1 = 1.7746, \\ \pi_2 &= \left( 1 - \sqrt{(c^2 - 1)/(c^2 + 1)} \right) / 2 = 0.1127, & \mu_2 &= 2\pi_2 = 0.2254. \end{aligned}$$

With these parameters, we can easily compute the mean and variance of the stationary-excess distribution  $H_2^e$ . We expect the variance of time between successive idle periods to be close to that of  $H_2^e$  for  $s$  large enough. Using the expression as shown above, we obtain the first two moments of  $H_2^e$ :

$$E[H_2^e] = \pi_1 / \mu_1^2 + \pi_2 / \mu_2^2 = 2.5 \quad \text{and} \quad E[(H_2^e)^2] = 2\pi_1 / \mu_1^3 + 2\pi_2 / \mu_2^3 = 20$$

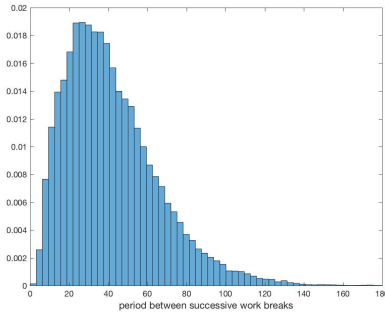
from which it follows

$$\text{std}(H_2^e) \equiv \sqrt{\text{Var}(H_2^e)} = \sqrt{20^2 - 2.5 \times 2.5} = 3.7081.$$

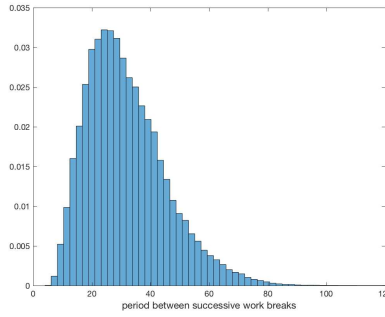
We elaborate in Table 5 by showing the 95% confidence intervals based on the data collected from the simulation experiments or the mean and the standard deviation of the interval between successive breaks, which we denoted by  $T_n$ . Consistent with the fluid limit with general service time distribution, the mean  $\mathbb{E}[T_n]$  decreases as  $n$  grows and converges from above to the limit  $\mathbb{E}[T] = 16.6667$ ; the standard deviation  $std(T_n)$  also decreases in  $n$  and has a tendency to converge to the correct limit, e.g.,  $std(T_n) = 3.876$  for  $s = 5000$ , very close to the theoretic value 3.7081. As a supplement, Figure 8 displays the histograms of  $T_n$  estimated from computer simulation as a function in  $n$ . Consistent with the fluid limit, the distribution of  $T_n$  converges to the suggested form, i.e., a shifted excess-lifetime distribution  $x^* + R(x^*) + \theta$ .

	$E[A_B]$	$std(A_B)$	$E[A_I]$	$std(A_I)$	$E[T_n]$	$std(T_n)$
$n = 100$	$27.145 \pm 0.098$	$22.059 \pm 0.102$	$37.622 \pm 0.106$	$23.851 \pm 0.115$	$41.663 \pm 0.126$	$23.7531 \pm 0.131$
$n = 250$	$18.277 \pm 0.085$	$13.584 \pm 0.092$	$29.473 \pm 0.089$	$13.922 \pm 0.079$	$31.748 \pm 0.095$	$13.473 \pm 0.104$
$n = 500$	$13.748 \pm 0.082$	$9.654 \pm 0.089$	$24.058 \pm 0.084$	$9.049 \pm 0.077$	$25.102 \pm 0.087$	$8.587 \pm 0.098$
$n = 1000$	$10.813 \pm 0.062$	$7.249 \pm 0.071$	$20.031 \pm 0.075$	$5.883 \pm 0.058$	$20.495 \pm 0.047$	$5.568 \pm 0.072$
$n = 2500$	$9.594 \pm 0.043$	$6.288 \pm 0.051$	$18.513 \pm 0.055$	$4.439 \pm 0.039$	$18.407 \pm 0.034$	$4.228 \pm 0.053$
$n = 5000$	$8.765 \pm 0.022$	$5.789 \pm 0.030$	$17.017 \pm 0.028$	$4.150 \pm 0.025$	$16.725 \pm 0.024$	$3.876 \pm 0.030$

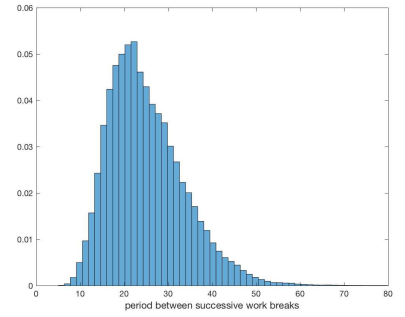
Table 5: Statistics for  $A_B, A_I$  and  $T_n$  for  $M/H_2/n$  model with rule  $D_1$ ,  $\rho = 0.9$  and  $\theta = 5/3$



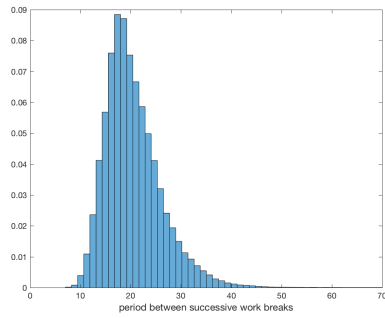
(a)  $n = 100$



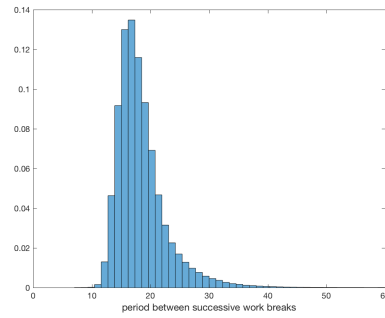
(b)  $n = 250$



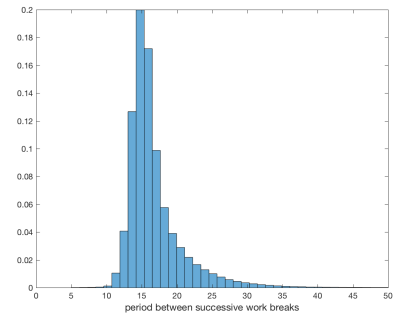
(c)  $n = 500$



(d)  $n = 1000$



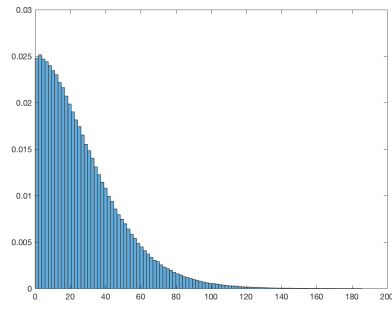
(e)  $n = 2500$



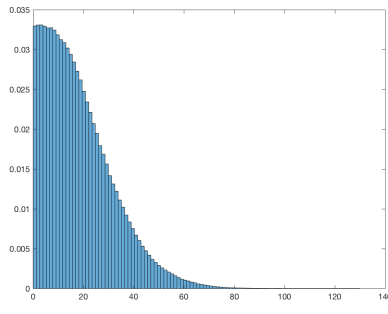
(f)  $n = 5000$

Figure 8: Histogram of  $T_n$  for  $M/H_2/n$  model with rule  $D_1$ ,  $\rho = 0.9$  and  $\theta = 5/3$

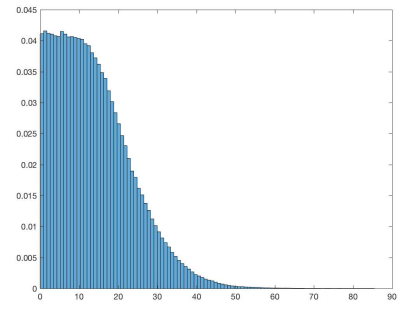




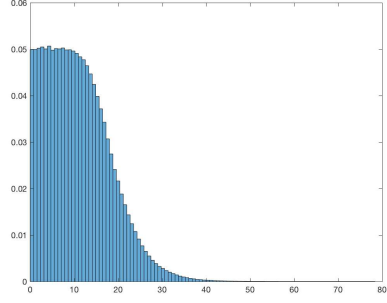
(a)  $n = 100$



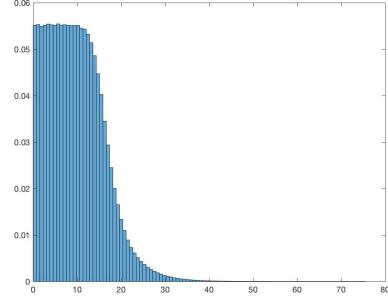
(b)  $n = 250$



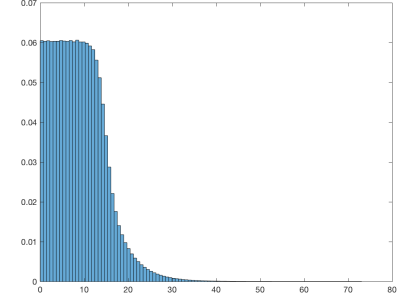
(c)  $n = 500$



(d)  $n = 1000$

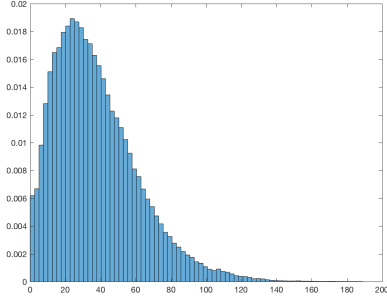


(e)  $n = 2500$

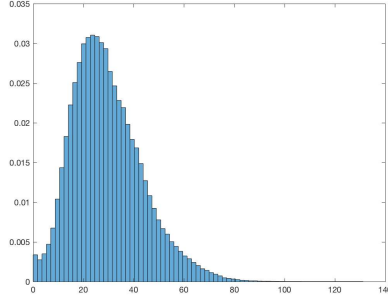


(f)  $n = 5000$

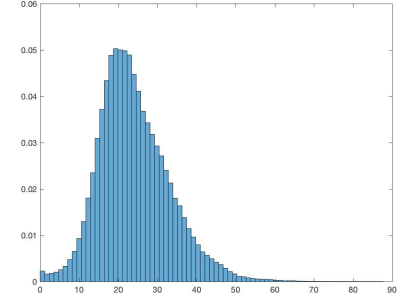
Figure 9: Histogram of age of a busy servers estimated from computer simulation for  $M/H_2/n$  model with rule  $D_1$ ,  $\rho = 0.9$  and  $\theta = 5/3$



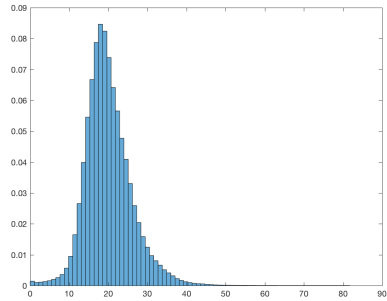
(a)  $n = 100$



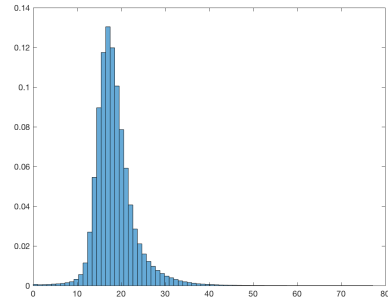
(b)  $n = 250$



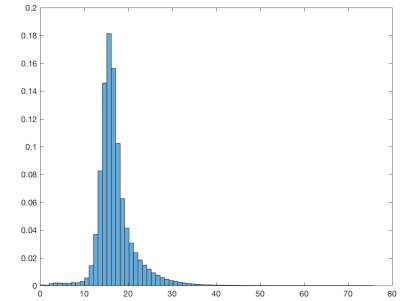
(c)  $n = 500$



(d)  $n = 1000$



(e)  $n = 2500$



(f)  $n = 5000$

Figure 10: Histogram of age of an idle server estimated from computer simulation for  $M/H_2/n$  model with rule  $D_1$  for  $\rho = 0.9$  and  $\theta = 5/3$

### 3 Additional Results on the $D_2$ Assignment Rule

The present section is to supplement Section 5 of the main paper. In §3.1 we examine the impact of the threshold parameter  $\eta$  on the system performance via appropriate sample paths. We present the simulation results for a large  $M/M/n$  queue in §3.3.

#### 3.1 Impact of the Parameter $\eta$

In Section 5.3 of the paper, we exposed the tradeoff in the choice of the parameter  $\eta$  through tables. In this section we show how the impact can be visualized through appropriate sample paths.

For greater insight, let  $X(t)$  the number of customers in system at time  $t$  and  $I_d(t) \equiv n - X(t)$  the number of idle servers at time  $t$ , allowing it to be negative as well as positive. Thus  $-I_d(t) = Q(t)$ , the queue length, when  $I_d(t) < 0$ , and  $I(t) = I_d(t)^+$ . We let  $S_b(t)$  be the number of servers on break at time  $t$ . Figure 11 displays sample paths of the number of servers on break,  $S_b(t)$ , and the number of idle servers,  $I_d(t)$ , for the base  $M/M/n$  model with  $n = 100$ ,  $\rho = 0.9$ , and four different values of  $\eta$  when  $\tau = 20$  and  $\theta = 5/3$ . Panel (c) with  $\eta = 8$  shows a severe performance degradation for customers because we often observe a big downward spike in  $I_d(t)$ , which suggests a buildup of large queue.

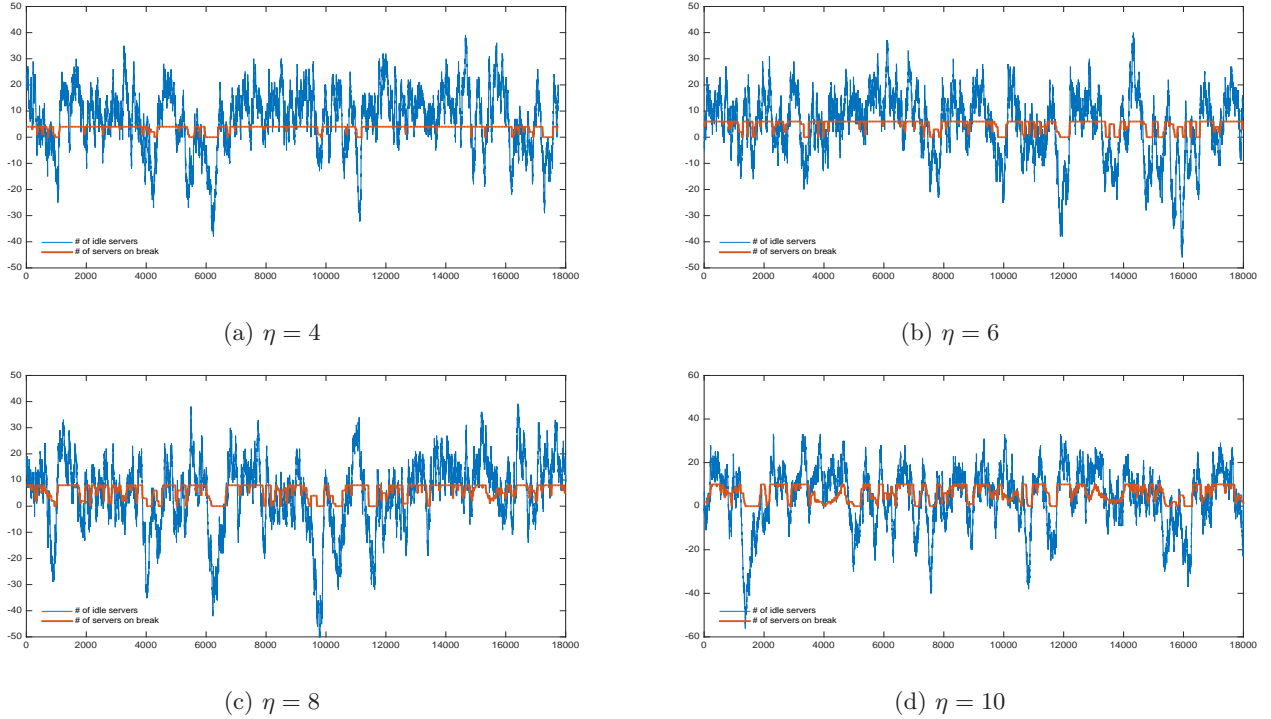


Figure 11: Sample paths of the number of servers on break,  $S_b(t)$  (red), and the number of idle servers,  $I_d(t) \equiv s - N(t)$  (blue), with rule  $D_2$  as a function of  $\eta$  for  $\theta = 5/3$  and  $\tau = 20$

### 3.2 A Small System

In the main paper, we formulate an optimization to choose the parameters  $\tau$  and  $\eta$ . In particular, we suggest performing a simple optimization with a cost function that is a convex weighted sum of  $1 - p_A$  and  $p_D - p_D^*$ , i.e.,

$$C \equiv C(p_A, p_D) \equiv w(1 - p_A) + (1 - w)(p_D - p_D^*), \quad 0 \leq w \leq 1, \quad (3.1)$$

where  $p_D^*$  is the LISF value, which is 0.223 for  $n = 100$  and 0.001 for  $n = 1000$  and the weight  $w$  reflects the relative cost we wish to attribute to  $p_A$  versus  $p_D$ . In the main paper, we displayed and discussed the simulation results for the base model with  $n = 100$ ,  $\rho = 0.9$  and weight  $w = 0.5$ . The present section provides addition numerical results with different weight parameters.

Figure 12 shows the cost  $C$  as a function of  $\tau$  and  $\eta$  for  $n = 100$ ,  $\theta = 5/3$  and four weights  $w = 0.3, 0.4, 0.6, 0.7$ . Panel 12a shows that for  $w = 0.3$  the optimal  $(\tau^*, \eta^*)$  is attained at  $(\tau = 40, \eta = 4)$ . Panel 12b shows that for  $w = 0.4$   $(\tau = 25, \eta = 6)$  is the optima. Panel 12c suggests for  $w = 0.6$  that  $(\tau = 15, \eta = 8)$  is the optimal combination while 12d show that for  $w = 0.7$ , the choice  $(\tau = 15, \eta = 10)$  is favorable.

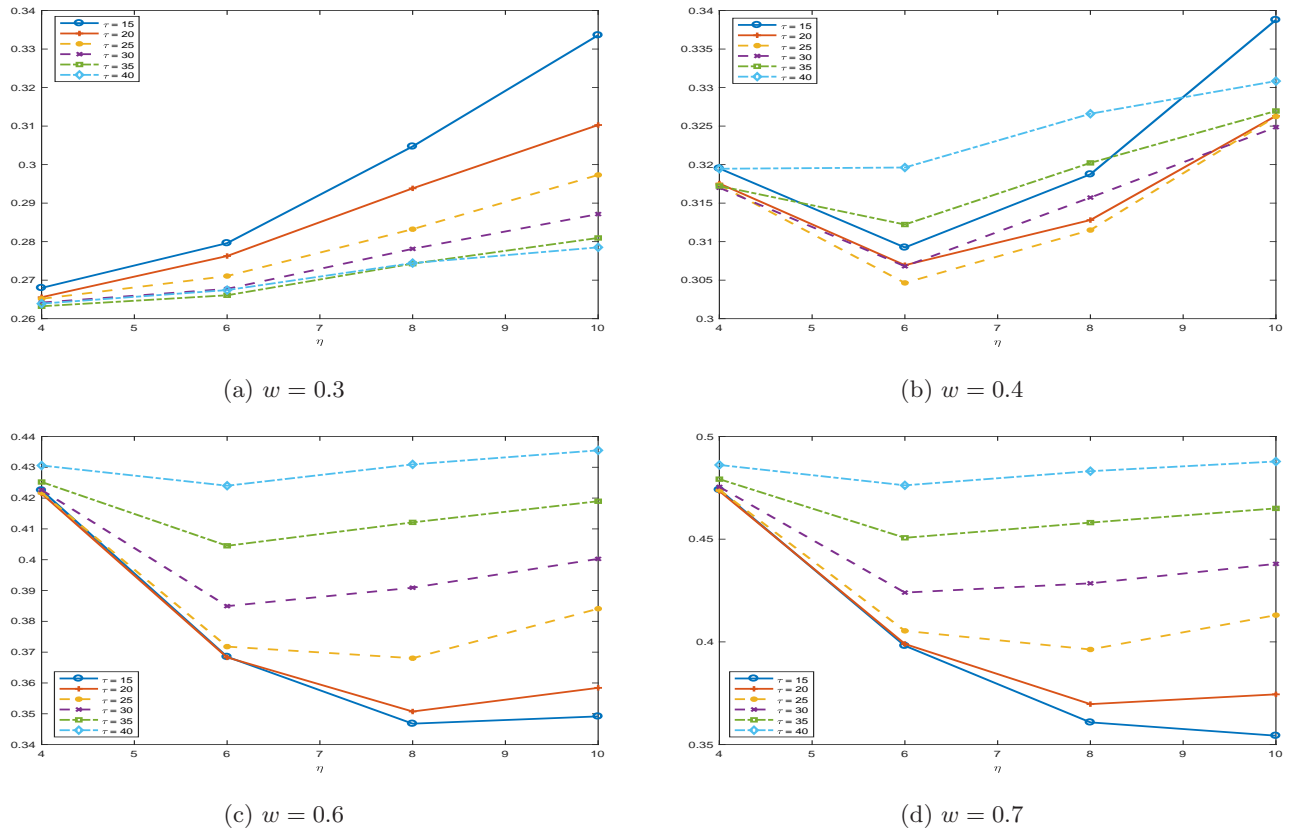


Figure 12: Cost function for  $D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $n = 100$ ,  $\theta = 5/3$  and  $w = 0.3, 0.4, 0.6, 0.7$

### 3.3 A Large System

We now amplify Remark 5.1 in the main paper by simulating a large  $M/M/n$  queue with  $n = 1000$  and  $\rho = 0.9$ .

Table 6 and 7 display two basic performance measures, namely the delay probability  $p_D$  and the proportion of idle on work breaks  $p_A$  as a function of  $\tau$  and  $\eta$ . Estimates of the 95% confidence intervals are displayed in the tables.

As the system size grows larger, both types of costs are significant reduced, and yet there is tradeoffs in the choice of  $(\tau, \eta)$ . Again both  $p_D$  and  $p_A$  are monotone in  $\tau$ , but  $p_A$  is not monotone in the bound  $\eta$ . For  $\tau$  ranging from 15 to 30, the largest value of  $p_A$  is attained at  $\eta = 90, 80, 70$  and 60 respectively. These values are highlighted in Table 7.

	$\eta = 4$	$\eta = 6$	$\eta = 8$	$\eta = 10$
$\tau$	$p_D$	$p_D$	$p_D$	$p_D$
$\tau = 15$	$0.0976 \pm 0.0024$	$0.1630 \pm 0.0029$	$0.2427 \pm 0.0052$	$0.3203 \pm 0.0030$
$\tau = 20$	$0.0892 \pm 0.0019$	$0.1364 \pm 0.0019$	$0.1487 \pm 0.0021$	$0.1611 \pm 0.0033$
$\tau = 25$	$0.0699 \pm 0.0016$	$0.0742 \pm 0.0015$	$0.0808 \pm 0.0023$	$0.0854 \pm 0.0023$
$\tau = 30$	$0.0422 \pm 0.0010$	$0.0529 \pm 0.0014$	$0.0612 \pm 0.0016$	$0.0720 \pm 0.0011$

Table 6: 95% confidence intervals of delay probability for rule  $D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$  and  $\theta = 5/3$ .

	$\eta = 4$	$\eta = 6$	$\eta = 8$	$\eta = 10$
$\tau$	$p_A$	$p_A$	$p_A$	$p_A$
$\tau = 15$	$0.7061 \pm 0.0013$	$0.7700 \pm 0.0011$	$0.8114 \pm 0.0009$	<b><math>0.8231 \pm 0.0018</math></b>
$\tau = 20$	$0.7033 \pm 0.0012$	$0.7631 \pm 0.0012$	<b><math>0.7808 \pm 0.0006</math></b>	$0.7773 \pm 0.0011$
$\tau = 25$	$0.6851 \pm 0.0015$	<b><math>0.6938 \pm 0.0014</math></b>	$0.6860 \pm 0.0024$	$0.6780 \pm 0.0024$
$\tau = 30$	<b><math>0.5850 \pm 0.0017</math></b>	$0.5593 \pm 0.0026$	$0.5422 \pm 0.0028$	$0.5368 \pm 0.0021$

Table 7: 95% confidence intervals for proportion of idle time spent on announced work breaks for rule  $D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$  and  $\theta = 5/3$ .

Solve the optimization problem (3.1) gives Figure 14 shows the cost  $C$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$ ,  $\theta = 5/3$  and four weights  $w = 0.3, 0.4, 0.6, 0.7$ . Panel 14b shows that for  $w = 0.4$  the combination  $(\tau = 80, \eta = 70)$  is optimal whereas Panel 14c shows that for  $w = 0.6$ , the choice  $(\tau = 70, \eta = 80)$  is more desirable.

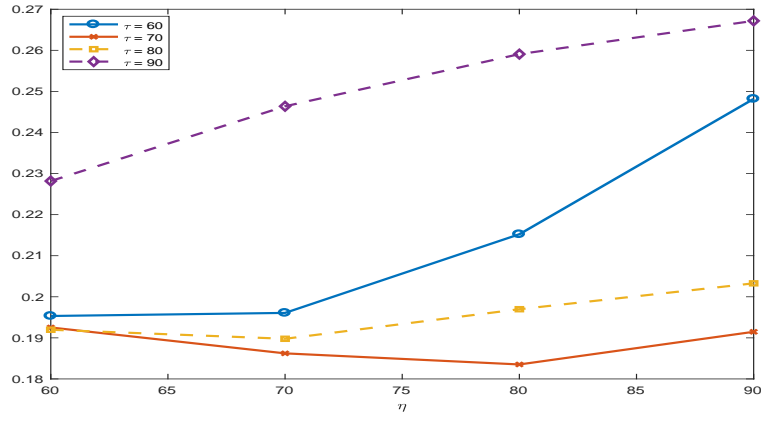
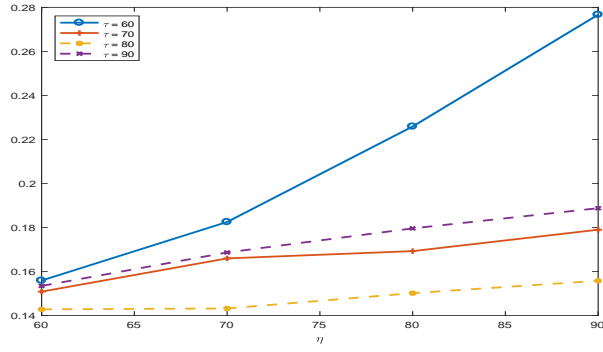
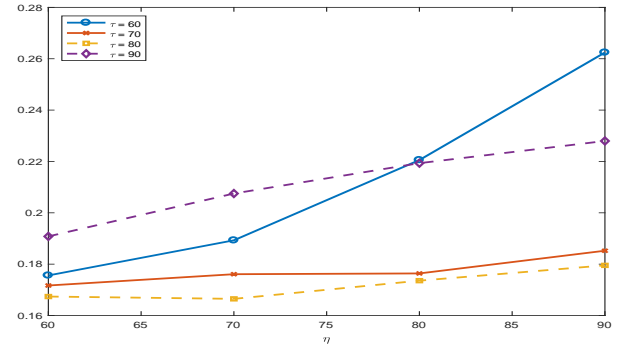


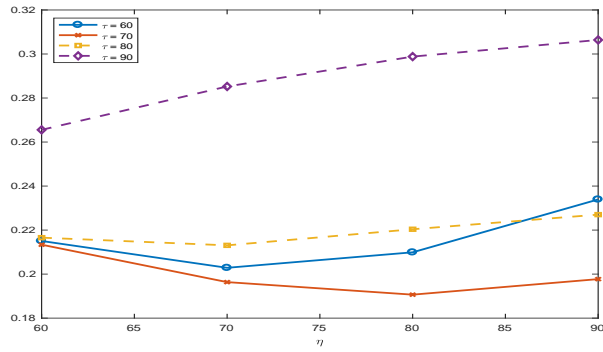
Figure 13: Cost function for  $D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $n = 1000$ ,  $\theta = 5/3$  and  $w = 0.5$



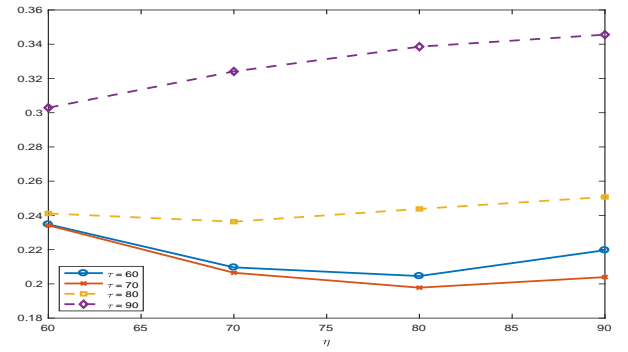
(a)  $w = 0.3$



(b)  $w = 0.4$



(c)  $w = 0.6$



(d)  $w = 0.7$

Figure 14: Cost function for  $D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $n = 1000$ ,  $\theta = 5/3$  and  $w = 0.3, 0.4, 0.6, 0.7$

### 3.4 Comparison with the Standard $M/M/(n - b)$ Model

An alternative way to obtain work breaks is to place a constant number of servers on break. If we put  $b$  servers on break at all times, then we obtain an  $M/M/(n - b)$  model with the customary LIFS server assignment rule. It is useful to compare  $D_2(\theta, \tau, \eta)$  to an  $M/M/(s - b)$  LIFS model by considering a range of  $b$  from  $\lfloor E[S_b] \rfloor$ , the greatest integer less than or equal to  $E[S_b]$ , to  $\eta$ . Table 8 displays  $E[S_b]$  as a function of  $\eta$  estimated from computation simulations.

We find that  $D_2(\theta, \tau, \eta)$  outperforms  $M/M/(s - b)$  LIFS for  $\lfloor E[S_b] \rfloor \leq b \leq \eta$ . For example, when  $\theta = 5/3$ ,  $\tau = 20$  and  $\eta = 8$ , Table 8 shows that  $E[S_b] = 5.610$ . Table 6 of the main paper shows that  $p_D = 0.460$  which is less than the values for  $b = 5$  in Table 9, and far less than the value for  $b = \eta = 8$ .

$\eta$	1	2	3	4	5	6	7	8	9
$E[S_b]$	0.908	1.789	2.630	3.422	4.166	4.834	5.358	5.610	5.708

Table 8: Estimated mean of the number of servers on announced breaks for the  $D_2$  rule as a function of  $\eta$  for  $\theta = 5/3$ ,  $\tau = 20$

$b$	0	1	2	3	4	5	6	7	8
$p_D$	0.216	0.257	0.304	0.358	0.420	0.488	0.564	0.648	0.737
$E[Q]$	1.90	2.51	3.33	4.45	6.00	8.23	11.54	16.68	25.19
$std(Q)$	5.62	6.69	8.01	9.65	11.75	14.49	18.20	23.48	31.00

Table 9: Performance measures for the standard  $M/M/(100 - b)$  queue with  $\rho = 0.9$

## 4 The $LISF - D_2$ Assignment Rule

In Remark 5.2 of the paper we mentioned an alternative to the  $D_2$  assignment rule that is easier to implement, and has similar performance. We now amplify this remark and describe in detail the alternative rule which we refer to as the  $LISF - D_2$  assignment rule. We show that its performance is similar to the  $D_2$  rule elaborated in the main paper.

Under the LISF rule there is a *FIFO queue for assignment*, i.e., whoever becomes idle first gets assigned first. Now we maintain two FIFO queues for assignment, a high priority queue (HPQ) and a low priority queue (LPQ). The rule stipulates that servers join the back of the HPQ once finishing a break.

### 4.1 Implementation of the $LISF - D_2$ Rule

There four types of events: customer arrival, customer departure (service completion), due for a break and work-break completion. We first explain how to treat the control parameter  $\tau$  with  $\eta = \infty$ , so it plays no role. Afterwards, we discuss the modifications to include  $\eta$ .

*At each arrival epoch*, we look for idle servers in the HPQ. If any, assign the server at the head of the HPQ. Otherwise we look for idle servers in the LPQ. If any, assign the server at the head of the LPQ. For the selected idle server, the algorithm generates a service requirement  $S$  and resets its service completion time to  $t + S$ . Then we find the minimum service-completion time among all busy servers and update the departure time accordingly. If there are no idle servers, the arriving customer waits in queue.

*At each departure epoch*, we look for customers in queue. If there is customer waiting, assign the server to the head-of-line customer. Otherwise the server either becomes idle or starts a break depending on whether or not a high priority designation was given. If a high priority designation was given, the break is announced and the server is not available to provide service for the duration  $\theta$  after that time. Otherwise it joins the back of the LPQ.

*At each break due time* (when a server's age reaches  $\tau$ ), if the server is busy, then we give the server a high priority designation indicating that that its next idle period will be replaced by an announced break. If the server is idle, then the server starts a break and goes off duty for the duration  $\theta$ .

*At each break-end time*, we first reset the server's age to zero. We assign to it a customer if there are customers in queue. Otherwise the server joins the back of HPQ. This prevents work break from being much greater than  $\theta$  since we always make assignment from the HPQ first.

We now discuss modifications to treat the bound  $\eta$ .

*Each time a break is due*, if the server is busy, we assign it a high priority designation. Meanwhile, we keep track of the elapsed time since this high priority designation has been assigned. If the server is idle and the number of off-duty servers is less than  $\eta$ , then a break is announced and the server is not available to provide service for the duration  $\theta$ . On the other hand, if the server is idle and the number of off-duty servers equals  $\eta$ , then we give the server a high-priority designation and do not make break announcement; and again keep track of the elapsed time since this high priority designation has been assigned..

*At each departure epoch*, if the queue is non-empty, then the server is assigned to the customer at the head of the queue. Hence suppose that the queue is empty. If a high priority designation was given *and* to the server and

the number off-duty servers is less than  $\eta$ , the break is announced and the server no longer provides service for the duration  $\theta$ . Otherwise the server joins the back of the LPQ.

At each break-end time, we first reset the server's age to zero. We assign to it a customer if there are customers in queue. Otherwise the server joins the back of HPQ. Meanwhile, we look for idle servers with high priority designation. If any, choose the one with the longest elapsed time since it receives this high priority level and announce the break.

## 4.2 A Small System

We now study the impact of the control parameter  $\tau$  and  $\eta$  for  $n = 100$  and  $\theta = 5/3$  with the  $LISF - D_2$  rule. Table 10 and 11 display two basic performance measures, namely the delay probability  $p_D$  and the proportion of idle on work breaks  $p_A$  as a function of  $\tau$  and  $\eta$ . In particular, estimates of the 95% confidence intervals are shown in the tables.

We see that there is a strong tradeoff in the choice of  $\eta$ , for a given  $\tau$ , between the effectiveness of the breaks for the servers and the performance experienced by customers. That tradeoff is dramatically in the two tables. For  $\tau = 20$  and  $\eta = 4$ , there is moderate performance degradation for customers with a delay probability 0.3353, but the algorithm for work breaks is ineffective, e.g., only a third of the available idleness is turned into work breaks. On the other hand for  $\tau = 20$  and  $\eta = 10$ , the algorithm is very effective in generating work breaks, i.e., more than half of the total idleness is turned into work breaks, but there is severe performance degradation for customers, e.g.,  $p_D$  increases by 50% reaching 0.4841.

Both  $p_D$  and  $p_A$  are monotone in  $\tau$ , but  $p_A$  is not monotone in the bound  $\eta$ . For  $\tau = 4 - 6, 7 - 8$  and  $9 - 10$ , the largest value of  $p_A$  is attained at  $\eta = 10, 8$  and  $6$  respectively. These values are highlighted in Table 11.

	$\eta = 4$	$\eta = 6$	$\eta = 8$	$\eta = 10$
$\tau$	$p_D$	$p_D$	$p_D$	$p_D$
$\tau = 15$	$0.3363 \pm 0.0023$	$0.4130 \pm 0.0025$	$0.4873 \pm 0.0023$	$0.5397 \pm 0.0021$
$\tau = 20$	$0.3353 \pm 0.0016$	$0.4065 \pm 0.0020$	$0.4594 \pm 0.0031$	$0.4841 \pm 0.0026$
$\tau = 25$	$0.3325 \pm 0.0020$	$0.3934 \pm 0.0025$	$0.4198 \pm 0.0022$	$0.4342 \pm 0.0024$
$\tau = 30$	$0.3323 \pm 0.0021$	$0.3728 \pm 0.0023$	$0.3853 \pm 0.0029$	$0.3979 \pm 0.0030$
$\tau = 35$	$0.3243 \pm 0.0021$	$0.3531 \pm 0.0017$	$0.3583 \pm 0.0016$	$0.3670 \pm 0.0026$
$\tau = 40$	$0.3195 \pm 0.0026$	$0.3314 \pm 0.0026$	$0.3381 \pm 0.0026$	$0.3446 \pm 0.0020$

Table 10: 95% confidence intervals for delay probability of  $LISF - D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 100$  and  $\theta = 5/3$ .

In order to choose the parameters  $\tau$  and  $\eta$ , we again solve the optimization with the cost function as shown in (3.1) Figure 16 shows the cost  $C$  as a function of  $\tau$  and  $\eta$  for  $n = 100$ ,  $\theta = 5/3$  and four weights  $w = 0.3, 0.4, 0.6, 0.6$ . Panel 16a shows that for  $w = 0.3$  the optimal  $(\tau^*, \eta^*)$  is attained at  $(\tau = 40, \eta = 4)$ . Panel 16b shows that for  $w = 0.4$  both  $(\tau = 25, \eta = 6)$  and  $(\tau = 30, \eta = 6)$  are optimal. Panel 16c - 16d show that for  $w \geq 0.6$ , the choice  $(\tau = 15, \eta = 10)$  is favorable.

Comparing Table 10 with Table 6 in the main paper, we see that the two approaches, i.e., the  $LISF - D_2$  and the  $D_2$  rule, are comparable in terms of performance degradation caused by enforced breaks; but putting Table 11



	$\eta = 4$	$\eta = 6$	$\eta = 8$	$\eta = 10$
$\tau$	$p_A$	$p_A$	$p_A$	$p_A$
$\tau = 15$	$0.3343 \pm 7 \times 10^{-4}$	$0.4548 \pm 6 \times 10^{-4}$	$0.5357 \pm 7 \times 10^{-4}$	<b><math>0.5726 \pm 7 \times 10^{-4}</math></b>
$\tau = 20$	$0.3323 \pm 7 \times 10^{-4}$	$0.4475 \pm 8 \times 10^{-4}$	$0.5081 \pm 9 \times 10^{-4}$	<b><math>0.5171 \pm 9 \times 10^{-4}</math></b>
$\tau = 25$	$0.3300 \pm 7 \times 10^{-4}$	$0.4308 \pm 9 \times 10^{-4}$	<b><math>0.4558 \pm 9 \times 10^{-4}</math></b>	$0.4527 \pm 9 \times 10^{-4}$
$\tau = 30$	$0.3270 \pm 4 \times 10^{-4}$	$0.3998 \pm 6 \times 10^{-4}$	<b><math>0.4035 \pm 7 \times 10^{-4}</math></b>	$0.3998 \pm 9 \times 10^{-4}$
$\tau = 35$	$0.3205 \pm 6 \times 10^{-4}$	<b><math>0.3620 \pm 6 \times 10^{-4}</math></b>	$0.3598 \pm 9 \times 10^{-4}$	$0.3542 \pm 8 \times 10^{-4}$
$\tau = 40$	$0.3097 \pm 7 \times 10^{-4}$	<b><math>0.3262 \pm 9 \times 10^{-4}</math></b>	$0.3225 \pm 9 \times 10^{-4}$	$0.3186 \pm 9 \times 10^{-4}$

Table 11: 95% confidence intervals for proportion of idle time spent on announced work breaks for rule  $LISF - D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 100$  and  $\theta = 5/3$ .

in contrast with Table 5 in the main paper, we observe that the  $D_2$  rule consistently outperforms the  $LISF - D_2$  rule by a small margin in terms of their effectiveness in generating announced breaks.

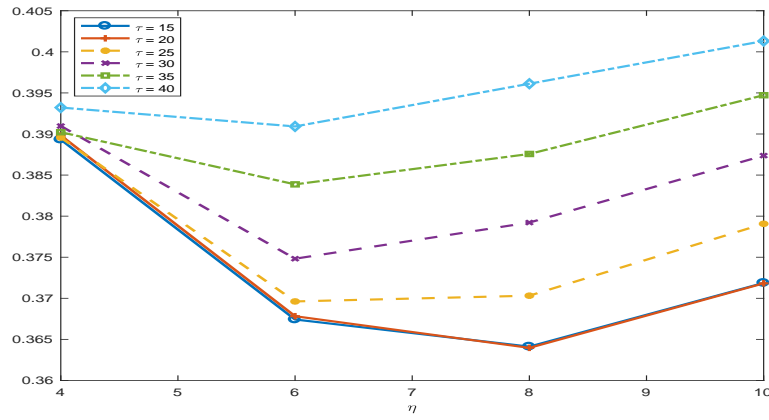


Figure 15: Cost function for  $LISF - D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $n = 100$ ,  $\theta = 5/3$  and  $w = 0.5$

### 4.3 A Large System

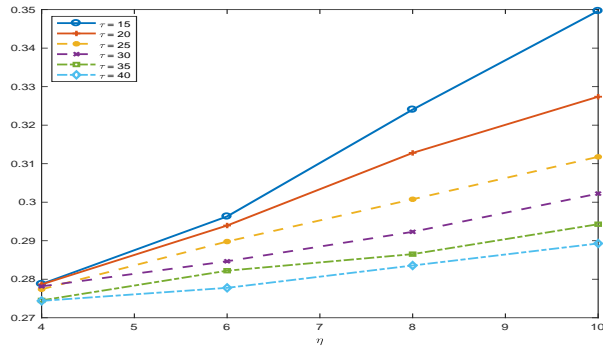
Here we fix traffic intensity  $\rho = 0.9$  and let system size  $n$  grows. Particularly we consider  $n = 1000$  and hence  $\lambda = 900$ .

Table 12 and 13 display two basic performance measures, namely the delay probability  $p_D$  and the proportion of idle on work breaks  $p_A$  as a function of  $\tau$  and  $\eta$  with the  $LISF - D_2$  rule. Estimates of the 95% confidence intervals are displayed in the tables.

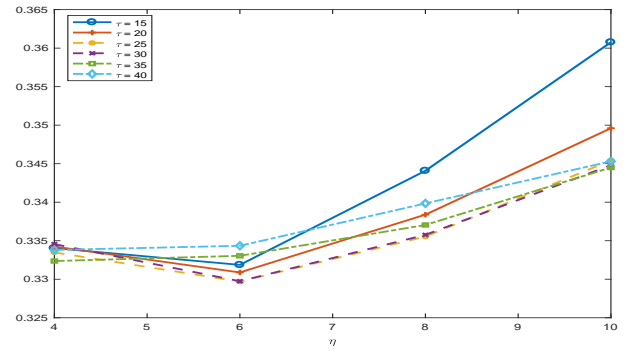
We see that with the system size being greater, both types of costs are significant reduced, and yet there is tradeoffs in the choice of  $(\tau, \eta)$ . When  $\tau = 15$ , for example, the choice of  $\eta$  exert a great influence on the effectiveness of the breaks for the servers and the performance experienced by customers.

Again both  $p_D$  and  $p_A$  are monotone in  $\tau$ , but  $p_A$  is not monotone in the bound  $\eta$ . For  $\tau$  ranging from 15 to 30, the largest value of  $p_A$  is attained at  $\eta = 90, 80, 70$  and  $60$  respectively. These values are highlighted in Table 13.

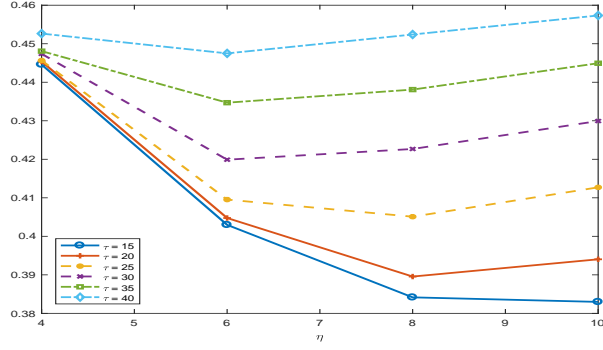
To determine the optimal  $(\tau^*, \eta^*)$ , we again solve the optimization problem with a cost function as given in (3.1).



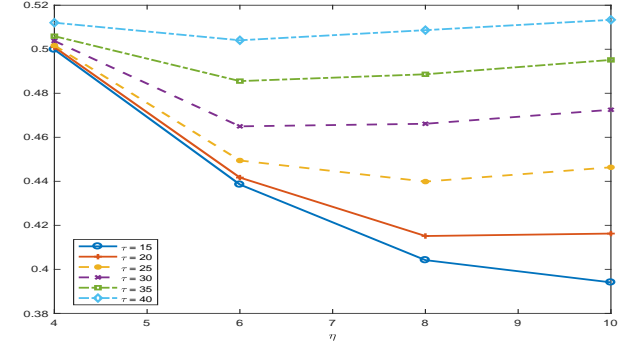
(a)  $w = 0.3$



(b)  $w = 0.4$



(c)  $w = 0.6$



(d)  $w = 0.7$

Figure 16: Cost function for  $LISF - D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $s = 100$ ,  $\theta = 5/3$  and  $w = 0.3, 0.4, 0.6, 0.7$

Figure 18 shows the cost  $C$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$ ,  $\theta = 5/3$  and four weights  $w = 0.3, 0.4, 0.6, 0.7$ . Panel 18a and 18d correspond to two extreme cases which are not very interesting. Panel 18b shows that for  $w = 0.4$  the combination  $(\tau = 80, \eta = 70)$  is optimal. whereas Panel 18c shows that for  $w = 0.6$ , the choice  $(\tau = 70, \eta = 80)$  is more desirable.

	$\eta = 60$	$\eta = 70$	$\eta = 80$	$\eta = 90$
$\tau$	$p_D$	$p_D$	$p_D$	$p_D$
$\tau = 15$	$0.0946 \pm 0.0022$	$0.1628 \pm 0.0032$	$0.2471 \pm 0.0035$	$0.3162 \pm 0.0034$
$\tau = 20$	$0.0907 \pm 0.0022$	$0.1348 \pm 0.0023$	$0.1512 \pm 0.0034$	$0.1624 \pm 0.0031$
$\tau = 25$	$0.0706 \pm 0.0018$	$0.0731 \pm 0.0016$	$0.0765 \pm 0.0018$	$0.0803 \pm 0.0018$
$\tau = 30$	$0.0407 \pm 0.0011$	$0.0430 \pm 0.0017$	$0.0433 \pm 0.0015$	$0.0435 \pm 0.0014$

Table 12: 95% confidence intervals for delay probability of  $LISF - D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$  and  $\theta = 5/3$ .

	$\eta = 60$	$\eta = 70$	$\eta = 80$	$\eta = 90$
$\tau$	$p_A$	$p_A$	$p_A$	$p_A$
$\tau = 15$	$0.5951 \pm 0.0015$	$0.6766 \pm 0.0016$	$0.7359 \pm 0.0011$	<b><math>0.7663 \pm 0.0015</math></b>
$\tau = 20$	$0.5940 \pm 0.0013$	$0.6642 \pm 0.0011$	<b><math>0.6822 \pm 0.0010</math></b>	$0.6797 \pm 0.0009$
$\tau = 25$	$0.5771 \pm 0.0009$	<b><math>0.5823 \pm 0.0009</math></b>	$0.5807 \pm 0.0018$	$0.5800 \pm 0.0016$
$\tau = 30$	<b><math>0.5015 \pm 0.0014</math></b>	$0.5007 \pm 0.0017$	$0.5001 \pm 0.0014$	$0.4994 \pm 0.0013$

Table 13: 95% confidence intervals for proportion of idle time spent on announced work breaks for rule  $LISF - D_2(\theta, \tau, \eta)$  as a function of  $\tau$  and  $\eta$  for  $n = 1000$  and  $\theta = 5/3$ .

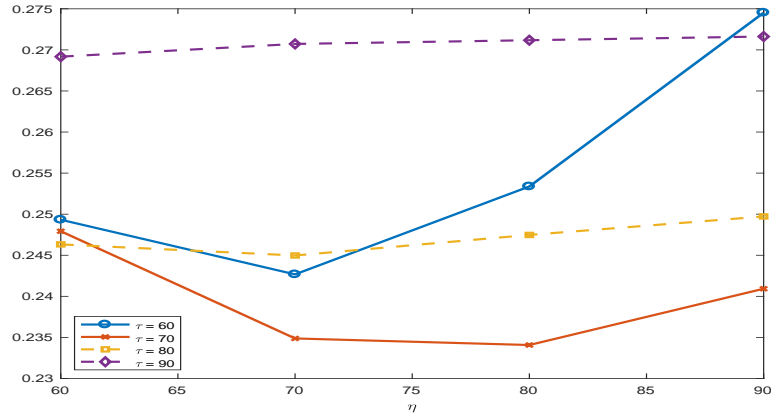
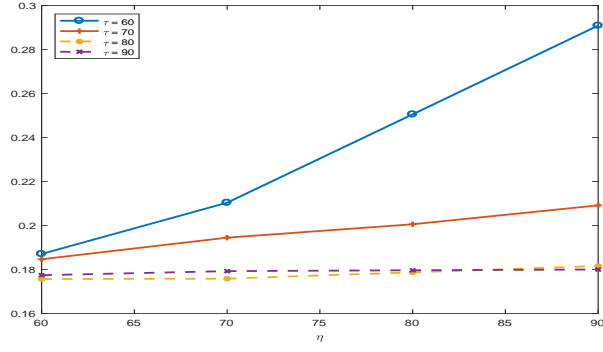
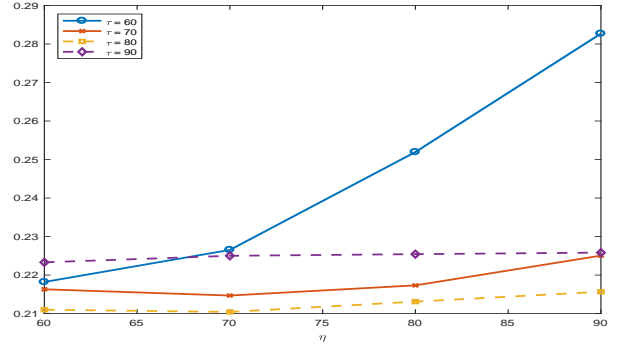


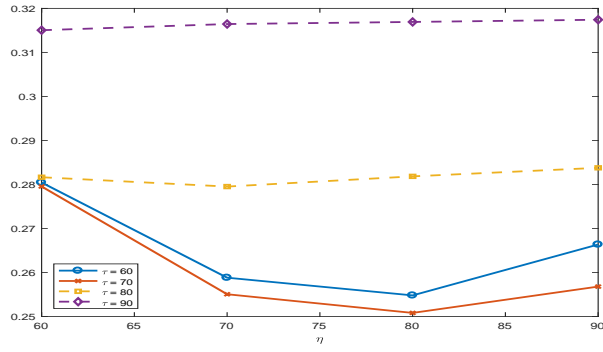
Figure 17: Cost function for  $LISF - D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $n = 1000$ ,  $\theta = 5/3$  and  $w = 0.5$



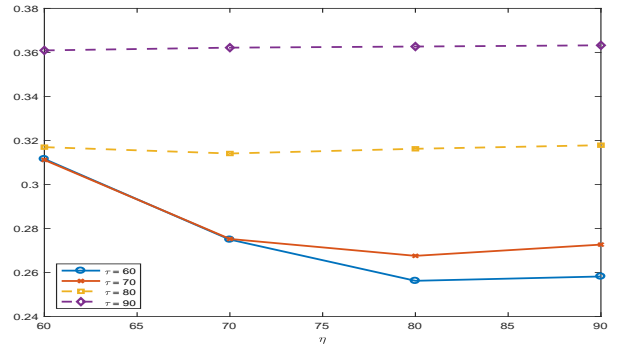
(a)  $w = 0.3$



(b)  $w = 0.4$



(c)  $w = 0.6$



(d)  $w = 0.7$

Figure 18: Cost function for  $LISF - D_2(\theta, \tau, \eta)$  in (3.1) as a function of  $\tau$  and  $\eta$  for  $s = 1000$ ,  $\theta = 5/3$  and  $w = 0.3, 0.4, 0.6, 0.7$

## References

W. Whitt. Approximating a point process by a renewal process, I: two basic methods. *Oper. Res.*, 30:125–147, 1982.