# Speed Scaling under QoS constraints
# with Finite Buffer

Parikshit Hegde
*Indian Institute of Technology Madras, India*
ee14b123@ee.iitm.ac.in

Akshit Kumar
*Indian Institute of Technology Madras, India*
ee14b127@ee.iitm.ac.in

Rahul Vaze
*Tata Institute of Fundamental Research, Mumbai, India*
vaze@tcs.tifr.res.in

*Abstract*—**A single server with variable speed and a finite buffer is considered under a maximum packet drop probability constraint. The cost of processing by the server is a convex function of the speed of the server. If a packet arrives when the buffer is full, it is dropped instantaneously. Given the finite server buffer, the objective is to find the optimal dynamic server speed to minimize the overall cost subject to the maximum packet drop probability constraint. Finding the exact optimal solution is known to be hard, and hence algorithms with provable approximation bounds are considered. We show that if the buffer size is large enough, the proposed algorithm achieves the optimal performance. For arbitrary buffer sizes, constant approximation guarantees are derived for a large class of packet arrival distributions such as Bernoulli, Exponential, Poisson etc.**

*Index Terms*—**Speed Scaling, Congestion Control, Queueing**

## I. INTRODUCTION

In a speed scaling problem, a server with tunable speed is considered, and the objective is to minimize an appropriate cost function subject to a quality of service (QoS) constraint (e.g., packet drop probability). In this paper, we consider a single server that is equipped with a finite buffer, and any packet that encounters a full buffer on its arrival is dropped instantaneously. The specific speed scaling problem is to minimize the server energy cost (convex function of the speed) under a packet drop probability constraint. This explicit problem was first studied in [1], where some structural results were obtained via results on Markov decision processes. The main analytical results were, however, derived under a slightly simpler QoS constraint, since the hard packet drop probability constraint is challenging to analyze.

One direction for analyzing speed scaling/job scheduling (choosing server speed depending on queue length) with finite but large buffers has been via large deviations exponents, where the rate-function optimality is the objective [2], [3]. Even though this provides theoretical answers, however, the asymptotic limits have limited scope in terms of practicality. For infinite buffer size, the tradeoff between average server speed and average delay can be found in [4]–[7], where the problem is cast as a Markov decision process, and structural results are obtained to reduce the search space for finding the optimal policy.

The authors are listed in last name alphabetical order.

Speed scaling has been also been considered in computer science theory literature [8]–[13], where a single/multiple servers have to choose their speed in order to maximize the utility (profit-service cost) subject to jobs having hard deadlines. Most of the results in [8]–[13] are derived under a worst case input, where arrival times and job sizes can be chosen by an adversary rather than being drawn from a stochastic process.

The speed scaling problem with finite buffers is also related to a more modern problem in energy harvesting (EH) paradigm, where a transmitter harvests random amount of energy from renewable sources, and its objective is to maximize its rate of transmission [14], [15]. The transmitter is equipped with a finite battery, and similar to a server with a finite buffer, any energy arriving when the battery is full is lost. The transmission rate of the transmitter is a concave function of the transmit power, and it has to decide the optimal rate of power transmission to maximize its rate, subject to the finite battery capacity and random energy arrivals. The rate maximization problem with EH, however, does not have any energy loss constraint. An additional constraint of limiting the probability of battery state hitting the full or empty state is also studied in [15].

In this paper, we consider the problem of choosing the optimal server speed for a single server with finite buffer under a hard packet drop probability constraint of $\alpha$, with convex server cost. Since finding the optimal policy is known to be hard [1], we take an alternate approach of deriving algorithms that have provable theoretical gap guarantees on their performance compared to the unknown optimal algorithm.

We consider a slotted time system, where $A_t$ packet arrive in slot $t$, and $A_t$ is assumed to be i.i.d. across slots $t$. We first consider arbitrary buffer sizes, and propose a simple policy that serves $(1-\alpha)A_t$ packets at time $t$ and forcibly drops the rest of the $\alpha A_t$ packets. This policy trivially satisfies the packet drop probability constraint, and we show that for large class of arrival distributions, that includes, uniform, exponential etc., this policy has at most a constant multiplicative gap in the server cost compared to the optimal policy. This policy can be shown to have a poor performance when the average arrival rate $\mathbb{E}[A_t]$ is low for certain distributions such as Bernoulli. For low $\mathbb{E}[A_t]$ regime, we propose an alternate algorithm and derive a tighter lower bound, and show that they have a constant

multiplicative gap between them for Bernoulli distributions. Finally, we consider the finite but large buffer regime, where we propose an algorithm similar to [15] and show that it has an additive gap of $\max\{\Theta\left(B^{-\beta}\right), \Theta\left(\frac{(\log B)^2}{B^2}\right)\}$[1] in terms of server cost from the optimal policy, when the packet drop probability can be violated by $\Theta\left(B^{-\beta}\right)$, where $B$ is the buffer size. Thus, with a reasonably sized buffers, the performance of the proposed policy is very close to the optimal with very small violation of the packet drop probability constraint. Since finding the exact optimal policy has remained open for long, this is the closest approximation result one can expect.

## II. SYSTEM MODEL

We consider a slotted time system, where in slot $t$, $A_t$ packet arrive to the buffer of the single server. The arrival process $\{A_t, t \geq 1\}$ is assumed to be an ergodic stochastic process with a long term mean given as $\lim_{\tau \to \infty} \frac{1}{\tau} \sum_{t=1}^{\tau} A_t = \mu$. We assume that the support of $A_t$ is in $[0, B]$, because otherwise it will not be possible for any policy to be feasible for sufficiently low values of packet drop probability. Packets are admitted subject to the finite buffer size of $B$, i.e., if the buffer state (number of packets in the buffer) at slot $t$ is $b_t$, then only $B - b_t$ packets among the $A_t$ packets are admitted into the buffer and the rest (without distinguishing) are dropped instantaneously. The server uses a policy/algorithm $\mathbf{g} = \{g_t\}_{t=1}^{\infty}$, where it serves $g_t$ packets in slot $t$, incurring a cost of $f_c(g_t)$, where $f_c(.)$ is a convex function. For most of the paper, we will consider $f_c(s) = s^2$, a common choice in literature [16]. We also assume that a fraction of a packet can be served in a slot $t$. The packet drop probability is denoted as $P_{\text{Drop}} = P(\text{an arriving packet is dropped})$. Policy $\mathbf{g}$ directly controls the overall cost and $P_{\text{Drop}}$, and the objective is to derive an optimal policy $\mathbf{g}^\star$ that minimizes the overall cost

$$\mathcal{J}_{\mathbf{g}} = \lim_{n \to \infty} \mathbb{E}\left[\frac{1}{n} \sum_{t=1}^{n} f_c(g_t)\right]. \tag{1}$$

subject to a hard constraint of $\alpha$ on $P_{\text{Drop}}$, where the expectation is over the arrival process $\{A_t\}_{t=1}^{\infty}$. Formally, we want to solve,

$$\begin{aligned}
\underset{\mathbf{g}}{\text{minimize}} \quad & \mathcal{J}_{\mathbf{g}} = \lim_{n \to \infty} \mathbb{E}\left[\frac{1}{n} \sum_{t=1}^{n} f_c(g_t)\right] \\
\text{subject to} \quad & 0 \leq g_t \leq b_t, t \geq 1, \\
& b_t = \min\{b_{t-1} + A_t - g_{t-1}, B\}, \\
& P_{\text{Drop}} \leq \alpha.
\end{aligned} \tag{2}$$

As described in the Introduction, problem (2) is hard to solve, and instead we seek to find policies that have a guaranteed multiplicative gap from the unknown optimal policy $\mathbf{g}^\star$. The multiplicative gap is generally referred to as the competitive ratio of the policy $\mathbf{g}$, that is defined as

$$\text{CR}_{\mathbf{g}} = \frac{\mathcal{J}_{\mathbf{g}}}{\mathcal{J}_{\mathbf{g}^\star}}. \tag{3}$$

[1] Notation used throughout the paper: $f_n = \Theta(g_n)$ if $f_n$ and $g_n$ go to zero at the same rate, $f_n = \Omega(g_n)$ if $f_n$ goes to zero no faster than $g_n$, $f_n = o(g_n)$ if $f_n$ goes to zero strictly faster than $g_n$

In the rest of the paper, we propose feasible policies $\mathbf{g}$ that have constant competitive ratios. Towards that end, we first rewrite the packet drop probability as the ratio of the expected number of packets dropped and the expected number of packet arrivals, i.e.,

$$P_{\text{Drop}} = \lim_{n \to \infty} \frac{\mathbb{E}\left[\sum_{t=1}^{n} \max\{0, A_t - B + b_t\}\right]}{\mathbb{E}\left[\sum_{t=1}^{n} A_t\right]}. \tag{4}$$

In the next subsection, we derive a lower bound on the cost $\mathcal{J}_{\mathbf{g}}$ incurred by any *feasible policy* $\mathbf{g}$ including the optimal policy $\mathbf{g}^\star$.

## III. GENERAL LOWER BOUND

**Theorem 1.** *Consider any continuous, convex, non decreasing cost function $f_c(.)$, the minimum cost incurred by any* feasible *policy $\mathbf{g}$ under any i.i.d. packet arrival process $\{A_t, t \geq 1\}$ is lower bounded as*

$$\mathcal{J}_{\mathbf{g}} \geq f_c\left(\mu\left(1 - \alpha\right)\right),$$

*where $\mu = \mathbb{E}[A_t]$, and $\alpha$ is the packet drop probability constraint*

*Proof.* The total number of packets served by a policy $\mathbf{g}$ in $n$ slots is $\sum_{t=1}^{n} g_t$, and the maximum number of packets remaining in the buffer after $n$ slots is $B$. Since the server is allowed to drop at most $\alpha$ fraction of the arrived packets, on average, the sum of total number of packets serviced $\sum_{t=1}^{n} g_t$ and the number of packets left in the buffer ($\leq B$) at the end of the packet arrival sequence need to be greater than $(1 - \alpha)$ fraction of the number of packet arrivals. Therefore, we get the lower bound

$$\lim_{n \to \infty} \mathbb{E}\left[\sum_{t=1}^{n} g_t + B\right] \geq (1 - \alpha) \mathbb{E}\left[\sum_{t=1}^{n} A_t\right],$$

$$\lim_{n \to \infty} \mathbb{E}\left[\frac{\sum_{t=1}^{n} g_t}{n} + \frac{B}{n}\right] \overset{(a)}{\geq} (1 - \alpha) \mathbb{E}\left[\frac{\sum_{t=1}^{n} A_t}{n}\right],$$

$$\lim_{n \to \infty} \mathbb{E}\left[\frac{\sum_{t=1}^{n} g_t}{n}\right] \overset{(b)}{\geq} (1 - \alpha) \mathbb{E}\left[\frac{\sum_{t=1}^{n} A_t}{n}\right],$$

$$\overset{(c)}{=} (1 - \alpha)\mu. \tag{5}$$

where (a) follows from dividing both sides by $n$, (b) follows since as $n \to \infty$, $\frac{B}{n} \to 0$ and, (c) follows since $A_t$ is i.i.d. and by linearity of expectation.

Since the cost function $f_c$ is convex, we next use Jensen's inequality to obtain a lower bound on the expected cost as follows.

$$\begin{aligned}
\mathcal{J} &\overset{(a)}{=} \mathbb{E}\left[\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} f_c(g_t)\right], \\
&\overset{(b)}{\geq} f_c\left(\mathbb{E}\left[\lim_{n \to \infty} \frac{1}{n} \sum_{t=1}^{n} g_t\right]\right), \\
&\overset{(c)}{\geq} f_c((1 - \alpha)\mu) \triangleq \mathcal{J}^*. 
\end{aligned} \tag{6}$$

where (a) is the definition of the cost of a policy, (b) follows from the convexity of $f_c$ and Jensen's inequality, and (c)

follows from (5). We define the lower bound on optimal cost $\mathcal{J}^*$. $\qquad\square$

**Corollary 1.** *For the cost function considered in this paper $f_c(s) = s^2$, the minimum cost incurred by any feasible policy* **g** *under any i.i.d. packet arrival process $\{A_t, t \geq 1\}$ is lower bounded as*

$$\mathcal{J} \geq (\mu\,(1 - \alpha))^2,$$

*where $\mu = \mathbb{E}[A_t]$ and $\alpha$ is the packet drop probability constraint.*

In the next subsection, we propose a simple greedy policy and bound its competitive ratio for $f_c(s) = s^2$.

## IV. GREEDY POLICY

Greedy policy: Out of $A_t$ newly arrived packets in slot $t$, service $g_t = (1 - \alpha)A_t$ packets in slot $t$ and drop the rest of $\alpha A_t$ packets. Greedy policy clearly transmits a $(1-\alpha)$ fraction of all the packets that arrive, since no packets are accumulated in the buffer at any time slot, and therefore satisfies the $P_{\text{Drop}} \leq \alpha$ constraint.

**Theorem 2.** *The competitive ratio of the greedy policy for $f_c(s) = s^2$ is*

$$CR_{greedy} \leq 1 + \frac{var(A_t)}{\mu^2}.$$

*Proof.* Let $\mathcal{J}_{\text{greedy}}$ denote the cost incurred by the greedy policy.

$$\begin{aligned}
\mathcal{J}_{\text{greedy}} &\overset{(a)}{=} \lim_{n\to\infty} \mathbb{E}\left[\frac{1}{n}\sum_{t=1}^{n} g_t^2\right], \\
&\overset{(b)}{=} \lim_{n\to\infty} \mathbb{E}\left[\frac{1}{n}\sum_{t=1}^{n} (1-\alpha)^2 A_t^2\right], \\
&\overset{(c)}{=} (1-\alpha)^2 \lim_{n\to\infty} \frac{1}{n}\sum_{t=1}^{n} \mathbb{E}\left[A_t^2\right], \\
&\overset{(d)}{=} (1-\alpha)^2 \left(\mu^2 + \text{var}(A_t)\right). \qquad (7)
\end{aligned}$$

where (a) follows by definition of cost of policy (b) follows from the definition of the greedy policy (c) follows from the linearity of expectation (d) follows from definition of variance.

To get the competitive ratio, divide the cost incurred by our online greedy policy $\mathcal{J}_{\text{greedy}}$ by the lower bound on the optimal cost $\mathcal{J}^*$ (Corollary 1),

$$CR_{\text{greedy}} \leq \frac{\mathcal{J}_{\text{greedy}}}{\mathcal{J}^*} = 1 + \frac{\text{var}(A_t)}{\mu^2}.$$

$\qquad\square$

**Remark 3.** *If $f_c(s) = s^\gamma$ for some $\gamma > 2$, we will get a competitive ratio for the greedy policy involving higher order moments of $A_t$.*

Next, we consider some well known distributions with a finite support $[0, B]$, and calculate the competitive ratio of the greedy policy given by Theorem 2 under these arrival distributions for $f_c(s) = s^2$.

### A. Uniform Distribution

Consider a continuous uniform packet arrival distribution $\{A_t\}$ with a probability density function defined as:

$$f_A(x) = \begin{cases} \frac{1}{B}, & 0 \leq x \leq B, \\ 0, & \text{otherwise.} \end{cases}$$

Under the uniform distribution, the competitive ratio of the greedy policy is

$$\text{CR}_{\text{greedy}} \leq 1 + \frac{\text{var}(A_t)}{\mu^2} = 1 + \frac{B^2/12}{(B/2)^2} = 1 + \frac{1}{3} = \frac{4}{3}.$$

Moreover, for any Uniform$[a, b]$ distribution where $0 \leq a \leq b \leq B$, it is easy to see that the competitive ratio of the greedy policy is no greater than $\frac{4}{3}$.

### B. Truncated Exponential Distribution

Consider a truncated exponential packet arrival distribution $\{A_t\}$ with a finite support on $[0, B]$, with the probability density function defined as

$$f_A(x) = \begin{cases} \frac{\frac{1}{\mu}e^{-x/\mu}}{1-e^{-B/\mu}}, & 0 \leq x \leq B, \\ 0, & \text{otherwise.} \end{cases}$$

Without loss of generality (WLOG), let $B = k\mu$ for some $k > 1$, then $\mathbb{E}[X] = \mu\left[\frac{1-(k+1)e^{-k}}{1-e^{-k}}\right]$ and $\mathbb{E}\left[X^2\right] = 2\mu^2\left[\frac{1-\frac{1}{2}\left(k^2+2k+2\right)e^{-k}}{1-e^{-k}}\right]$. Hence, the competitive ratio of the greedy policy is

$$\text{CR}_{\text{greedy}} \leq 2\frac{\left(1 - \frac{1}{2}\left(k^2 + 2k + 2\right)e^{-k}\right)\left(1 - e^{-k}\right)}{\left(1 - (k+1)e^{-k}\right)^2} \leq 2.$$

### C. Truncated Poisson Distribution

Consider a truncated Poisson arrival distribution $\{A_t\}$ with parameter $\nu$ with the support on $[0, B]$. Then the competitive ratio of the greedy policy can be upper bounded by

$$\text{CR}_{\text{greedy}} \leq 1 + \frac{2}{\nu}.$$

Proof is omitted for lack of space, which essentially follows from first principles. Notice that for $\nu \geq 1$, the greedy algorithm is at most 3-competitive. However, in the $\nu \to 0$ regime, the competitive ratio of the greedy policy grows unbounded. We illustrate a similar behavior for the Bernoulli distribution next, where the competitive ratio of the greedy policy is small when $\mu$, the expected number of packets arriving to the server, is moderate, but grows as $\mu$ decreases.

### D. Bernoulli Distribution

Consider a class of Bernoulli arrival distribution $\{A_t\}$ with probability mass function defined as

$$A_t = \begin{cases} m, & \text{w.p. } \frac{\mu}{m}, \\ 0, & \text{w.p. } 1 - \frac{\mu}{m}, \end{cases}$$

where $0 \leq m \leq B$, $\mathbb{E}[A_t] = \mu$ and $\mathbb{E}\left[A_t^2\right] = m\mu$. Thus, the competitive ratio of the greedy policy is given by

$$\text{CR}_{\text{greedy}} \leq 1 + \frac{\text{var}[X]}{\mu^2} = \frac{m}{\mu}. \tag{8}$$

When $\mu \geq \frac{m}{2}$, (8) shows that the greedy policy is at most 2-competitive, where as if $\mu \to 0$, we can see that competitive ratio of the greedy policy goes unbounded. This is similar to the truncated Poisson arrival distribution where $\nu \to 0$. To address this question of competitive ratio going unbounded when $\mu$ decreases, in the next section, we consider a separate policy for Bernoulli distribution and derive a stronger lower bound compared to Corollary 1 for any policy under Bernoulli packet arrivals. We are unable to derive a similar lower bound for online policy for all distributions because of technical challenges.

## V. Low $\mu$ REGIME

To improve the competitive ratio performance in the $\mu \to 0$ regime, in this section we restrict our attention to the extreme Bernoulli arrival distribution

$$A_t = \begin{cases} B, & \text{w.p. } p \implies \text{buffer gets full}, \\ 0, & \text{w.p. } 1-p, \end{cases}$$

and $A_t$ are i.i.d. Thus, with the extreme Bernoulli distribution, on each arrival, the buffer gets full, and all the left over packets are dropped. Essentially, a new arrival is a renewal event for the system. The motivation behind studying the extreme Bernoulli arrival distribution follows from [14] that shows that for a similar rate maximization problem under random energy arrival distributions, the worst case input is the extreme Bernoulli arrival distribution. Thus, if the competitive ratio of any policy is bounded under the extreme Bernoulli arrival distribution, the same bound holds for all distributions [14]. It is with similar hope that we consider the extreme Bernoulli arrival distribution, where the server speed has to be chosen very delicately in order to satisfy the packet drop probability constraint, since on any arrival all the left over packets are dropped. At this time, however, we are unable to show that it is indeed the worst case distribution for (2), because of the hard packet drop probability constraint, unlike the problem considered in [14].

We first propose a $\lambda$-*fraction policy* that is similar to the policy first proposed in [14], and then improve upon the lower bound (Corollary 1) to show that the competitive ratio of the $\lambda$-fraction policy is constant for the extreme Bernoulli distribution even when $\mu \to 0$.

We will now derive an alternative expression for the packet drop probability for the extreme Bernoulli arrival distribution that will be easier to analyse.

**Lemma 1.** *For extreme Bernoulli arrival distribution,*

$$P_{Drop} = \frac{\mathbb{E}\left[\left(B - \sum_{t=1}^N g_t\right)^+\right]}{B},$$

*where expectation is over $N$, the inter-arrival time random variable for the extreme Bernoulli distribution, which is geometric with the same parameter $p$, and $(x)^+ = \max\{0, x\}$*

*Proof.* Recall that each new arrival is a renewal event with the extreme Bernoulli distribution, where all the left over packets are dropped at the new arrival. WLOG assume that two consecutive arrivals happen at slots $1$ and $N+1$, $N$ is the inter-arrival time random variable. Thus, the buffer state $b_1 = B$. Then from the Renewal Reward Theorem (RRT), we have

$$P_{\text{Drop}} = \frac{\mathbb{E}_N\left[\left(B - \sum_{t=1}^N g_t\right)^+\right]}{\mathbb{E}_N\left[\sum_{t=1}^N A_t\right]},$$

$$\overset{(a)}{=} \frac{\mathbb{E}_N\left[\left(B - \sum_{t=1}^N g_t\right)^+\right]}{\mathbb{E}_N[B]},$$

$$\overset{(b)}{=} \frac{\mathbb{E}_N\left[\left(B - \sum_{t=1}^N g_t\right)^+\right]}{B}.$$

where (a) follows since the two consecutive arrivals happen at slots $1$ and $N+1$, i.e., $A_1 = B$ and $A_t = 0$, $\forall\, 2 \leq t \leq N$, and (b) follows since $B$ is a constant. $\square$

### A. $\lambda$-Fraction Policy:

Let $b_t$ be the current state of the buffer, then the $\lambda$-fraction policy services $\lambda$-fraction of packets in the buffer, where $\lambda$ is determined by the arrival distribution and the packet drop probability constraint $\alpha$, i.e.,

$$g_t = \lambda b_t, t \geq 1. \tag{9}$$

A necessary condition for $\lambda$-fraction policy to be feasible is that $0 \leq \lambda \leq 1$ because $0 \leq g_t \leq b_t$.

For the $\lambda$-fraction policy, the number of packets in the buffer after $n$ time slots since the last arrival will be $B(1-\lambda)^n$, which get dropped when a new arrival happens with inter-arrival time $N = n$. Therefore, to compute the packet drop probability via Lemma 1, we calculate

$$\mathbb{E}\left[\left(B - \sum_{t=1}^N g_t\right)^+ | N = n\right] = B(1-\lambda)^n,$$

$$\mathbb{E}\left[\left(B - \sum_{t=1}^N g_t\right)^+\right] \overset{(a)}{=} \mathbb{E}_N\left[\mathbb{E}\left[\left(B - \sum_{t=1}^N g_t\right)^+ | N\right]\right],$$

$$= \Sigma_{n=1}^\infty B(1-\lambda)^n p(1-p)^{n-1},$$

$$\overset{(b)}{=} \frac{Bp(1-\lambda)}{\lambda + p - \lambda p}. \tag{10}$$

where (a) follows from taking the conditional expectation, and (b) follows by taking the expectation with respect to $N$. Thus, using Lemma 1 and (10), the packet drop probability is $P_{\text{Drop}} = \frac{p(1-\lambda)}{\lambda+p-\lambda p}$. Setting $P_{\text{Drop}} = \alpha$, we get a feasible

$$\lambda = \frac{p(1-\alpha)}{p + \alpha(1-p)}, \tag{11}$$

for the $\lambda$-fraction policy. Note, in (11), $\lambda \leq 1$. Next, we compute the cost incurred by the $\lambda$-fraction policy when $\lambda$ is given by (11).

### B. Cost Incurred by the $\lambda$-Fraction Policy

**Lemma 2.** *The cost incurred by the $\lambda$-fraction policy (where $\lambda$ satisfies (11)) is $\mathcal{J}_\lambda = \frac{\mu^2(1-\alpha)^2}{1-(1-p)(1-\alpha)^2}$ for the extreme Bernoulli distribution.*

*Proof.* Under the extreme Bernoulli, using the fact that each arrival is a renewal event, using RRT, the expected cost of the $\lambda$-fraction policy is

$$
\begin{aligned}
\mathcal{J}_\lambda &= \frac{\mathbb{E}\left[\sum_{t=1}^N g_t^2\right]}{\mathbb{E}[N]}, \\
&\stackrel{(a)}{=} \frac{B^2\lambda^2\,\mathbb{E}_N\left[1+(1-\lambda)^2+\cdots+(1-\lambda)^{2(N-1)}\right]}{\mathbb{E}[N]}, \\
&\stackrel{(b)}{=} \frac{pB^2\lambda^2\,\mathbb{E}_N[1-(1-\lambda)^{2N}]}{1-(1-\lambda)^2}, \\
&\stackrel{(c)}{=} \frac{pB^2\lambda^2}{1-(1-p)(1-\lambda)^2}, \\
&\stackrel{(d)}{=} \frac{B^2p^2(1-\alpha)^2}{1-(1-p)(1-\alpha)^2}, \\
&\stackrel{(e)}{=} \frac{\mu^2(1-\alpha)^2}{1-(1-p)(1-\alpha)^2}.
\end{aligned}
\tag{12}
$$

where (a) follows from (9) and expanding the summation, (b) follows from the fact that $\mathbb{E}[N] = \frac{1}{p}$, and summing up the geometric series in (b), (c) follows from taking expectation with respect to $N$, (d) follows from substituting (11), and (e) follows from the fact that $\mu = Bp$ from the extreme Bernoulli arrival process. $\square$

**Corollary 2.** *The competitive ratio of the $\lambda$-fraction policy is at most $\frac{1}{1-(1-p)(1-\alpha)^2}$.*

*Proof.* From Corollary 1, we have the lower bound $\mathcal{J} \geq \mu^2(1-\alpha)^2$. Thus, from Lemma 2, we get the following bound on the competitive ratio of the $\lambda$-fraction policy,

$$
\mathrm{CR}_\lambda \leq \frac{\mathcal{J}_\lambda}{\mathcal{J}^*} = \frac{1}{1-(1-p)(1-\alpha)^2}. \tag{13}
$$

$\square$

For the case of small $\alpha$ and small $p$, the competitive ratio of the $\lambda$-fraction policy can be approximated using binomial approximation as

$$
\mathrm{CR}_\lambda \approx \frac{1}{2\alpha + p}. \tag{14}
$$

From (14), as $\alpha \to 0$ and $p \to 0$, the competitive ratio of the $\lambda$-fraction policy also becomes unbounded similar to the greedy policy, even though the cost of $\lambda$-fraction policy is always lower than the greedy policy.

**Lemma 3.** *Given $p$ and $\alpha$, $\mathcal{J}_{greedy} \geq \mathcal{J}_\lambda$.*

*Proof.*

$$
\begin{aligned}
\mathcal{J}_{\text{greedy}} - \mathcal{J}_\lambda &\stackrel{(a)}{=} (1-\alpha)^2\left[B^2p - \frac{B^2p^2}{1-(1-p)(1-\alpha)^2}\right], \\
&= (1-\alpha)^2\,B^2p\left[\frac{(1-p)(1-\alpha)^2}{1-(1-p)(1-\alpha)^2}\right], \\
&\stackrel{(b)}{\geq} 0.
\end{aligned}
\tag{15}
$$

where (a) follows from (12) and (7) and the fact that for extreme Bernoulli distribution $\mu = Bp$ and $\mathrm{var}[A_t] = B^2p(1-p)$; (b) follows from the fact that $1-p \geq 0$ and $1-\alpha \geq 0$. $\square$

The reason for the large competitive ratio $\lambda$-fraction policy when $\alpha \to 0$ and $p \to 0$, is that the lower bound derived in Corollary 1 is weak, as we show next, via tightening the lower bound for extreme Bernoulli arrival distribution.

### C. Improved Lower Bound for small $\alpha$ and $p$

**Proposition 1.** *For $\alpha \ll 1, p \ll 1$, the cost of any online policy is $\Omega\left(\frac{B^2p^2}{p+\alpha}\right)$ for the extreme Bernoulli packet arrival process with parameter $p$.*

*Proof.* We consider two cases for the proof as follows: Case 1: $\alpha, p \ll 1$ and $\alpha \leq \frac{p}{2}$. From Lemma 1, because of packet drop probability constraint of $\alpha$, we have

$$
\begin{aligned}
\mathbb{E}\left[\left(B - \sum_{t=1}^N g_t\right)^+\right] &\leq B\alpha, \\
\sum_{i=1}^\infty P(N=i)\left(B - \sum_{t=1}^i g_t\right)^+ &\stackrel{(a)}{\leq} B\alpha, \\
P(N=1)\,(B-g_1) &\stackrel{(b)}{\leq} B\alpha, \\
g_1 &\stackrel{(c)}{\geq} B\left(1-\frac{\alpha}{p}\right), \\
&\stackrel{(d)}{\geq} \frac{B}{2}.
\end{aligned}
\tag{16}
$$

where (a) follows from expanding the expectation, (b) follows from taking only the first term of the series in the LHS, (c) is a rearrangement of terms, and (d) follows from the fact that $\frac{\alpha}{p} \leq \frac{1}{2}$. Thus, using the RRT again, the cost of any policy $\mathbf{g}$

$$
\begin{aligned}
\mathcal{J}_{\mathbf{g}} &= \frac{\mathbb{E}\left[\sum_{t=1}^N g_t^2\right]}{\mathbb{E}[N]}, \\
&\stackrel{(a)}{=} \frac{g_1^2 + P(N \geq 2)g_2^2 + P(N \geq 3)g_3^2 + \cdots}{\mathbb{E}[N]}, \\
&\stackrel{(b)}{\geq} \frac{g_1^2}{\mathbb{E}[N]}, \\
&\stackrel{(c)}{\geq} \frac{B^2p}{4}, \\
&\stackrel{(d)}{=} \Omega\left(\frac{B^2p^2}{p+\alpha}\right).
\end{aligned}
\tag{17}
$$

where (a) follows from expanding the expectation into summation, (b) follows from considering only the first term of the

series, (c) follows from the fact that $\mathbb{E}[N] = \frac{1}{p}$ and (16), and (d) follows from the fact that $p > \frac{p^2}{p+\alpha}$.

Case 2: $\alpha, p << 1$ and $\alpha > \frac{p}{2}$. Once again from Lemma 1, because of packet drop probability constraint of $\alpha$,

$$\mathbb{E}\left[\left(B - \sum_{t=1}^{N} g_t\right)^+\right] \le B\alpha,$$

$$\sum_{i=1}^{\infty} P(N = i)\left(B - \sum_{t=1}^{i} g_t\right)^+ \overset{(a)}{\le} B\alpha,$$

$$P(N \le k)\left(B - \sum_{t=1}^{k} g_t\right) \overset{(b)}{\le} B\alpha,$$

$$\sum_{t=1}^{k} g_t \overset{(c)}{\ge} B\left(1 - \frac{\alpha}{pk}\right),$$

$$\sum_{t=1}^{k} g_t \overset{(d)}{\ge} B/2. \tag{18}$$

where (a) is expansion of expectation into summation, (b) follows from considering only the first $k$ terms and taking the smallest drop $(B - \sum_{t=1}^{k} g_t)$, (c) follows from the approximation we make for $P(N \le k) = 1 - (1-p)^k \approx pk$ which is a reasonable approximation for the choice of $k = \left\lceil \frac{2\alpha}{p} \right\rceil$, and (d) follows since $\alpha > \frac{p}{2}$.

Thus, the cost for any policy $\mathbf{g}$

$$\mathcal{J}_{\mathbf{g}} \overset{(a)}{=} \frac{\mathbb{E}\left[\sum_{t=1}^{N} g_t^2\right]}{\mathbb{E}[N]},$$

$$\overset{(b)}{\ge} p \cdot P(N \ge k)[g_1^2 + g_2^2 + \cdots + g_k^2],$$

$$\overset{(c)}{\ge} p(1 - pk)[g_1^2 + g_2^2 + \cdots + g_k^2],$$

$$\overset{(d)}{\ge} p(1 - pk) \cdot \frac{B^2}{4k},$$

$$\overset{(e)}{=} \Omega\left(\frac{B^2 p^2}{\alpha + p}\right). \tag{19}$$

where (a) follows from the RRT, (b) follows since $\mathbb{E}[N] = 1/p$ and from taking terms corresponding to $N \ge k$ (we choose $k = \left\lceil \frac{2\alpha}{p} \right\rceil$ as taken above) and within each of them only considering the $\sum_{i=1}^{k} g_i^2$ term, (c) follows from the approximation $P(N \le k) = 1 - (1-p)^k \approx pk$ as described above, (d) follows from Lemma 4 since $\sum_{t=1}^{k} g_t \ge B/2$ from (18), finally, (e) follows from the choice of $k = \left\lceil \frac{2\alpha}{p} \right\rceil$ since the lower bound is valid for all $k$, and the fact that $\frac{p^2}{\alpha} > \frac{p^2}{p+\alpha}$.

From (17) and (19), $\mathcal{J}_{\mathbf{g}} = \Omega\left(\frac{B^2 p^2}{\alpha + p}\right)$, $\forall \alpha, p << 1$. $\square$

**Lemma 4.** *Minimum of $\sum_{i=1}^{k} g_i^2$ subject to the constraint that $\sum_{i=1}^{k} g_i \ge \frac{B}{2}$ is achieved for $g_i = \frac{B}{2k}, 1 \le i \le k$.*

*Proof.* The Lagrangian for the convex optimization program is $\mathcal{L} = \sum_{i=1}^{k} g_i^2 - \gamma\left(\sum_{i=1}^{k} g_i - \frac{B}{2}\right)$, which on differentiating and equating to zero, we get $\frac{\partial \mathcal{L}}{\partial g_i} = 2g_i - \gamma \overset{(b)}{=} 0$. Using the

constraint, $\sum_{i=1}^{k} g_i \ge \frac{B}{2}$, we get $\gamma = B/k$, and the optimal $g_i^{\star} = \frac{B}{2k}$. $\square$

Combining the lower bound Proposition 1 and the upper bound for the $\lambda$-fraction policy (Lemma 2) for the small $\alpha, p$ regime, we get the following result.

**Theorem 4.** *For any $\alpha << 1, p << 1$, the competitive ratio cost of the $\lambda$-fraction policy is constant under the extreme Bernoulli packet arrival process with parameter $p$.*

We can generalize the stronger lower bound that we derived in Proposition 1 that is valid only for extreme Bernoulli distribution to distributions that have significant mass to the right of $\frac{B}{2}$ as follows.

**Proposition 2.** *Let the arrival distribution be such that $P\left(A_t \ge \frac{B}{2} + \Theta(B)\right) = \Omega\left(\frac{\mu}{B}\right)$. Then, for $\alpha << 1$ and $\frac{\mu}{B} << 1$, the cost of any online algorithm is $\Omega\left(\frac{\mu^2}{\frac{\mu}{B} + \alpha}\right)$.*

The proof is omitted for lack of space, which is similar to proof of Proposition 1. Thus, comparing Proposition 2 and Lemma 2 it follows that the competitive ratio of the $\lambda$-fraction policy is constant for distributions where $P\left(A_t \ge \frac{B}{2} + \Theta(B)\right) = \Omega\left(\frac{\mu}{B}\right)$.

After deriving limited results for arbitrary buffer size $B$, we next turn to finite but large buffer capacities, for which we can derive universal results independent of the packet arrival distributions, when the packet drop probability is violated by a small margin.

## VI. CASE OF LARGE BUFFER CAPACITY

In this section, we consider that the buffer size is finite but large, and propose algorithms that are shown to achieve the optimal cost as the buffer size is increased. Towards that end, we will make some mild assumptions about the arrival process $A_t$. We assume that the asymptotic semi-invariant log moment generating function,

$$\Lambda_A(s) = \lim_{\tau \to \infty} \frac{1}{\tau} \log \mathbb{E}\left[e^{s \sum_{t=1}^{\tau} A_t}\right] \tag{20}$$

of $\{A_t\}$ exists for $s \in (-\infty, s_{\max})$ for some $s_{\max} > 0$ and the asymptotic variance $\sigma_A^2 = \lim_{\tau \to \infty} \frac{1}{\tau} \mathrm{var}\left(\sum_{t=1}^{\tau} A_t\right)$ of $\{A_t\}$ exists as assumed in [15].

**Remark 5.** *For the extreme Bernoulli distribution, asymptotic semi-invariant log moment generating function does not exist.*

Our policy is motivated by the energy management policy for an energy-harvesting node to maximise utility that was proposed in [15]. The policy that we propose is as follows: At each time slot, where $A_t$ packets arrive, we immediately drop $\alpha A_t$ number of those packets. The remaining $(1-\alpha)A_t$ number of packets are added to the buffer, i.e., buffer state evolves as $b_t = \min\{b_{t-1} + (1-\alpha)A_t - g_{t-1}, B\}$. Following which, $g_t$ number of packets are served at time slot $t$ according to the following rule:

$$g_t = \begin{cases} \min\{\mu(1-\alpha) - \delta, b_t\} & \text{if } b_t < B/2, \\ \mu(1-\alpha) + \delta & \text{if } b_t \ge B/2, \end{cases} \tag{21}$$

where $\delta = \beta \sigma_A^2 \frac{\log B}{B}$, and $\sigma_A^2$ is the variance of $A_t$.

**Theorem 6.** *For an arrival process $\{A_t\}_{t=1}^{\infty}$ for which the semi-invariant log moment generating function and the asymptotic variance exist, policy (21) ensures the following:*

$$\mathcal{J} \leq \mathcal{J}^* + \Theta\left(B^{-\beta}\right) + \Theta\left(\frac{(\log B)^2}{B^2}\right), and$$

$$P_{Drop} \leq \alpha + \Theta(B^{-\beta}),$$

*for any $\beta \geq 1$, where $\mathcal{J}^*$ is the lower bound as derived in Theorem 1.*

Theorem 6 states that if packet drop probability of $\alpha$ is violated by a margin of $\Theta(B^{-\beta})$, then the proposed policy can achieve a cost that is very close to the optimal. Note that it is tempting to make $\beta$ arbitrarily large so that the packet drop probability violation is negligible, however, higher value of $\beta$ leads to slower rate of convergence of the cost of the proposed policy to the optimal cost. Clearly, from Theorem 6 $\beta = 2$ gives the best performance in terms of cost.

*Proof.* Buffer overflow is said to occur if at time $t$ on the arrival of $A_t$ packets, $b_{t-1} + (1-\alpha)A_t > B$, where $b_{t-1}$ are the number of packets in the buffer in the time slot $t-1$. The buffer is said to be empty at time slot $t$ if $b_t = 0$. Let $P_{\text{overflow}}$ and $P_{\text{empty}}$ be the probability of buffer overflow and buffer is empty, respectively. It has been shown in [15] that for policy (21), $P_{\text{overflow}} = \Theta(B^{-\beta})$ and $P_{\text{empty}} = \Theta(B^{-\beta})$ using tools from large deviations theory and stochastic process limits. For a proof of these claims, refer [15].

A packet drop event occurs when either of the two events happen : (a) the $\alpha$ fraction of packets that get dropped immediately upon arrival by the definition of the policy, (b) the packets that get dropped because of an overflow in the buffer, i.e., $b_{t-1} + (1-\alpha)A_t > B$. Therefore, the packet drop probability can be upper bounded as:

$$P_{\text{Drop}} \overset{(a)}{\leq} \alpha + P_{\text{overflow}}\, \mathbb{E}[A_t],$$
$$\overset{(b)}{=} \alpha + \mu P_{\text{overflow}},$$
$$\overset{(c)}{=} \alpha + \Theta(B^{-\beta}). \qquad (22)$$

where (a) follows from the union bound over the two disjoint events that lead to packet drop, (b) follows from the fact that $\mathbb{E}[A_t] = \mu$, and (c) follows from $P_{\text{overflow}} = \Theta\left(B^{-\beta}\right)$ and the fact that $\mu$ is constant since the arrival distribution is fixed.

Let $\mathcal{J}^+$ be the cost incurred by our policy when $b_t \geq B/2$ by serving $\mu(1-\alpha) + \delta$ packets, i.e. $\mathcal{J}^+ = f_c(\mu(1-\alpha)+\delta)$. Similarly, when $b_t < B/2$, the policy serves at most $\mu(1-\alpha) - \delta$ packets, and let $\mathcal{J}^- = f_c(\mu(1-\alpha)-\delta)$. Since $f_c(.)$

is an analytic function, the Taylor series expansion of the cost function ($\mathcal{J}^+$ and $\mathcal{J}^-$) about $\mu(1-\alpha)$ can be written as

$$\mathcal{J}^+ = f_c\left(\mu\left(1-\alpha\right)\right) + f_c^{(1)}\left(\mu\left(1-\alpha\right)\right)\delta$$
$$+ f_c^{(2)}\left(\mu\left(1-\alpha\right)\right)\delta^2 + o(\delta^2), \qquad (23)$$

$$\mathcal{J}^- = f_c\left(\mu\left(1-\alpha\right)\right) - f_c^{(1)}\left(\mu\left(1-\alpha\right)\right)\delta$$
$$+ f_c^{(2)}\left(\mu\left(1-\alpha\right)\right)\delta^2 + o(\delta^2). \qquad (24)$$

where $f_c^{(i)}$ is the $i^{th}$ derivative of $f_c$. Let $\mathcal{J}$ denote the average cost incurred by policy (21). Define $\rho^+$ as the fraction of time that $b_t > \frac{B}{2}$ and $\rho^-$ as the fraction of time that $b_t \leq \frac{B}{2}$. Then the average cost function can be written as

$$\mathcal{J} \overset{(a)}{\leq} \rho^+ \mathcal{J}^+ + \left(\rho^- - P_{\text{empty}}\right)\mathcal{J}^-,$$
$$\overset{(b)}{=} f_c(\mu(1-\alpha)) + f_c^{(1)}(\mu(1-\alpha))(\rho^+\delta - \rho^-\delta)$$
$$+ \Theta\left(\frac{(\log B)^2}{B^2}\right), \qquad (25)$$

where (a) follows from the fact that at most $\mu(1-\alpha) - \delta$ packets are served when $b_t < B/2$ and $\mathcal{J}^-$ is overestimated, (b) follows from (23) and (24) and the fact that $P_{\text{empty}} = \Theta\left(B^{-\beta}\right)$, $\delta = \Theta\left(\frac{\log B}{B}\right)$, $\rho^+ + \rho^- = 1$, and $f_c^{(2)}(\mu(1-\alpha))$ is constant.

Using the fact that the number of serviced packets is equal to the number of packet arrivals except for the packets that get dropped, we get $\rho^+ \left(\mu(1-\alpha)+\delta\right) + \left(\rho^- - P_{\text{empty}}\right)\left(\mu(1-\alpha)-\delta\right) = \mu(1-\alpha)(1-P_{\text{overflow}})$. Rearranging this, and substituting expressions for $P_{\text{overflow}}$ and $P_{\text{empty}}$, and considering $\mu$ as a constant we get

$$\rho^+\delta - \rho^-\delta = \Theta\left(B^{-\beta}\right). \qquad (26)$$

Thus, the average cost of the policy can be written as

$$\mathcal{J} \overset{(a)}{\leq} f_c(\mu(1-\alpha)) + \Theta\left(B^{-\beta}\right) + \Theta\left(\frac{(\log B)^2}{B^2}\right),$$
$$\overset{(b)}{=} \mathcal{J}^* + \Theta\left(B^{-\beta}\right) + \Theta\left(\frac{(\log B)^2}{B^2}\right), \qquad (27)$$

where (a) follows from (25), (26) and the fact that $f_c^{(1)}(\mu(1-\alpha))$ is a constant, and (b) follows from Theorem 1. $\qquad \square$

## VII. SIMULATIONS

In this section, we present some simulations results on the competitive ratio of the proposed policies. We first consider the arbitrary buffer size $B$ case, and plot the competitive ratio of the $\lambda$-fraction policy in Fig. 1 with respect to the two lower bounds (Corollary 1 and Proposition 1) for extreme Bernoulli distribution with small $p$ and small packet drop probability constraint $\alpha$. We see that the competitive ratio of the $\lambda$-fraction policy is extremely large with the loose lower bound (Corollary 1), while it is $\sim 4$ with the improved lower bound Proposition 1. Next, we consider, large $B$ case, and plot in Fig. 2, the competitive ratio of the policy (21) together with packet drop probability violation (sub-figure on top right corner) for Poisson distribution with average packet arrival
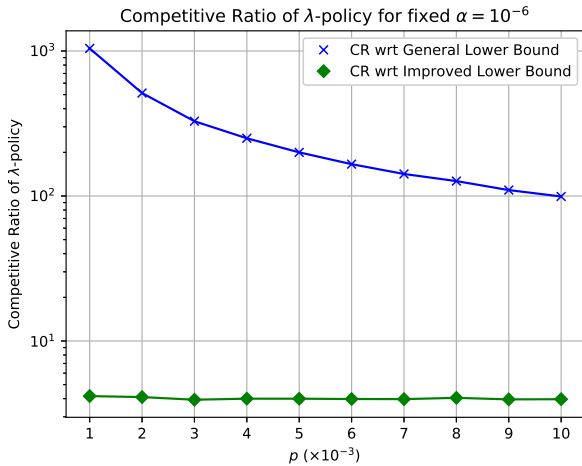
Fig. 1. Competitive ratio performance of the $\lambda$-fraction policy with extreme Bernoulli distribution as a function of $p$ with $B = 100$.
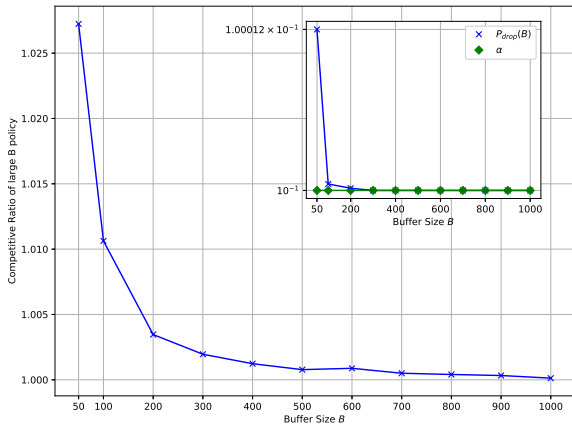


Fig. 2. Competitive Ratio Performance of the large $B$ policy for truncated Poisson distribution with parameter $\nu = 10$ and $\beta = 2$.

parameter $\nu = 10$ and $\beta = 2$. We see that as $B$ increases, the policy rapidly approaches the optimal performance, with negligible packet drop probability violation.

## VIII. Conclusions

In this paper, we considered a classical but challenging resource allocation problem, of finding optimal dynamic service speed of a server equipped with a finite buffer, under a hard constraint on the probability of dropping an arriving packet. This problem has been studied in literature before, however, with limited success and most useful results are available only in asymptotic regimes. In this paper, we considered a different approach and sought policies with provable guarantees on their gap from the optimal performance. We proposed a simple greedy policy and show that its competitive ratio is a small constant for large class of distributions. One limitation of the

greedy policy is that its competitive ratio is unbounded when the expected rate of arrival is very small. To address this issue, we then proposed a $\lambda$-fraction policy, that services a fixed fraction of the number of existing packets. We showed that the competitive ratio of the $\lambda$-fraction policy is bounded even when the expected rate of arrival is very small, when the arrival distribution is extreme Bernoulli. It is expected that the similar result will be applicable for other distributions, however, it remains open at this time. Finally, we considered the finite but large buffer case, where the tradeoff between packet drop probability violation versus the cost is studied. We showed that a policy inspired from [15], can achieve close to optimal performance when small packet drop probability violation is allowed, for all arrival distribution that satisfy a technical condition.

## References

[1] J. R. Perkins and R. Srikant, "The role of queue length information in congestion control and resource pricing," in *Decision and Control, 1999. Proceedings of the 38th IEEE Conference on*, vol. 3, pp. 2703–2708, IEEE, 1999.

[2] S. Bodas, S. Shakkottai, L. Ying, and R. Srikant, "Scheduling in multi-channel wireless networks: Rate function optimality in the small-buffer regime," in *ACM SIGMETRICS Performance Evaluation Review*, vol. 37, pp. 121–132, ACM, 2009.

[3] K. Jagannathan, E. Modiano, and L. Zheng, "On the role of queue length information in network control," *IEEE Transactions on Information Theory*, vol. 57, no. 9, pp. 5884–5896, 2011.

[4] J. M. George and J. M. Harrison, "Dynamic control of a queue with adjustable service rate," *Operations research*, vol. 49, no. 5, pp. 720–731, 2001.

[5] S. Stidham Jr and R. R. Weber, "Monotonic and insensitive optimal policies for control of queues with undiscounted costs," *Operations research*, vol. 37, no. 4, pp. 611–625, 1989.

[6] B. Ata and S. Shneorson, "Dynamic control of an m/m/1 service system with adjustable arrival and service rates," *Management Science*, vol. 52, no. 11, pp. 1778–1791, 2006.

[7] V. B. Sukumaran and U. Mukherji, "Tradeoff of average service cost and average delay for the state dependent m/m/1 queue," in *Communications (NCC), 2013 National Conference on*, pp. 1–5, IEEE, 2013.

[8] N. Bansal, T. Kimbrel, and K. Pruhs, "Speed scaling to manage energy and temperature," *Journal of the ACM (JACM)*, vol. 54, no. 1, p. 3, 2007.

[9] N. Bansal, H.-L. Chan, T.-W. Lam, and L.-K. Lee, "Scheduling for speed bounded processors," in *International Colloquium on Automata, Languages, and Programming*, pp. 409–420, Springer, 2008.

[10] N. Bansal, H.-L. Chan, K. Pruhs, and D. Katz, "Improved bounds for speed scaling in devices obeying the cube-root rule," *Automata, languages and programming*, pp. 144–155, 2009.

[11] H.-L. Chan, J. W.-T. Chan, T.-W. Lam, L.-K. Lee, K.-S. Mak, and P. W. Wong, "Optimizing throughput and energy in online deadline scheduling," *ACM Transactions on Algorithms (TALG)*, vol. 6, no. 1, p. 10, 2009.

[12] X. Han, T.-W. Lam, L.-K. Lee, I. K. To, and P. W. Wong, "Deadline scheduling and power management for speed bounded processors," *Theoretical Computer Science*, vol. 411, no. 40-42, pp. 3587–3600, 2010.

[13] S. Albers, F. Müller, and S. Schmelzer, "Speed scaling on parallel processors," *Algorithmica*, vol. 68, no. 2, pp. 404–425, 2014.

[14] D. Shaviv and A. Ö. Aydin, "Universally near optimal online power control for energy harvesting nodes," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 12, pp. 3620–3631, 2016.

[15] R. Srivastava and C. E. Koksal, "Basic performance limits and tradeoffs in energy-harvesting sensor nodes with finite data and energy storage," *IEEE/ACM Trans. Netw.*, vol. 21, pp. 1049–1062, Aug. 2013.

[16] A. Wierman, L. L. H. Andrew, and A. Tang, "Power-aware speed scaling in processor sharing systems: Optimality and robustness," *Perform. Eval.*, vol. 69, pp. 601–622, Dec. 2012.