# Energy-Optimal Scheduling with Variable Processing Speed: The Role of Task Size Variability

Jonatha ANSELMI and Bruno GAUJAL

{jonatha.anselmi,bruno.gaujal}@inria.fr

Univ. Grenoble Alpes, CNRS, Inria, Grenoble INP, LIG, 38000 Grenoble, France.

## Abstract

In this paper, we study the execution of a single task with an unknown size on a server with variable processing speed. Our goal is to analyze structural properties of the optimal energy consumption under the optimal speed profile that minimizes the expected energy consumption while meeting a hard deadline constraint. Specifically, we investigate how the task size probability distribution impacts the overall energy.

Under mild assumptions, our main result shows that the expected energy consumption induced by the optimal speed profile preserves the convex increasing order with respect to the task size distribution. Then, we leverage this property to derive simple bounds and conduct a worst-case analysis. In particular, we derive a simple, general formula for the energy gap induced by the 'best' and 'worst' task size distributions, expressed in terms of the support and expectation of the task size.

*Keywords:* Energy minimization, real time, optimal speed profile, stochastic comparison, increasing convex order

## 1. Introduction

Energy consumption is a critical factor in modern computing, impacting both environmental sustainability and operational costs. Large-scale GPU clusters, widely used for AI training, scientific computing, and cloud-based services, require substantial power to execute workloads efficiently. The high energy demands of these systems contribute significantly to carbon emissions, raising concerns about their environmental footprint. Moreover, energy costs represent a major expense for data centers, influencing overall profitability and resource allocation strategies. As computing workloads continue to grow in scale and complexity, optimizing energy efficiency is essential for reducing costs, mitigating environmental impact, and ensuring the long-term sustainability of large-scale computing infrastructures.

Optimizing the energy consumption of a task while meeting strict performance requirements has been a key research focus for decades, particularly in applications with real-time constraints, leading to an extensive body of literature. Most existing work on energy-efficient scheduling has focused on deterministic settings where task characteristics are fully known in advance or, equivalently, revealed upon task arrival. More specifically, when a task consists of a sequence of sub-tasks with *known* sizes, arrival times, and deadlines, the problem of determining the optimal speed profile has been studied in [10, 6, 5], where efficient algorithms for computing or approximating optimal speed profiles have been developed. This class of problems can be approached using the theory of Markov decision processes; see, for example, [4, 2].

However, in many practical scenarios, task sizes are uncertain, and only probabilistic information is available to the scheduler [7, 9]. Understanding the impact of this uncertainty on energy consumption is crucial for designing robust and efficient scheduling policies, and this is the approach taken in this paper. Within this framework, an elegant "closed-form" expression for the optimal speed profile has been proposed in [7]. This work considers the execution of a single task under the assumptions that (i) power consumption at speed $s$ follows $P(s) = ks^3$, where $k > 0$ is a constant, and (ii) the set of available speeds is $\mathbb{R}_+$. In contrast, we use more realistic assumptions, that is the power function $P(s)$ is an arbitrary convex and increasing function and the speed is bounded by $s_{\max}$. We give an implicit equation that determines the optimal speed function. Closed-form solutions are also obtained when the power function is $P(s) = ks^a$ with $a > 1$. In contrast, the analysis

in [1] considers *discontinuous* speed profiles, a scenario that naturally arises in multi-core or multi-GPU systems, where the computational platform provides only a finite set of speeds. Under mild assumptions on the structure of the power function, it is shown in [1] that an optimal speed profile corresponds to the unique solution of a strictly convex optimization problem. Additionally, an ad-hoc algorithm is developed to solve this optimization problem in $O(\log_2 n)$ steps, where $n$ is the number of available processing speeds.

Finally, let us mention that optimal speed profiles for executing a single task can be integrated into broader frameworks to develop heuristics for scheduling multiple jobs arriving at random times [3].

### 1.1. Contribution

In this work, we consider the expected energy consumption induced by an optimal speed profile as a starting point. Then, our objective is to investigate how the probability distribution of task sizes influences the total energy consumed.

Our main result establishes that the expected energy consumption induced by the optimal speed profile preserves the convex increasing order with respect to the task size distribution. This fundamental property allows us to show that the minimal energy consumption to execute a task (with $P(s) = ks^a$) can be written as $\frac{S^a}{D^{a-1}}$ where $S$ can be interpreted as the "energetic size" of the task and $D$ is the deadline. For example if the task is uniformly distributed between 0 and $W_{max}$, its energetic size is $S = \frac{a}{a+1} W_{max}$. Using the convex increasing order property, we can derive explicit upper and lower bounds for this energetic size: For any random task, $M \leq S \leq W_{max}$, where $M$ is the expected size of the task and $W_{max}$ its maximal size. This induces a simple general formula for the energy gap between the most and least favorable task size distributions that only depends on the support of the task size and its mean.

### 1.2. Organization

The rest of the paper is organized as follows. In Section 2, we introduce the system model and formalize the energy minimization problem. All of our theoretical results are presented in Section 3. Finally, Section 4 summarizes our findings and discusses potential research directions.

## 2. Energy Minimization Framework

The energy minimization framework under investigation is composed of a *server* (or processor), a *task* (or

job) and a *scheduler*. The server represents the processing unit responsible for executing the task, while the scheduler represents the control unit that decides at what speed the task must be processed at any point in time. The objective is to find a *speed profile* that minimizes the expected energy consumption induced by the processing of the task while ensuring that the task completes before a given deadline.

### 2.1. Server

The server is a processing unit with variable processing speed capabilities. It can model a Dynamic Voltage and Frequency Scaling (DVFS) processor or a more complicated system. The set of its available speeds is given by the interval $\mathcal{S} := [0, s_{max})$, where $s_{max} \in \mathbb{R}_+ \cup \{+\infty\}$.

The power dissipated by the server when running at speed $s \in \mathcal{S}$ is denoted by $P(s)$.

### 2.2. Task

A task is a sequence of operations of random size $W$ that arrives at time zero. The random variable $W$ represents the total amount of *work*. Only statistical information is known about the task size $W$. In particular, its probability distribution function is given by $F(w) := \mathbb{P}(W < w)$, which represents the probability that the task size does not exceed $w$. Let also $F^c(w) := \mathbb{P}(W > w)$.

We assume that the support of $W$ is $[W_{min}, W_{max}]$, where $W_{min} \geq 0$ and $W_{max} < \infty$ are the minimum and maximum task sizes, respectively. Finally, we impose the following *hard* real-time constraint: the task must be completed before a deadline $D$. We assume that $D$ satisfies the constraint $W_{max}/s_{max} \leq D$, which ensures that the server can complete the job if it operates at its maximum speed from the start.

### 2.3. Scheduler

The role of the scheduler is to control the speed at which the task is processed at any point in time. We assume that any change in processing speed is immediate and do not incur any additional energy cost; for instance, ramp-up and/or speed change effects are neglected. The information available to the scheduler is:

- The deadline $D$ (the task must be completed before $D$);

- The task size distribution (with support $[W_{min}, W_{max}]$);

- The set of available speeds $\mathcal{S}$;

- The dynamic power dissipation function $P$;

- The current amount of executed work: at time $t$, this is denoted by $w(t)$.

The task execution process clearly stops as soon as the task is complete (here, the scheduler puts the server at rest). Of course, the scheduler does not know when this is going to happen since it does not know the actual size of the task $W$.

Since the scheduler does not know the size of the task in advance, it is natural to take decisions as a function of $w$ (the work already executed) instead of as a function of time, as done in [7]. For this reason, in the following, we denote by $s(w)$ the speed used once exactly $w$ work units have been executed, $w \in [0, W_{\max}]$. Also, we let $Q(s)$, the dynamic power consumption *per unit of work* when operating at speed $s$, i.e.,

$$Q(s) := P(s)\frac{\mathrm{d}t}{\mathrm{d}w} \qquad (1)$$

where derivatives are always intended as right derivatives.

In the remainder of the paper, we require the following assumption, which is satisfied for most electronic circuits.

**Assumption 1.** *The dynamic power consumption function $Q$ is convex and increasing.*

### 2.4. Speed Profiles and Energy Consumption

Let $s : [0, W_{\max}] \to \mathcal{S}$ denote a speed profile. The mean energy consumption under speed profile $s$ is defined by

$$\mathcal{E}(s) := \mathbb{E}\left[\int_0^W Q(s(u))\,\mathrm{d}u\right]. \qquad (2)$$

By conditioning on the size ($W = w$) of the task and changing the order of integration, we obtain

$$\mathcal{E}(s) = \int_0^{W_{\max}} \left(\int_0^x Q(s(w))\,\mathrm{d}w\right)\mathrm{d}F(x)$$
$$= \int_0^{W_{\max}} Q(s(w))\left(\int_w^{W_{\max}} \mathrm{d}F(x)\right)\mathrm{d}w$$
$$= \int_0^{W_{\max}} Q(s(w))\,F^c(w)\,\mathrm{d}w. \qquad (3)$$

Within the foregoing assumptions, we denote by $s^*$ the speed profile that minimizes $\mathcal{E}(s)$ subject to the constraint

$$\int_0^{W_{\max}} \frac{\mathrm{d}w}{s(w)} \leq D, \qquad (4)$$

which states that the job has to be completed before the deadline $D$.

We recall that in the literature, the structure of $s^*$ has been studied under certain assumptions. If $s_{\max} = +\infty$ and the power function has the form $Q(s) = s^\alpha$, it is shown in [7] that

$$s^*(w) = K\,F^c(w)^{-\frac{1}{\alpha+1}} \qquad (5)$$

where $K$ is a (normalizing) constant such that the deadline constraint (4) is satisfied with equality, i.e., $\int_0^{W_{\max}} \frac{\mathrm{d}w}{s^*(w)} = D$.

The interpretation in Formula (5) shows that the optimal speed profile is increasing in the amount of completed work $w$ (and thus also time). This is intuitive because if the deadline approaches and the job has not completed, it is natural to accelerate.

If $\mathcal{S}$ is composed of a finite number of speeds and $Q(s)$ is increasing, it is shown in [1] that $s^*$ is the solution of a strictly convex optimization problem with $|\mathcal{S}| - 1$ decision variables.

## 3. Main Results

In this section, we present our main results. Specifically, (i) we generalize Formula (5) to a more general setting, (ii) we show our stochastic comparison results, and finally, (iii) we develop an efficiency analysis to evaluate the impact of the task size distribution.

### 3.1. Optimal Speed Profile

While (5) provides an elegant solution to the energy minimization problem under deadline constraint, it has an important drawback. Since $F^c(w) \to 0$ as $w \uparrow W_{\max}$, Formula (5) implies that the optimal speed profile relies on arbitrarily large speeds. In other words, it does not work when the maximal server speed is bounded.

The following result provides an extension of (5) that works even when $s_{\max} < \infty$.

While (5) was originally proven by using Jensen's inequality, our approach here is more general and relies on Pontryagin maximum principle. This also lets us deal with the case where $Q$ is "general".

**Theorem 1.** *The optimal speed profile $s^*$ satisfies*

$$s^*(w)^2 Q'(s^*(w)) = \frac{\lambda^*}{F^c(w)} \qquad (6)$$

*if $s^*(w) \in [0, s_{\max}]$ and $s^*(w) = s_{\max}$ otherwise.*

3

*Proof.* We start by characterizing $s^*$ using Pontryagin maximum principle as follows. With respect to the infinite dimensional optimization problem $\min_s \mathcal{E}(s)$ subject to (4) (with $\mathcal{E}(s)$ given by the expression in (3)), we define the Hamiltonian

$$H(\theta(w), s(w), \lambda(w), w) := F^c(w)Q(s(w)) - \frac{\lambda(w)}{s(w)},$$

where $\lambda(w)$ is the Lagrangian multiplier. The Pontryagin maximum principle says that the optimal solution $(s^*, \lambda^*)$ satisfies the following conditions:

1. $H(\tau^*(w), s^*(w), \lambda^*(w), w) \leq H(\theta(w), s(w), \lambda(w), w)$ for all $w \in [0, W_{\max}]$.
2. $\frac{d\lambda^*(w)}{dw} = \lambda^*(w)\frac{\partial(1/s^*(w))}{\partial \tau} + \frac{\partial(F^c(w)Q(s^*(w)))}{\partial \tau}$.

The second condition implies that $\lambda^*(w)$ is a constant denoted $\lambda^*$, and the first condition gives the optimal speed under an implicit form. Differentiating with respect to $s$, we obtain that the optimal speed profile satisfies (9). □

Note that the equation (9) has a unique non-negative solution because $Q$ is convex increasing. In the classical case, where

$$Q(s) = ks^\alpha, \quad \alpha > 1, \tag{7}$$

this gives

$$s^*(w) = K^{-1} F^c(w)^{-1/(\alpha+1)} \wedge s_{\max},$$

where $K$ is a constant such that $\int_0^{W_{\max}} \frac{dw}{s^*} = D$, which ensures the validity of the deadline constraint.

Figure 1 displays the shape of the optimal solution when $\alpha = 2$ and the task size is uniformly distributed over $[0, W_{\max}]$, i.e., $F^c(w) = 1 - \frac{w}{W_{\max}}$. In this case,

$$\int_0^{W_{\max}} \left( K \left( 1 - \frac{w}{W_{\max}} \right)^{1/3} \vee \frac{1}{s_{\max}} \right) dw = D$$

.

If $s_{\max} = +\infty$, then the resulting speed profile corresponds to the solution given in [7], and one can indeed see that $s^*(w) \to \infty$ as $w \to W_{\max}$.

### 3.2. Comparing Task Size Distributions

We are interested in investigating how changes on the task size $W$ with probability distribution $F$ impact the *optimal* mean energy consumption function, which we define by

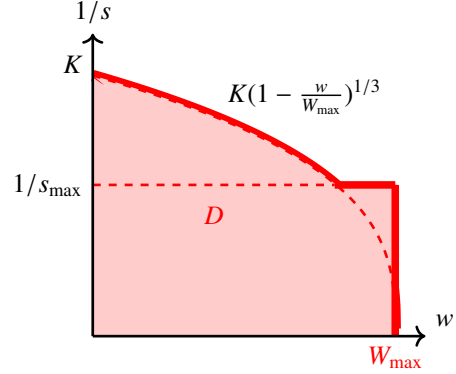$$\mathcal{E}_W^* := \mathcal{E}(s_W^*), \tag{8}$$



Figure 1: Optimal speed profile when the job size follows a uniform distribution ($F^c(w) = 1 - \frac{w}{W_{\max}}$): $1/s^*(w) = K(1 - \frac{w}{W_{\max}})^{1/3} \vee 1/s_{\max}$. The starting point $K$ is chosen such that the area of the shaded zone is $D$.

where $s_W^*$ is the optimal speed profile to execute a task of size $W$. In the following result, we use the increasing convex order (denoted by $\leq_{icx}$) [8]. Given two random variables $W_1$ and $W_2$, we recall that $W_1 \leq_{icx} W_2$ if $\mathbb{E}[h(W_1)] \leq \mathbb{E}[h(W_2)]$ for all convex increasing function $h : \mathbb{R} \to \mathbb{R}$, provided the expectations exist.

When comparing different task sizes, the next result shows that the optimal mean energy consumption is monotone w.r.t. the *icx* order.

**Theorem 2.** *Consider two tasks, let $W_1$ and $W_2$ denote their random sizes, and let $\mathcal{E}_{W_1}^*$ and $\mathcal{E}_{W_2}^*$ their optimal mean energy consumptions, respectively. If $W_1 \leq_{icx} W_2$, then $\mathcal{E}_{W_1}^* \leq \mathcal{E}_{W_2}^*$.*

*Proof.* We have shown in Theorem 1 that if $s^*$ is an optimal speed profile, then

$$g(s^*(w)) = \frac{\lambda^*}{F^c(w)}, \text{ if } s^*(w) \in [0, s_{\max}], \tag{9}$$

where we have defined the function $g : \mathbb{R} \to \mathbb{R}$ by $g(x) = x^2 Q'(x)$. The equation (9) has a unique non-negative solution because $Q$ is convex increasing and this implies that $g$ is increasing. Let $g^{-1}$ denote the inverse function of $g$. Note that

$$s^*(w) = g^{-1}\left(\frac{\lambda^*}{F^c(w)}\right) \wedge s_{\max}$$

and given that $g^{-1}$ is increasing (because $g$ is increasing) and $F^c(w)$ is non-decreasing, we obtain that $s^*(w)$ is non-decreasing in $w$.

Now, by definition of (2), we can interpret the optimal mean energy consumption as given by $\mathcal{E}^*(s) = E[h_W(W)]$ where $h_W$ is convex increasing and depends

4

on the probability distribution of $W$ (via (9)). To see this, we notice that we have just shown that $s^*$ is increasing, which means that $Q(s^*(u))$ is also increasing and positive, which in turn proves that $h(x) = \int_0^x Q(s^*(u))\, du$ is convex increasing.

The proof is thus concluded by the following chain inequalities:

$$\mathcal{E}^*_{W_1} = \mathbb{E}[h_{W_1}(W_1)] \le \mathbb{E}[h_{W_2}(W_1)] \le \mathbb{E}[h_{W_2}(W_2)] = \mathcal{E}^*_{W_2}.$$

The first inequality holds because $h_{W_1}$ is the function induced by the speed profile that is optimal for $W_1$; note that $\mathbb{E}[h_{W_2}(W_1)]$ is interpreted as the energy consumption induced by a task with size $W_1$ under the speed profile that is optimal for a task of size $W_2$. In turn, the second inequality holds because $W_1 \le_{icx} W_2$ (by hypothesis) and because we have shown that $h_{W_2}$ is convex increasing. $\qquad\square$

### 3.3. Efficiency Analysis

We now rely on the previous result to find the probability distribution that maximizes /minimizes $\mathcal{E}^*$. Towards this purpose, let $\mathcal{T}(W_{\min}, M, W_{\max})$ denote the set of task sizes $W$ with support $[W_{\min}, W_{\max}]$ and mean $M$.

In the set $\mathcal{T}(W_{\min}, M, W_{\max})$, we single out the task size $W_2$ with only two atoms at extreme points $W_{\min}$ and $W_{\max}$]. Its distribution is

$$F_2(w) = \frac{W_{\max} - M}{W_{\max} - W_{\min}}, \quad \forall w \in (W_{\min}, W_{\max}].$$

We also focus on the deterministic task size $W_1 \in= \mathcal{T}(W_{\min}, M, W_{\max}) : W_1 = M$ with probability 1.

The following result, a corollary of Theorem 2, says that among all task sizes in $\mathcal{T}(W_{\min}, M, W_{\max})$, the best one (minimizing the convex increasing order) is $W_1$ and the bast one is $W_2$

**Corollary 1.** Let $\mathcal{E}^*_1$ and $\mathcal{E}^*_2$ be the optimal mean energy consumption when the task size is $W_1$ and $W_2$, respectively. Then, for all $W$ in $\mathcal{T}(W_{\min}, M, W_{\max})$,

$$\mathcal{E}^*_1 \le \mathcal{E}^*_W \le \mathcal{E}^*_2.$$

*Proof.* Let $\le_{cx}$ denote the convex order (by definition, $X \le_{cx} Y$ if $\mathbb{E}(h(X)) \le \mathbb{E}(h(Y))$ for all $h$ convex). Theorem 3.A.44 in [8] says that for two random variables $X$ and $Y$ with equal means and distribution functions $F$ and $G$ respectively on $[a, b]$, $X \le_{cx} Y$ if there exists a critical value $h \in [a, b]$ such that function $F^c(w) \le G^c(w)$ for all $w \le h$ and $F^c(w) > G^c(w)$ otherwise. Consider any $W$ in $\mathcal{T}(W_{\min}, M, W_{\max})$ with distribution function $F$. Let us first show that the critical value $h_1$ in the comparison between $F$ and $F_1$ (the distribution function of

the deterministic size $W_1 = M$), is $h_1 = M$. Indeed, by definition, $F_c(w) \le 1 = F^c_1(w)$, for all $w \le M$ and $F_c(w) > 0 = F^c_1(w)$, for all $w > M$ because its mean cannot be smaller than $M$. As for the comparison between $W$ with $W_2$, remark that the distribution function $F^c$ of $W$ starts with value 1 in $W_{\min}$, is decreasing and ends at a value strictly smaller than $F_2(W_{\min})$ (otherwise its mean would be larger than $M$). This implies the existence of a critical value $h_2$.

Therefore,

$$W_1 \le_{cx} W \le_{cx} W_2.$$

By definition of the $\le_{cx}$ and $\le_{icx}$ orders, this implies

$$W_1 \le_{icx} W \le_{icx} W_2,$$

and the proof is concluded by applying Theorem 2. $\quad\square$

Our objective is now to examine the impact of the distribution function of the random task on the magnitude of $\mathcal{E}^*$. In view of Corollary 1, we define the efficiency ratio

$$\phi(W) := \frac{\mathcal{E}^*_W}{\mathcal{E}^*_1} \ge 1, \quad W \in \mathcal{T}(W_{\min}, M, W_{\max})$$

and leverage the previous two corollaries to establish bounds on $\phi(F)$. Before proceeding, we introduce the following intermediate result that characterizes the optimal speed profile when the task size distribution is $F_2$.

**Proposition 1.** *Assume that the task size distribution is $F_2$. If $W_{\min} > 0$, then*

$$s^*(w) = \begin{cases} W_{\min} \left( D - \frac{W_{\max} - W_{\min}}{s_M} \right)^{-1} & w \le W_{\min} \\ s_M & W_{\min} < w \le W_{\max} \end{cases}$$

*where $s_M$ minimizes over $\mathcal{T}$ the function*

$$s \mapsto W_{\min}\, Q\left( W_{\min} \left( D - \frac{W_{\max} - W_{\min}}{s} \right)^{-1} \right) + (M - W_{\min})\, Q(s). \tag{10}$$

*If $W_{\min} = 0$, then $s^*(w) = W_{\max}/D$.*

*Proof.* Let $p := \frac{W_{\max} - M}{W_{\max} - W_{\min}}$. Substituting $F_2$ in (2), we obtain

$$\mathcal{E}(s) = p \int_0^{W_{\min}} Q(s(u))\, du + (1 - p) \int_0^{W_{\max}} Q(s(u))\, du$$

$$= \int_0^{W_{\min}} Q(s(u))\, du + (1 - p) \int_{W_{\min}}^{W_{\max}} Q(s(u))\, du \tag{11}$$

and the optimal speed profile $s^*$ is thus given by minimization of (11) over $s$ subject to (4). Note that the

integrals in (11) are defined over disjoint sets, so minimizing $\mathcal{E}(s)$ over $s : [0, W_{max}] \to \mathbb{R}_+$ is equivalent to minimizing

$$E(x, y) := \int_0^{W_{min}} Q(x(u)) \, du + (1 - p) \int_{W_{min}}^{W_{max}} Q(y(u)) \, du$$

over

$$\mathcal{D} := \left\{ (x : [0, W_{min}] \to \mathbb{R}_+, y : [W_{min}, W_{max}] \to \mathbb{R}_+) : \int_0^{W_{min}} \frac{dw}{x(w)} + \int_{W_{min}}^{W_{max}} \frac{dw}{y(w)} \leq D \right\}.$$

Since $Q$ is convex, the structure of this optimization implies that $x$ and $y$ are constants. So

$$\mathcal{E}_2^* = \min_{(x,y) \in \mathcal{D}} E(x, y)$$
$$= \min_{s_m, s_M} W_{min} Q(s_m) + (1 - p)(W_{max} - W_{min}) Q(s_M)$$
$$\text{s.t.: } \frac{W_{min}}{s_m} + \frac{W_{max} - W_{min}}{s_M} \leq D \qquad (12)$$
$$s_m \geq 0, s_M \geq 0$$

Since $Q$ is increasing, the inequality in (12) can be replaced by equality. This gives

$$s_m = W_{min} \left( D - \frac{W_{max} - W_{min}}{s_M} \right)^{-1}$$

and the previous optimization becomes in dimension one with the objective function given by (10). Finally, when $W_{min} \to 0$, $s_m \to 0$ and the deadline constraint gives $s_M = W_{max}/D$. □

We can now present our results on the worst case efficiency ratio $\phi(F)$.

**Theorem 3.** *Assume that $W_{min} = 0$. Then,*

$$\sup_{W \in \mathcal{T}(W_{min}, M, W_{max})} \phi(W) = \frac{\mathcal{E}_2^*}{\mathcal{E}_1^*} = \frac{Q\left(\frac{W_{max}}{D}\right)}{Q\left(\frac{M}{D}\right)} \qquad (13)$$

*Proof.* The first equality is trivial given Corollary 1. If the task size is $M$, then it is optimal to run at the constant speed that completes the task exactly at $D$. This is given by $s^* = M/D$, and thus

$$\mathcal{E}_1^* = M \, Q\left(\frac{M}{D}\right). \qquad (14)$$

If the task size distribution is $F_2$ and $W_{min} = 0$, then $s^* = \frac{W_{max}}{D}$ (by Proposition 1), and this gives $\mathcal{E}_2^* = M Q\left(\frac{W_{max}}{D}\right)$. Dividing by $\mathcal{E}_1^*$, we obtain (15). □

This characterizes in simple terms the efficiency loss in energy consumption caused by perturbations in the task size probability distribution function. Depending on the support, it is clear that this loss can be arbitrarily large. If the power function follows the form (7) then the previous result yields (provided that $W_{min} = 0$)

$$\sup_{W \in \mathcal{T}(W_{min}, M, W_{max})} \phi(W) = \left(\frac{W_{max}}{M}\right)^{\alpha},$$

which indicates that efficiency is primarily influenced by the ratio between the maximum task size and its expected value.

We now consider the specific case where $Q(s) = k s^{\alpha}$ but $W_{min}$ is not necessarily zero. In this case, the next result provides an explicit expression for the efficiency ratio.

**Theorem 4.** *Assume that (7) holds. Then,*

$$\sup_{W \in \mathcal{T}(W_{min}, M, W_{max})} \phi(W) = \left(\frac{W_{min} + W_\star}{M}\right)^{\alpha+1} \qquad (15)$$

*where $W_\star := (M - W_{min})^{\frac{1}{\alpha+1}} (W_{max} - W_{min})^{\frac{\alpha}{\alpha+1}}$.*

*Proof.* By Proposition 1, the optimal speed profile induced by $W_2$ is composed by two speeds only. Let us refer to the speed to be used in $[0, W_{min}]$ as $\frac{W_{max}}{d}$. Then, the speed to be used in $[W_{min}, W_{max}]$ needs to be $(W_{max} - W_{min})/(D - d)$, as otherwise the deadline constraint would not be respected. Let $p := \frac{W_{max} - M}{W_{max} - W_{min}}$. Therefore,

$$\mathcal{E}_2^* = \min_{0 \leq d \leq D} \left( \int_0^{W_{min}} Q\left(\frac{W_{min}}{d}\right) dw \right.$$
$$\left. + (1 - p) \int_{W_{min}}^{W_{max}} Q\left(\frac{W_{max} - W_{min}}{D - d}\right) dw \right)$$
$$= \min_{0 \leq d \leq D} W_{min} \left(\frac{W_{min}}{d}\right)^{\alpha} + (M - W_{min}) \left(\frac{W_{max} - W_{min}}{D - d}\right)^{\alpha}$$

Differentiating, we obtain that the minimizing $d$ satisfies

$$\frac{W_{min}^{\alpha+1}}{d^{\alpha+1}} = (M - W_{min}) \frac{(W_{max} - W_{min})^{\alpha}}{(D - d)^{\alpha+1}}.$$

After some algebra, this gives

$$d = D \frac{W_{min}}{W_{min} + \underbrace{(M - W_{min})^{\frac{1}{\alpha+1}} (W_{max} - W_{min})^{\frac{\alpha}{\alpha+1}}}_{:= W_\star}}$$

and, substituting back in $\mathcal{E}_2^*$,

$$\mathcal{E}_2^* = \frac{W_{min}^{\alpha+1}}{d^{\alpha}} + (M - W_{min}) \frac{(W_{max} - W_{min})^{\alpha}}{(D - d)^{\alpha}}$$

6

$$= \frac{(W_{\min} + W_\star)^{\alpha+1}}{D^\alpha}.$$

Using Corollary 1 and $\mathcal{E}_1^* = \frac{M^{\alpha+1}}{D^\alpha}$ (see (14)), we obtain

$$\sup_{W \in \mathcal{T}(W_{\min}, M, W_{\max})} \phi(W) = \frac{\mathcal{E}_2^*}{\mathcal{E}_1^*} = \frac{(W_{\min} + W_\star)^{\alpha+1}}{M^{\alpha+1}}$$

as desired. $\qquad\square$

## 4. Conclusion

This work analyzed the impact of task size uncertainty on energy-optimal scheduling in systems with variable processing speed. We established fundamental properties of the optimal speed profile and leveraged them to derive energy bounds and a worst-case analysis, showing that the energy gap between the best and worst task size distributions essentially depends on the ratio between its maximal size and its average size, when the dynamic power follows the classical form.

Our findings contribute to the broader understanding of energy-efficient scheduling under uncertainty and highlight key structural properties that can inform practical scheduling policies. In fact, when scheduling multiple tasks, our results may be leveraged to determine the optimal processing *order* based on their size distributions. We leave this issue as future research, together with extensions to cases where the set of available speeds is finite.

## References

[1] J. Anselmi and B. Gaujal. Energy optimal activation of processors for the execution of a single task with unknown size. In *2022 30th International Symposium on Modeling, Analysis, and Simulation of Computer and Telecommunication Systems (MASCOTS)*, pages 65–72, 2022.

[2] J. Anselmi, B. Gaujal, and L. Sébastien Rebuffi. Optimal Speed Profile of a DVFS Processor under Soft Deadlines. *Performance Evaluation*, 152, Dec. 2021.

[3] F. Filippini, J. Anselmi, D. Ardagna, and B. Gaujal. A stochastic approach for scheduling AI training jobs in gpu-based systems. *IEEE Trans. Cloud Comput.*, 12(1):53–69, 2024.

[4] B. Gaujal, A. Girault, and S. Plassart. Dynamic Speed Scaling Minimizing Expected Energy Consumption for Real-Time Tasks. *Journal of Scheduling*, pages 1–25, July 2020.

[5] B. Gaujal, A. Girault, and S. Plassart. A Pseudo-Linear Time Algorithm for the Optimal Discrete Speed Minimizing Energy Consumption. *Discrete Event Dynamic Systems*, 31:163–184, 2021.

[6] M. Li, F. F. Yao, and H. Yuan. An $O(n^2)$ algorithm for computing optimal continuous voltage schedules. In *TAMC'17*, volume 10185 of *LNCS*, pages 389–400, Bern, Switzerland, Apr. 2017.

[7] J. R. Lorch and A. J. Smith. Improving dynamic voltage scaling algorithms with PACE. In *Joint International Conference on Measurements and Modeling of Computer Systems, SIGMETRICS'01*, pages 50–61, Cambridge (MA), USA, June 2001. ACM.

[8] M. Shaked and J. G. Shanthikumar. *Stochastic orders and their applications*. Academic Pr, 1994.

[9] R. Xu, D. Mossé, and R. Melhem. Minimizing expected energy in real-time embedded systems. In *Proceedings of the 5th ACM International Conference on Embedded Software*, EMSOFT '05, page 251–254, New York, NY, USA, 2005. Association for Computing Machinery.

[10] F. Yao, A. Demers, and S. Shenker. A scheduling model for reduced CPU energy. In *Proceedings of IEEE Annual Foundations of Computer Science*, pages 374–382, 1995.