# A Theoretical Foundation for Scheduling and Designing Heterogeneous Processors for Interactive Applications

Shaolei Ren[1] and Yuxiong He[2] and Kathryn S. McKinley[2]

[1] Florida International University, Miami, FL
[2] Microsoft Research, Redmond, WA

**Abstract**

To improve performance and meet power constraints, vendors are introducing heterogeneous multicores that combine high performance and low power cores. However, choosing which cores and scheduling applications on them remain open problems. This paper presents a scheduling algorithm that provably minimizes energy on heterogeneous multicores and meets latency constraints for interactive applications, such as search, recommendations, advertisements, and games. Because interactive applications must respond quickly to satisfy users, they impose multiple constraints, including average, tail, *and* maximum latency. We introduce SEM (Slow-to-fast, Energy optimization for Multiple constraints), which minimizes energy by choosing core speeds and how long to execute jobs on each core. We prove SEM minimizes energy without *a priori* knowledge of job service demand, satisfies multiple latency constraints simultaneously, and only migrates jobs from slower to faster cores. We address practical concerns of migration overhead and congestion. We prove optimizing energy for *average* latency requires homogeneous cores, whereas optimizing energy for *tail* and *deadline* constraints requires heterogeneous cores. For interactive applications, we create a formal foundation for scheduling and selecting cores in heterogeneous systems.

## 1 Introduction

Power constraints are forcing computer architects to turn to heterogeneous multicore hardware to improve performance. For instance, smartphones are shipping with Qualcomm's Snapdragon and ARM's Cortex-A15 [17], which include high performance and low power cores with the same instruction set, called big/little and Asymmetric Multicore Processors (AMP). Design principles for selecting cores in heterogeneous system and scheduling algorithms that optimize their energy consumption, however, remain open problems. This paper presents a scheduling algorithm that provably minimizes energy on heterogeneous processors serving interactive applications. We prove and establish scheduling insights and design principles with practical implications for heterogeneous core selection.

Interactive applications are latency-sensitive. Examples include serving web pages, games, search, advertising, recommendations, and mobile applications. Since interactive applications must be responsive to attract and please users, they

must meet *latency* requirements. Furthermore, they must be energy efficient. In the data center, power is an increasingly higher fraction of total costs [21,26,40]. A 1% energy saving may translate to millions of dollars. On mobile, energy efficiency translates directly into longer battery life and happier users.

Prior schedulers that optimize for energy efficiency and heterogeneity have major limitations. (1) They must predict demand for each request, scheduling high demand jobs to high performance fast cores and other jobs to low power slow cores [1,11,13,39,42,44]. Unfortunately, the service demand of individual requests in interactive applications is usually unknown and difficult to predict [25]. (2) For unknown service demand, prior work only optimizes for a *single simple* latency constraint [25,38,41], such as average latency or maximum latency, and is inadequate for two reasons. First, many applications strive for consistency by reducing tail latency (e.g., 95th- and 99th-percentile) or variance [15,19], which average and maximum latency do not model. Second, some applications *require* a combination of low average, tail, and worst-case latencies [14,18]. For example, search, finance applications, ads, and commerce have customer requirements and expectations for average and tail latency [14,15,19,26,40].

This paper shows how to optimize energy efficiency of interactive workloads subject to multiple latency constraints by exploiting heterogeneous multicores, addressing the aforementioned challenges as follows.

*Unknown service demand.* Instead of predicting individual job demand, we exploit the service demand *distribution* measured online or offline, which changes slowly over time [25,29]. We schedule incoming jobs to appropriate cores without knowing their individual service demands.

*Multiple latency constraints.* The scheduling literature typically optimizes for average or maximum latency only. To generalize and combine latency constraints, we use $L_p$ norms [2–5,20,28,34,43]. The $L_p$ norms encapsulate maximum latency ($p \to \infty$) and average latency ($p = 1$) as special cases. Optimizing for larger values of $p$ places more emphasis on the latency of longer jobs. Appropriate values of $p$ effectively mitigate unfairness and extreme outliers for long jobs [2,5]. Optimizing the $L_1$ and $L_2$ norms together reduce latency variance, which makes latency more *predictable* and improves user experience [35].

This paper presents an optimal algorithm that minimizes energy on heterogeneous processors given a demand distribution and latency constraints. We quantitatively characterize the optimal schedule and the ratio of fast to slow core speeds in a heterogeneous system. We present an *optimal* scheduling algorithm, called SEM (Slow-to-fast, Energy optimization for Multiple constraints). Given a service demand distribution, SEM schedules interactive jobs on heterogeneous multicore processors to minimize energy consumption while simultaneously satisfying multiple $L_p$ norm latency constraints.

We show *an optimal schedule migrates jobs from slower to faster cores.* Ideally, we want to schedule high demand (long) jobs on fast cores to meet latency requirements and short jobs on slow cores to save energy *without* a prior knowledge of service demand. SEM exploits this observation by scheduling short jobs

on energy efficient slow cores where they complete with high probability and then migrating long jobs to fast cores to meet the latency constraints.

We show *more heterogeneity is desirable for higher p*, where $p$ is the $L_p$ norm moment and the heterogeneity degree is the ratio of the fastest to slowest core speed. Given a single average latency constraint ($p = 1$), the energy optimal schedule requires a homogeneous processor. For all other latency constraints ($p > 1$) and multiple constraints, the optimal schedule requires heterogeneous processors.

We show *bounds on the ratio of the fastest and slowest core speeds for an optimal heterogeneous processor*. The result indicates that the more heterogeneous workload is and/or the less power additional core performance consumes, the more heterogeneous the hardware needs to be. Our result provides a formal and quantitative guide for selecting core speeds while designing heterogeneous processors. For practical choices of $p$ and measured service load distributions, the ratio ranges from two to eight. Systems with this degree of heterogeneity are thus quite practical to assemble from current server, client, and mobile cores.

This technical report supplements a conference publication at DISC 2014 [30]. The Appendix contains the proofs and algorithms and the main body contains the theorems and intuitions. We leave to future work experimental evaluation of energy. Our own prior work exploits the slow-to-fast insight to optimize *performance* (not energy) of interactive applications [29]. We achieved substantial performance improvements in simulation and on real systems by configuring Simultaneous Multi-Threading (SMT) hardware as a dynamic heterogeneous multicore [29]. No prior work presents an optimal algorithm or theory for energy efficiency under multiple latency constraints, nor provides guidelines for selecting core speeds. This work is the first formal analysis to deliver these properties for scheduling interactive workloads on heterogeneous multicore processors for energy minimization subject to multiple latency constraints.

## 2 Job, Processor, and Scheduling Models

This section and Table 1 describe our job, processor, and scheduling model.

*Job model* We focus on CPU intensive interactive services such as search, ads, finance option pricing, games, and serving dynamic web page content [7, 21, 26, 44]. Each interactive service request is a *job*. Each job has *work w* (service demand), which represents the number of CPU cycles the job takes to complete. Since it is often impossible to accurately predict a job's service demand [25], we model $w$ as a discrete random variable whose value is unknown until the job completes. We divide the service demand into $N$ bins and the *size* of the $i$-th bin is denoted by $w_i$, which we obtain by measuring the distribution of work for the application. The choice of "bin" sizes is determined by the measurement accuracy, and our model is not restricted to any particular choices. The job service demand $w$ follows a distribution that only takes values out of the set $\mathcal{W} = \{\tilde{w}_1, \tilde{w}_2, \cdots, \tilde{w}_N\}$, where we define $\tilde{w}_i = \sum_{j=1}^{i} w_j$, for $i = 1, 2, \cdots, N$. This assumption is not restrictive. In practice, a job's service demand cannot be continuous and is typically grouped into a finite number of bins [38].

| Definition | | Definition | |
|---|---|---|---|
| $w$ | CPU service demand | $x_i$ | Speed of core $i$ |
| $w_i$ | Size of the $i$-th demand bin | $z(x)$ | Power consumption |
| $f_i$ | Probability of demand $\tilde{w}_i = \sum_{j=1}^{i} w_j$ | $e(x)$ | Energy function |
| $F_i$ | Cumulative distribution | $L_p$ | $L_p$ norm with moment $p$ |
| $F_i^c$ | Complementary cumulative distribution | $\tilde{D}(p)$ | $L_p$ norm latency constraint |

**Table 1:** Symbols and definitions

Let $\{f_1, f_2, \cdots, f_N\}$ and $\{F_0, F_1, \cdots, F_N\}$ be the probability distribution and cumulative distribution of the job's service demand, respectively: $f_i = \Pr(w = \tilde{w}_i)$ and $F_i = \sum_{j=1}^{i} f_j$, for $i = 1, 2, \cdots, N$. While the service demand of any single job is unknown *a priori*, we assume the aggregate service demand distribution of jobs is measured with online or offline profiling as in previous work [25].

*Processor model* We adopt a standard processor model. With speed $x > 0$, a core will consume a power of $z(x)$. Correspondingly, the energy consumption per unit work is $e(x) = z(x)/x$. The processing time for a unit work increases linearly with respect to the inverse of core speed. Given a particular application, the effective speed $x$ and power $z(x)$ can be obtained by system measurements. Consequently, the effective speed $x$ may differ from the clock rate of CPU and both clock speed an power may vary depending on the application [16, 22].

We assume the energy function $e(x)$ is continuously differentiable, increasing, and strictly convex in $x \geq 0$. This assumption is validated extensively by both analytical models and measurement studies [16, 25, 38]. In practice, if a slower core consumes more power and thus energy than a fast one, it wont be built. Because of CMOS circuit characteristics, energy is well approximated as $e(x) = b \cdot x^{\alpha-1} + c$ for core speed $x$, where the power exponent $\alpha \geq 2$ and static energy $c \geq 0$ [10, 25]. We concentrate on heterogeneous multicores which consists of multiple diverse cores, but our approach applies to cores with multiple speeds realized with DVFS.

We refer to the core executing the $i$-th bin of a job's demand as core $i$, for $i = 1, 2, \cdots, N$. We denote the core speed and power consumption of core $i$ by $x_i$ and $z_i = z(x_i)$, respectively. The energy consumption per unit work of core $i$ is given by $e_i = e(x_i) = z(x_i)/x_i$. Two cores $i$ and $j$ may be equivalent in some cases, i.e., $x_i = x_j$, for $i, j = 1, 2, \cdots, N$. For example, one core will execute multiple bins when demand for this core differs between two or more jobs.

## 3  Scheduling objective —Energy

Our scheduling objective is minimize average energy on a heterogeneous processor when scheduling interactive jobs that are subject to multiple latency constraints. The scheduler determines the core speeds $x_i$ for each bin $i = 1, 2, \cdots, N$. We express the average energy consumption of a job as

$$\bar{e}(\mathbf{x}) = \sum_{i=1}^{N} \left[ \sum_{j=1}^{i} z_j \cdot \frac{w_j}{x_j} \right] \cdot f_i = \sum_{i=1}^{N} [1 - F_{i-1}] \cdot e(x_i) \cdot w_i, \tag{1}$$

where $e_i = e(x_i) = z(x_i)/x_i$ is the energy per unit work consumed by core $i$ and $\mathbf{x} = (x_1, x_2, \cdots, x_N)$ is a vector expression. The term "$\sum_{j=1}^{i} z_j \cdot \frac{w_j}{x_j}$" represents the energy consumption of a job with a service demand of $\sum_{j=1}^{i} w_j$ (which occurs with a probability of $f_i$), and hence we have the average energy consumption as $\sum_{i=1}^{N} \left[ \sum_{j=1}^{i} z_j \cdot \frac{w_j}{x_j} \right] \cdot f_i$. Equivalently, we can rewrite the average energy consumption as $\sum_{i=1}^{N} [1 - F_{i-1}] \cdot e(x_i) \cdot w_i$, where $(1 - F_{i-1})$ is the probability that the $i$-th bin of the service demand is processed (i.e., the probability that a job has at least a service demand of $\sum_{j=1}^{i} w_j$).

## 4  Scheduling constraints —Latency

Many prior studies mainly focused on *single* and *simple* latency constraints, such as maximum latency (deadline) or average latency [25, 38]. Motivated by recent work that addresses latency requirements in contexts such as load balancing [2, 28], we introduce the $L_p$ norm to generalize latency constraints. For concision, we sometimes abbreviate the $L_p$ norm with $L_p$. Specifically, given the core speeds $\mathbf{x} = (x_1, x_2, \cdots, x_N)$, we mathematically express the $L_p$ norm for latency as follows

$$D(p) = \left[ \sum_{i=1}^{N} (t_i)^p \cdot f_i \right]^{\frac{1}{p}} = \left\{ \sum_{i=1}^{N} \left[ \sum_{j=1}^{i} \frac{w_j}{x_j} \right]^p \cdot f_i \right\}^{\frac{1}{p}}, \tag{2}$$

where $p \geq 1$ and $t_i = \sum_{j=1}^{i} \frac{w_j}{x_j}$ is the latency of a job with a service demand of $\tilde{w}_i = \sum_{j=1}^{i} w_j$. The $L_p$ norm for latency generalizes over maximum and average latency. Given $p = \infty$, $L_\infty$ is maximum latency and given $p = 1$, $L_1$ is average latency. Intuitively, larger values of $p$ emphasize optimizing the latency of longer jobs, effectively mitigating unfairness and extreme outliers for long jobs [2, 5].

Latency variance determines the *predictability* of a scheduling algorithm [35] and depends on the $L_2$ and $L_1$ through the simple expression $L_2 - L_1$. For average latency and latency variance, we can apply various techniques, such as Chebyshev inequality, to bound tail distributions and estimate high-percentile latency. Thus, simultaneously considering multiple $L_p$ latency constraints, such as the $L_1$ and $L_2$ norms, well characterizes requirements on interactive applications [2–4, 28].

This paper focuses on interactive applications where the actual demand of individual jobs is unknown and hence all jobs have the same latency constraints, e.g., all web pages have similar latency constraints, since users will abandon the browser if responses are too slow. Differentiated services for different jobs are beyond the scope of this paper and could be interesting future work.

## 5  Problem Formulation and Algorithm

This section formalizes the energy minimization problem and presents the SEM scheduling algorithm, which minimizes energy subject to latency constraints.

The inputs to SEM are the probability distribution of service demand $f_i$, the size of each service demand bin $w_i$, and energy consumption per unit work $e(x)$ in terms of the processing speed $x$. SEM outputs the optimal job schedule, which

prescribes a sequence of core speeds $x_1, x_2, \cdots, x_N$, where $x_i$ is the core speed to process the $i$-th service demand bin. An incoming job with unknown service demand will execute on the prescribed sequence of core speeds until completion. For example, given an application that has jobs with service demands of $1, 2, 5$, or $10$ (units of work) and some probability distribution, then there are 4 service demand bins with the following sizes: $w_1 = 1, w_2 = 2 - 1 = 1, w_3 = 5 - 2 = 3, w_4 = 10 - 5 = 5$. Given a set of $L_p$ latency constraints, SEM determines the optimal core speed $x_i$ for executing each service demand bin $w_i$. For example, $x_1 = 1$ GHz, $x_2 = x_3 = 1.5$ Ghz, and $x_4 = 3$ GHz. This scheduling plan is determined offline and then used in deployment. In deployment, when a job arrives, it's service demand is unknown. SEM first processes the job on a 1 GHz core. If the job does not completed after 1 unit of work, SEM migrates the job to a 1.5 GHz core. If the job does not completed after processing another $w_2 + w_3 = 4$ units of work, SEM migrate it to a 3GHz core, and continue processing the job until it completes. Formally, this problem is stated as follows.

$$\textbf{P1}: \quad \min_{\mathbf{x}} \sum_{i=1}^{N} \{[1 - F_{i-1}] \cdot e(x_i) \cdot w_i\} \tag{3}$$

$$s.t., \quad \left\{ \sum_{i=1}^{N} \left[ \sum_{j=1}^{i} \frac{w_j}{x_j} \right]^{p_k} f_i \right\}^{\frac{1}{p_k}} \leq \tilde{D}(p_k), \tag{4}$$

$$\text{for } k = 1, 2, \cdots, K,$$

$$\mathbf{x} \succeq \mathbf{0}, \tag{5}$$

where $\succeq$ is an element-wise operator, constraining all the core speeds to be non-negative. This formulation assumes that the core speeds $x_1, x_2, \cdots, x_N$ can be continuously chosen from any non-negative values. In other words, here core speeds are unconstrained. (Section 8 shows how to handle the limited numbers of core speeds available in practice.) The objective function in (3) minimizes the average energy of all jobs. The latency constraints in (4) are imposed with $K$ different norms where $1 \leq p_1 < p_2 < \cdots < p_K \leq \infty$. Note that imposing a tail latency constraint of $L_\infty$ excludes outlier jobs, e.g., for 95-percentile latency, the 5% longest jobs are excluded by the $L_\infty$ norm.

This **P1** formulation is a convex optimization problem. The latency constraints in Inequality (4) are convex because $L_p$ norms are convex when $p \geq 1$. The speed constraints in Inequality (5) are linear. A linear combination of the energy consumption per unit work $e(x)$ is strictly convex in terms of the processing speed $x$ due to CMOS characteristics [25]. The objective function in (3) is also convex. Since **P1** is convex, there exist efficient algorithms that find the globally optimal solution, which we denote as $\mathbf{x}^*$.

We derive the solution to **P1** using a primal-dual iterative approach. Appendix A presents the algorithm and its proof. We set a threshold $\epsilon$ as a stopping criterion such that the iteration stops once the difference of the $L_2$ norm between two consecutively iterated values is below the threshold. The iterative approach has a iteration-complexity bounded by $\mathcal{O}(1\backslash\epsilon^2)$ [24].

Note that we analytically derive the solution to **P1** instead of using a convex solver. The analytical form exposes important properties of the optimal solution and has implications for hardware core choices that we discuss in Section 6 and Section 7. These properties cannot be derived using a convex solver.

Further note that we only compute an optimal schedule *once* offline for any given job service demand distribution and heterogeneous system. Our online scheduler simply applies the precomputed optimal schedule, executing a job on each core speed for the precomputed specified optimal time, until the job completes. Therefore, the computational overhead in deployment is negligible.

## 6 An Optimal Schedule Migrates from Slow to Fast Cores

Under the optimal schedule, core speeds monotonically increase as hardware processes more of the job's work. In other words, *an optimal scheduler need only migrate a job from slower to faster cores.* Theorem 1 formalizes this property. While prior studies [25, 29, 38, 41] show to use the "slow to fast" property under the maximum latency constraint in different contexts such as DVFS, in contrast, Theorem 1 is the first formal result that applies it to the more general case of *any* latency norm constraint and with *multiple* latency norm constraints.

**Theorem 1.** *The optimal core speeds that solve **P1** satisfy* $0 < x_1^* \leq x_2^* \leq \cdots \leq x_N^*$. *If only the $L_1$ latency constraint is imposed, then $x_1^* = x_2^* = \cdots = x_N^*$.*

*Proof.* Appendix B contains the proof. ∎

Theorem 1 tells us, without a priori knowledge of each job's service demand, an optimal schedule first processes a job on a slow core. If the job does not complete within some time interval (because it is long), SEM migrates it to faster cores. Thus, a short job completes on slower cores to save energy while a long job uses faster cores to meet the latency constraints. Consequently, the average energy consumption is minimized while satisfying latency constraints.

The intuition behind Theorem 1 is that long jobs have a greater impact on latency constraints. In particular, the latency norm constraint specified by Equation (2) is mostly dominated by long jobs (the larger $p_k$, the more dominated by long jobs, which can be seen by taking the partial derivative of (2) with respect to the latency experienced by jobs with various demands). In the extreme case, when $p_k \to \infty$, only the maximum latency incurred by the longest jobs is important. Thus, we want to process the long jobs fast enough to meet the latency constraints. On the other hand, processing short jobs using slower cores saves energy without penalizing the latency constraints.

If only the average latency constraint $(p = 1)$ is considered, Theorem 1 reduces to a special case where $x_1^* = x_2^* = \cdots = x_N^*$, i.e., the optimal schedule uses a homogeneous processor. Intuitively, this reduction holds because delaying short and long jobs have the same impact on the $L_1$ norm. More formally, Appendix A derives that $e(x_i)x_i^2$ is the same for all $i = 1, 2, \cdots, N$ and hence, homogeneous speeds are optimal when only satisfying the $L_1$ norm latency constraint. For all other latency constraints $(p > 1)$ and multiple constraints, the optimal energy-efficient schedule requires heterogeneous processors.

# 7 Implications for Cores in a Heterogeneous System

This section analyzes how latency constraints, workload characteristics, and core power/performance characteristics effect core choices in a heterogeneous system.

## 7.1 Effect of latency constraints on heterogeneity

Given Theorem 1, a key question is what core speeds to include in a heterogeneous system. In practice, the fastest cores are limited by physics and the software will be tuned such that the fastest core speed can satisfy the most demanding jobs. We therefore exploit this theorem to select the remaining lower power cores by investigating the ratio of the fastest $x_N^*$ to the slowest $x_1^*$ speed. We define this ratio as the *degree of heterogeneity*, giving a formal quantitative guideline for selecting core speeds in a heterogeneous processor. Our analysis shows that *more heterogeneity is desired for larger $p$* in the $L_p$ norm constraint.

We derive this result using a widely-used class of energy functions [25] expressed in the form $e(x) = b \cdot x^{\alpha-1}$, where $b > 0$ and $\alpha \geq 2$ (corresponding to a power function of $z(x) = b \cdot x^{\alpha}$ [25]). The lack of a closed-form expression of the optimal core speeds $\mathbf{x}^*$ makes it prohibitive to derive the exact value of the degree of heterogeneity. We instead exploit monotonicity to derive upper and lower bounds, using Theorem 2 to show that degree of heterogeneity is monotonically increasing in $p \geq 1$.

**Theorem 2.** *Given $e(x) = b \cdot x^{\alpha-1}$ and one $L_p$ latency constraint, then the degree of heterogeneity $\frac{x_N^*}{x_1^*}$ increases with increasing $p$ for $p \geq 1$.*
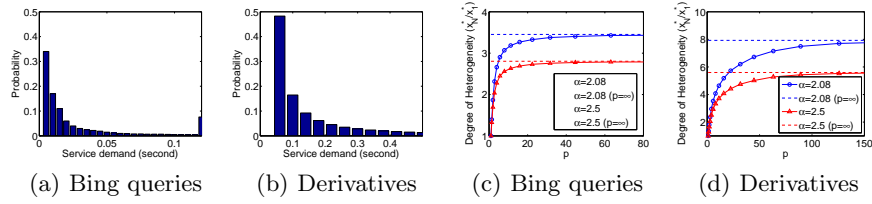
*Proof.* Appendix C contains the proof. ∎

Theorem 2 proves that as $p \geq 1$ increases, the optimal degree of heterogeneity also increases; the latency constraint thus imposes the optimal choice of core speeds. More precisely, given two different values of $p$, we can select two different latency constraints, under which the corresponding minimum core speeds are the same using the optimal job schedule. Under a latency constraint with a larger $p$ value, long jobs require faster cores, because larger values of $p$ place a more stringent requirement on the latency of longer jobs. Thus, if $p$ increases, so does $x_N^*/x_1^*$. Furthermore, we prove in Lemma 1 in Appendix C that the degree of heterogeneity is a constant for a given $p$ regardless of latency constraints, which establishes *hardware requires more heterogeneity for larger $p$.*

## 7.2 How much heterogeneity is desirable?

This section explores how much heterogeneity is desirable. We use Theorem 2 to derive both upper and lower bounds on degree of heterogeneity in Theorem 3. This result delivers quantitative guidance for selecting the cores in heterogeneous multicore processors for interactive applications.

(a) Bing queries    (b) Derivatives    (c) Bing queries    (d) Derivatives

**Fig. 1:** (a) (b) Service demand distributions of Bing and Financial derivative workloads. Most jobs are short, but long jobs are not negligible. (c) (d) Degree of heterogeneity as a function of $p$ given one $L_p$ constraint and power model: $z(x) = 21 \cdot x^{\alpha}$.

**Theorem 3.** *Given $e(x) = b \cdot x^{\alpha-1}$ and $K$ $L_p$ latency constraints specified by $1 \leq p_1 \leq p_2 \leq \cdots \leq p_K \leq \infty$, then the degree of heterogeneity $\frac{x_N^*}{x_1^*}$ satisfies:*

$$1 \leq \frac{x_N^*}{x_1^*} \leq \left(\frac{1}{f_N}\right)^{\frac{1}{\alpha}} \tag{6}$$

*where $f_N$ is the probability that a job has the maximum service demand of $\tilde{w}_N$.*

We call a latency constraint **dominant** *if and only if satisfying it ensures that all the other latency constraints, if any, are also satisfied under the optimal schedule. Thus, the dominant latency constraint is the most stringent requirement. When average latency is dominant, the first inequality above becomes an equality: $x_N^*/x_1^* = 1$. When maximum latency is dominant, the second inequality becomes equality: $\frac{x_N^*}{x_1^*} = \left(\frac{1}{f_N}\right)^{\frac{1}{\alpha}}$.*

*Proof.* Appendix D contains the proof. ∎

Theorem 3 has two interesting implications.

*1. Workload heterogeneity prefers hardware heterogeneity.* The upper bound on the degree of heterogeneity increases as $f_N$ decreases (i.e., with fewer long jobs). When the workload is homogeneous, all jobs have the same service demand and $f_N = 1$. In this case, Theorem 3 indicates that $x_N^*/x_1^* = 1$ and homogeneous hardware is optimal. For a heterogeneous workload where $f_N$ is small, the value of $x_N^*/x_1^*$ may become very large. When slow cores complete short jobs, they save energy, whereas with the optimal schedule, the fastest processors process long jobs to satisfy the maximum latency constraint without incurring too much average energy, since $f_N$ is small.

*2. Core power and performance influences on hardware heterogeneity.* When the speed of a core increases, so does its power consumption. We observe from (6) that the upper bound on the degree of heterogeneity decreases with $\alpha$. A larger $\alpha$ indicates power consumption grows faster than core speed and hence using fast cores will significantly increases average energy consumption and degree of heterogeneity will be smaller.

*Example* We consider two example interactive workloads, Bing web search and Monte Carlo financial pricing (see elsewhere for details [29]). They illustrate how latency constraints, workload, and core performance and power characteristics

affect the desired heterogeneity. Figure 1(a) and Figure 1(b) show the distributions of service demand for the two applications, measured in terms of the job processing time on an Intel i7-2600 Sandy Bridge core. The demand spike in Figure 1(a) occurs because the search engine caps job processing time at 120 ms and returns the top results found so far. Search engines often cap query processing time and return partial results to tradeoff quality and response time [19].

Figure 1(c) and 1(d) show how the degree of heterogeneity ($Y$-axis $x_n^*/x_1^*$) changes as a function of $p$ in $L_p$ with Bing and financial applications, respectively, when we can choose any core speed. We normalize speed to an i7-2600 Sandy Bridge core and use the power model: $z(x) = 21 \cdot x^\alpha$, because $z(1) = 21W$ is the power consumption of the i7-2600 Sandy Bridge core. Blue and red lines represent the cases of $\alpha = 2.08$ (a lower energy cost for performance) and $\alpha = 2.5$ (a higher energy cost for performance) respectively.

Figure 1(c) and 1(d) confirm Theorems 2 and 3. (1) When $p$ increases, the degree of heterogeneity increases and has an upper bound, as predicted. In particular, a homogeneous processor is optimal in terms of energy consumption when $p = 1$ (average latency), whereas the maximum degree of heterogeneity is desirable when $p = \infty$ (a deadline). (2) The degree of heterogeneity decreases with larger $\alpha$ because faster cores consume proportionally more energy. (3) Comparing Figure 1(c) and 1(d) shows financial derivative pricing requires a higher degree of heterogeneity than Bing web search given the same $p$ because the longest jobs are rarer in derivatives ($f_N$ is smaller). The rarer the long jobs, the faster the fastest core we can choose without compromising average energy because the prolific short jobs execute on the slowest low power cores.

## 8  Discrete Core Speeds, Migration, and Congestion

This section extends SEM to address the following practical considerations: (1) a limited selection of core speeds, (2) job migration overhead, and (3) congestion due to multiple jobs competing for the same core(s).

*Discrete core speeds* Given a set of core speeds, $0 < s_1 \leq s_2 \leq \cdots \leq s_M$, we formulate our problem as follows:

$$\mathbf{P2}: \ \min_{\mathbf{x}} \sum_{i=1}^{N} \{[1 - F_{i-1}] \cdot e(x_i) \cdot w_i\} \tag{7}$$

$$s.t., \ \text{Constraint (4)} \tag{8}$$

$$x_i \in \{s_1, s_2, \cdots, s_M\}, i = 1, 2, \cdots, N. \tag{9}$$

**P2** is a combinatorial optimization problem, which is notoriously difficult to solve [41]. We use an efficient branch-and-bound algorithm to produce solutions arbitrarily close-to-optimal. A greedy solution finds a schedule that will consume more energy than the optimal schedule (i.e., the upper bound), whereas the job schedule obtained by replacing "$x_i \in \{s_1, s_2, \cdots, s_M\}$" with $x_i \in [s_1, s_M]$ and then using convex optimization will produce an average energy consumption that is less than the optimal schedule (i.e., the lower bound). By iteratively finding and

refining the upper and lower bounds until the gap becomes sufficiently small, we identify a schedule arbitrarily close to the optimal schedule [8]. Appendix E contains the details of the solution and its derivation.

**P2** is an NP-hard problem, even if only the maximum latency constraint is considered [41]. Without specifying the maximum number of iterations, the proposed algorithm may iterate up to $M^N$ times, enumerating all the possible solutions in the worst case. Nevertheless, the beauty of branch-and-bound algorithm is that it typically converges much faster, which we also observe. In fact, with an appropriately-set stopping criterion, the number of iterations required for convergence is upper bounded, and in practice, the actual number of iterations is typically even much smaller than the upper bound. The complete analysis of convergence rate is beyond our scope, and interested readers are referred to the literature [8].

Moreover, as we discussed in Section 5, we only compute an optimal schedule *once* offline for any given job service demand distribution and heterogeneous processor. Our online scheduler simply applies the precomputed optimal schedule. Therefore, the computational overhead in deployment is negligible.

*Migration overhead* Migrating a job from one core to another incurs overhead from copying job state and warming up caches. Our experiments show that job migration overheads are fairly small on both web search [19] and interactive finance applications. One migration is less than 50 microseconds, less than 0.1% of the maximum latency requirement in the order of 100 milliseconds. Moreover, a job can only migrate up to $Q - 1$ times, where $Q$ is the number of different core speeds. Because $Q$ is very small $(2 \sim 4)$ in practice and many short jobs completed on slow cores, SEM often does not incur much migration overhead.

To extend our solution when migration costs are high, e.g., migrating a job between two servers, we describe a heuristic approach to incorporate migration overhead in the analytical model. This approach is conservative and assumes worst-case migration overhead. More specifically, let $\tau^o$ represent the migration overhead, quantified by the time during which a core cannot process any work. In the worst case, a job with a demand of $\tilde{w}_i = \sum_{j=1}^{i} w_j$ may migrate up to $(i-1)$ times, for $i = 1, 2, \cdots, N$. Thus, the new worst-case latency constraint becomes

$$\left\{ \sum_{i=1}^{N} \left[ \sum_{j=1}^{i} \frac{w_j}{x_j} + (i-1) \cdot \tau^o \right]^{p_k} f_i \right\}^{\frac{1}{p_k}} \leq \tilde{D}(p_k). \tag{10}$$

By neglecting the constant energy consumption incurred by the migration process in the worst case, we reformulate the energy minimization problem **P2** by replacing the latency constraint (8) with (10) to account for the migration overhead. The solution can be found in a similar way following our preceding analysis.

*Congestion* We briefly discuss how to apply SEM as a building block when congestion or queuing delay results in multiple jobs demanding the same core at the same time. A key observation is that the presence of congestion may cause a violation in the latency constraints if we directly apply SEM. To satisfy the

desired latency that includes both processing delay and queueing delay, we can impose a more stringent constraints for the processing delay which, if *appropriately* chosen and after adding the queueing delay, will satisfy the total latency constraints. To choose the appropriate $L_p$ norm constraint to handle this delay, we propose integral control to dynamically adjust the processing delay constraint based on the difference between the observed latency and the target latency (latency constraint). The control function is expressed as

$$\tilde{D}_i(p_k) = \tilde{D}_{i-1}(p_k) + V \cdot d_i(p_k), \text{ for } k = 1, 2, \cdots, K,$$

where $i = 1, 2, \cdots$ represents time steps, $\tilde{D}_i(p_k)$ is the output of the integral controller at time $i$ representing the augmented $L_p$ norm constraint on the processing delay. $V > 0$ defines the ratio of the control adjustment to the control error and $d_i(p_k)$ is the difference between the target and observed latency. Thus, if the observed latency is greater than the constraint, $d_i(p_k) < 0$, a more stringent processing delay constraint, $\tilde{D}_i(p_k)$, will be imposed for the next time step, and vice versa.

Finally, note that using the above method to address congestion will not alter the value of $p$. Thus, our slow to fast scheduling insight and the quantitative upper and lower bounds on the ratio of fast to slow core speeds still hold.

## 9   Related Work

*Heterogeneous multicore processors* As computer architects face the end of Dennard scaling, they are turning to heterogeneous multicore processors, which combine high performance but high power cores with lower power and lower performance cores to meet a variety of performance objectives, i.e., throughput, energy, power, etc. To effectively utilize these systems, a scheduler must match jobs to an appropriate core. Four types of schedulers have been proposed to allocate jobs or parts of jobs to different cores. (1) With known or predicted resource demand, incoming jobs are scheduled to the most appropriate core [11, 13, 42]. (2) With known performance requirements, latency-sensitive applications such as games or videos are processed by fast cores, whereas latency-tolerant applications such as background services are processed by slow cores [17,27,31]. (3) With known job characteristics, complementary job allocation is applied to maximize the server utilization while avoiding resource bottlenecks (e.g., memory-intensive jobs and CPU-intensive jobs are allocated to the same server [37]). (4) If a single job has different phases [23,32,33], such as parallel phases and sequential phases, schedulers map the sequential phase on a high-performance core and the parallel phase on a number of energy-efficient cores.

*$L_p$ norms and multiple latencies* Because the $L_p$ norms are a general class of constraints, researchers have applied them in various contexts, such as minimizing the total latency via online load balancing [5, 34] and multi-user scheduling of wireless networks [43]. Our study considers multiple $L_p$ norm latency constraints simultaneously for individual interactive services. Prior work mainly considers multiple latency constraints to provide differentiated performance guarantees to

different traffic classes [12,36], whereas we exploit the diversity of demand within the requests, without requiring knowledge about the demand of any individual request, to meet constraints for a variety of interactive applications.

*Latency sensitive and real-time scheduling* Related work also considers exploiting heterogeneous processors and DVFS to improve energy-efficiency for latency-sensitive and real-time jobs [1, 25, 38, 39, 41]. Some of them [1, 39, 42, 44] assume that the service demand of each job is either known or accurately predicted, which is not available for many applications. Other studies on DVFS and real-time systems assume unknown service demand [25, 38, 41], but they consider a hard deadline as the only latency constraint. Our prior work [29] studies scheduling interactive workloads on a heterogeneous processor for quality/throughput maximization (not energy minimization) subject to a single deadline constraint. While it also leverages the "slow to fast" insight, it always uses fast cores first whenever they are available for performance optimization. In contrast, SEM starts jobs on slow cores and migrates them to fast cores along the execution to minimize energy. Moreover, this prior work [29] does not address multiple latency constraints and it does not deliver quantitative insights for selecting cores in heterogeneous processors. To the best of our knowledge, we offer the first formal analysis to characterize the optimal schedule and hardware design for scheduling latency-sensitive jobs on heterogeneous processors with multiple latency constraints without requiring a priori knowledge of the service demand of each individual job.

## 10  Conclusion

This paper presents an efficient scheduling algorithm for interactive jobs on heterogeneous processors subject to multiple latency constraints expressed in the form of $L_p$ norms and optimizes energy. We introduce the SEM scheduling which advances the existing research in two key ways. (1) The SEM algorithm does not rely on the service demand of each individual job, which is difficult and even impossible to obtain in many interactive applications such as web search. (2) The SEM algorithm explicitly incorporates multiple $L_p$ norm latency constraints which, compared to prior work, more accurately characterize the explicit and implicit multiple service level agreements on the latency of interactive applications. We prove that an optimal schedule only migrates jobs from slower to faster cores. Moreover, we quantify how to select cores in heterogeneous hardware for interactive applications. The more the system needs to limit outliers, the more heterogeneous the hardware needs to be. The more heterogeneous the workload service demand is, the less power additional performance costs and the more heterogeneous the hardware needs to be.

## References

1. ALBERS, S., MULLER, F., AND SCHMELZER, S. Speed scaling on parallel processors. In *SPAA* (2007).
2. ANAND, S., GARG, N., AND KUMAR, A. Resource augmentation for weighted flow-time explained by dual fitting. In *SODA* (2012).

3. Azar, Y., and Epstein, A. Convex programming for scheduling unrelated parallel machines. In *STOC* (2005).

4. Azar, Y., Epstein, L., Richter, Y., and Woeginger, G. J. All-norm approximation algorithms. In *SWAT* (2002).

5. Bansal, N., and Pruhs, K. Server scheduling in the $l_p$ norm: A rising tide lifts all boat. In *STOC* (2003).

6. Bertsekas, D. P., and Tsitsiklis, J. N. *Parallel and Distributed Computation: Numerical Methods*. Athena Scientific, 1989.

7. Bornholt, J., Mytkowicz, T., and McKinley, K. S. The model is not enough: Understanding energy consumption in mobile devices. In *Hot Chips* (2012).

8. Boyd, S., Ghosh, A., and Magnani, A. Branch and bound methods, `http://www.stanford.edu/class/ee392o/bb.pdf`, 2003.

9. Boyd, S., and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2004.

10. Brooks, D., Bose, P., Schuster, S., Jacobson, H., Kudva, P., Buyuktosunoglu, A., Wellman, J., Zyuban, V., Gupta, M., and Cook, P. Power-aware microarchitecture: Design and modeling challenges for next generation microprocessors. In *Micro* (2000).

11. Cao, T., Blackburn, S. M., Goa, T., and McKinley, K. S. The yin and yang of power and performance for asymmetric hardware and managed software. In *ISCA* (2012).

12. Chao, H. J., and Uzun, N. An atm queue manager handling multiple delay and loss priorities. *IEEE/ACM Trans. Networking 3* (December 1995), 652–659.

13. Chen, J., and John, L. K. Efficient program scheduling for heterogeneous multicore processors. In *DAC* (2009).

14. Dean, J., and Barroso, L. A. The tail at scale. *CACM 56*, 2 (2013), 74–80.

15. DeCandia, G., Hastorun, D., Jampani, M., Kakulapati, G., Lakshman, A., Pilchin, A., Sivasubramanian, S., Vosshall, P., and Vogels, W. Dynamo: Amazon's highly available key-value store. In *SOSP* (2007).

16. Esmaeilzadeh, H., Cao, T., Xi, Y., Blackburn, S. M., and McKinley, K. S. Looking back on the language and hardware revolutions: Measured power, performance, and scaling. In *ASPLOS* (2011).

17. Greenhalgh, P. Big.LITTLE processing with ARM Cortex-A15 & Cortex-A7. *ARM Whitepaper* (September 2011).

18. Harchol-Balter, M. The effect of heavy-tailed job size distributions on computer system design. In *Applications of Heavy Tailed Distributions in Economics* (1999).

19. He, Y., Elnikety, S., Larus, J., and Yan, C. Zeta: Scheduling interactive services with partial execution. In *SOCC* (2012).

20. Im, S., and Moseley, B. An online scalable algorithm for minimizing $l_k$-norms of weighted flow time on unrelated machines. In *SODA* (2011).

21. Janapa Reddi, V., Lee, B. C., Chilimbi, T., and Vaid, K. Web search using mobile cores: Quantifying and mitigating the price of efficiency. In *ISCA* (2010).

22. Kotla, R., Devgan, A., Ghiasi, S., Keller, T., and Rawson, F. Characterizing the impact of different memory-intensity levels. In *WWC* (2004).

23. Kumar, R., Farkas, K. I., Jouppi, N. P., Ranganathan, P., and Tullsen, D. M. Single-ISA heterogeneous multi-core architectures: The potential for processor power reduction. In *MICRO* (2003).

24. Lan, G., Lu, Z., and Monteiro, R. D. Primal-dual first-order methods with $\mathcal{O}(1\backslash\epsilon)$ iteration-complexity for cone programming. *Mathematical Programming 126*, 1 (2011), 1–29.

25. LORCH, J. R., AND SMIT, A. J. Improving dynamic voltage scaling algorithms with PACE. In *SIGMETRICS* (2001).

26. MEISNER, D., SADLER, C. M., BARROSO, L. A., WEBER, W.-D., AND WENISCH, T. F. Power management of online data-intensive services. In *ISCA* (2011).

27. NATHUJI, R., ISCI, C., AND GORBATOV, E. Exploiting platform heterogeneity for power efficient data centers. In *ICAC* (2007).

28. PRUHS, K. Competitive online scheduling for server systems. In *SIGMETRICS* (2007).

29. REN, S., HE, Y., ELNIKETY, S., AND MCKINLEY, K. S. Exploiting processor heterogeneity in interactive systems. In *ICAC* (2013).

30. REN, S., HE, Y., AND MCKINLEY, K. S. A theoretical foundation for scheduling and designing heterogeneous processors for interactive applications. In *DISC* (2014).

31. SRINIVASAN, S., IYER, R., ZHAO, L., AND ILLIKKAL, R. HeteroScouts: Hardware assist for OS scheduling in heterogeneous CMPs. In *SIGMETRICS* (2011).

32. SULEMAN, M. A., MUTLU, O., QURESHI, M. K., AND PATT, Y. N. Accelerating critical section execution with asymmetric multi-core architectures. In *ASPLOS* (2009).

33. SULEMAN, M. A., PATT, Y. N., SPRANGLE, E., ROHILLAH, A., GHULOUM, A., AND CARMEAN, D. Asymmetric chip multiprocessors: Balancing hardware efficiency and programmer efficiency. *TR-HPS-2007-001* (2007).

34. SURI, S., TÓTH, C. D., AND ZHOU, Y. Selfish load balancing and atomic congestion games. In *SPAA* (2004).

35. WIERMAN, A., AND HARCHOL-BALTER, M. Classifying scheduling policies with respect to higher moments of conditional response time. In *SIGMETRICS* (2005).

36. XIE, Y., AND YANG, T. Cell discarding policies supporting multiple delay and loss requirements in atm networks. In *Globecom* (1997).

37. XIONG, W., AND KANSAL, A. Energy efficient data intensive distributed computing. In *IEEE Data Eng. Bull.* (2011).

38. XU, R., XI, C., MELHEM, R., AND MOSS, D. Practical PACE for embedded systems. In *EMSOFT* (2004).

39. YAO, F. F., DEMERS, A. J., AND SHENKER, S. J. A scheduling model for reduced CPU energy. In *FOCS* (1995).

40. YI, J., MAGHOUL, F., AND PEDERSEN, J. Deciphering mobile search patterns: A study of Yahoo! mobile search queries. In *WWW* (2008).

41. YUAN, W., AND NAHRSTEDT, K. Energy-efficient CPU scheduling for multimedia applications. *ACM Trans. Computer Systems 24*, 3 (2006), 292–331.

42. YUN, H., WU, P.-L., ARYA, A., KIM, C., ABDELZAHER, T. F., AND SHA, L. System-wide energy optimization for multiple DVS components and real-time tasks. *Real-Time Systems 47*, 5 (2011), 489–515.

43. ZENG, W., NG, C., AND MEDARD, M. Joint coding and scheduling optimization in wireless systems with varying delay sensitivities. In *SECON* (2012).

44. ZHU, Y., AND REDDI, V. J. High-performance and energy-efficient mobile web browsing on big/little systems. In *HPCA* (2013).

# Appendices

## A    Optimal algorithm to solve P1

This section presents an optimal algorithm to solve **P1** based on a primal-dual approach. Instead of resorting to a convex solver, we choose to analytically derive the solution to **P1**, which exposes important insights of the optimal solution that we explore in Section 6 and 7.

First, we rewrite the constraints (4) as follows

$$\sum_{i=1}^{N} \left[ \sum_{j=1}^{i} \frac{w_j}{x_j} \right]^{p_k} \cdot f_i \leq D_{p_k} \triangleq \left[ \tilde{D}(p_k) \right]^{p_k}, \tag{11}$$

for $k = 1, 2, \cdots, K$. Thus, we reformualte problem **P1** as an equivalent problem of minimizing (3) over $\mathbf{x} \succeq \mathbf{0}$ subject to (11). We refer to this equivalent problem as **P1'** and shall solve it in the remainder of this section. While the maximum latency constraint (i.e., $L_\infty$ norm) can be treated separately using "$\sum_{n=1}^{N} \frac{w_n}{x_n} \leq \tilde{D}(\infty)$" for rigorousness, we still use (11) to capture the maximum latency constraint by letting $p_k \to \infty$, slightly abusing the norm equation for succinct notation. We also define $F_0 = 0$ and $F_i^c = 1 - F_{i-1}$ as the complementary cumulative distribution.

Due to the strict convexity of (3), the solution to **P1'**, if any, is unique, and strong duality holds [9]. In the following, we solve **P1'** using a primal-dual approach in two steps: (1) derive the dual and then (2) maximize the dual.

### A.1    Deriving the dual of P1'

We derive the dual of **P1'** by first writing the Lagrangian of **P1'** as follows

$$\begin{aligned}
\mathcal{L}(\lambda, \mathbf{x}) = &\sum_{i=1}^{N} \left\{ [1 - F_{i-1}] \cdot e(x_i) \cdot w_i \right\} \\
&+ \sum_{k=1}^{K} \lambda_k \left\{ \sum_{i=1}^{N} \left[ \sum_{j=1}^{i} \frac{w_j}{x_j} \right]^{p_k} \cdot f_i - D_{p_k} \right\},
\end{aligned} \tag{12}$$

where $\lambda = (\lambda_1, \lambda_2, \cdots, \lambda_K)$ and $\lambda_k \geq 0$ is the Lagrangian multiplier associated with the $L_{p_k}$ norm constraint. The dual of **P1'**, denoted by $g(\lambda)$, can be derived by minimizing the Lagrangian (12) over $\mathbf{x} \succeq \mathbf{0}$, i.e., $g(\lambda) = \min_{\mathbf{x} \succeq \mathbf{0}} \mathcal{L}(\lambda, \mathbf{x})$.

Given a set of Lagrangian multipliers $\lambda$, we denote the optimal primal variables by $\mathbf{x}^*(\lambda)$. The nonlinearity and coupling of $\mathbf{x}$ in the Lagrangian (12) make it non-trivial to directly obtain $\mathbf{x}^*(\lambda) \succeq \mathbf{0}$. Therefore, we adopt the Gauss-Seidel method to iteratively find $\mathbf{x}^*(\lambda) \succeq \mathbf{0}$ [6]. The basic idea of the Gauss-Seidel method is that instead of minimizing the Lagrangian over all the core speeds simultaneously, we sequentially minimize the Lagrangian over one core speed at a time while fixing the other core speeds. Specifically, by taking the first-order partial derivative of (12) with respect to $x_i$ and setting it to zero, we obtain

$$e'(x_i)x_i^2 - \sum_{j=i}^{N} \left\{ \sum_{k=1}^{K} \left[ \lambda_k p_k \left( \sum_{m=1}^{j} \frac{w_m}{x_m} \right)^{p_k - 1} \frac{f_j}{F_i^c} \right] \right\} = 0, \tag{13}$$

---

**Algorithm 1** Ideal

---

1: $t_1 \leftarrow 1$; initialize $\lambda^0 \succeq \mathbf{0}$ and $\lambda^1 \succeq \mathbf{0}$
2: **while** $\| \lambda^{(t_1)} - \lambda^{(t_1-1)} \|_2 > \epsilon$ **do**
3:    $t_2 \leftarrow 1$; initialize $\mathbf{x}^{(0)}(\lambda^{(t_1)})$ and $\mathbf{x}^{(1)}(\lambda^{(t_1)})$
4:    **while** $\| \mathbf{x}^{(t_2)}(\lambda^{(t_1)}) - \mathbf{x}^{(t_2-1)}(\lambda^{(t_1)}) \|_2 > \epsilon$ **do**
5:       **for** $i = 1 \to N$ **do**
6:          Solve $\tilde{x}_i^{(t_2)}(\lambda^{(t_1)}, \tilde{x}_i^{(t_2)}, \cdots, \tilde{x}_{i-1}^{(t_2)}, \tilde{x}_{i+1}^{(t_2-1)}, \cdots, \tilde{x}_N^{(t_2-1)})$ based on (13)
7:       **end for**
8:       $t_2 \leftarrow t_2 + 1$
9:    **end while**
10:   Update $\lambda^{(t_1+1)}$ according to (14), and $t_1 \leftarrow t_1 + 1$
11: **end while**
12: **return** $\mathbf{x}^* = \mathbf{x}^*(\lambda^{(t_1-1)})$

---

where $F_i^c = \sum_{j=i}^N f_j = 1 - F_{i-1}$, for $i = 1, 2, \cdots, N$. It is easy to show that there exists a unique value, denoted by $\tilde{x}_i(\lambda) \geq 0$, that equates (13) and hence minimizes (12) under the assumption that the other variables $\mathbf{x} \backslash \{x_i\}$ are held constant. Because the lefthand side of (13) is monotonically increasing, the bi-section method computes $\tilde{x}_i(\lambda)$ in an efficient manner. To highlight the dependence of $\tilde{x}_i(\lambda)$ on $\mathbf{x} \backslash \{x_i\} = (x_1, \cdots, x_{i-1}, x_{i+1}, \cdots, x_N)$, we rewrite $\tilde{x}_i(\lambda)$ as $\tilde{x}_n(\lambda) = \tilde{x}_i(\lambda, \mathbf{x} \backslash \{x_i\})$ wherever applicable. Following the strict convexity of the Lagrangian (12) in $\mathbf{x} \succeq \mathbf{0}$, the Gauss-Seidel based iterative algorithm is guaranteed to converge to a unique and optimal point, denoted by $\mathbf{x}^*(\lambda)$, that minimizes the Lagrangian [6]. In practice, we test for convergence based on certain termination criteria, e.g., $\| \mathbf{x}^{(t)}(\lambda) - \mathbf{x}^{(t-1)}(\lambda) \|_2 \leq \epsilon$, where $\epsilon$ is a sufficiently small positive number. The formal description of the Gauss-Seidel based iterative algorithm is provided in Lines 5–7 of Algorithm 1.

While other techniques, such as parallel Jacobian updates, may apply in some cases and will speed up convergence, the potentially slow convergence of Gauss-Seidel is not a concern for our purposes. We *only* execute our algorithm when the service demand distribution or the latency constraints change, not when each job arrives. Moreover, Gauss-Seidel is sufficient for deriving the insights of the optimal solution.

## A.2 Maximizing the dual of P1'

This subsection maximizes the dual $g(\lambda) = \min_{\mathbf{x} \succeq \mathbf{0}} \mathcal{L}(\lambda, \mathbf{x})$ over $\lambda \succeq \mathbf{0}$. Due to the lack of a closed-form expression of $g(\lambda)$, we resort to the sub-gradient method to maximize $g(\lambda)$. Specifically, for $k = 1, 2, \cdots, K$, we iteratively update $\lambda_k^{(t)}$ over $t = 1, 2, \cdots$, following the recursion below.

$$\lambda_k^{(t)} = \left\{ \lambda_k^{(t-1)} + a(t) \left[ \sum_{i=1}^N \left\{ \left[ \sum_{j=1}^i \frac{w_j}{x_j^*(\lambda^{(t-1)})} \right]^{p_k} f_i \right\} - D_{p_k} \right] \right\}^+, \qquad (14)$$

where $\sum_{i=1}^N \left\{ \left[ \sum_{j=1}^i \frac{w_j}{x_j^*(\lambda)} \right]^{p_k} \cdot f_i \right\} - D_{p_k}$ is the sub-gradient of $g(\lambda)$ with respect to $\lambda_k$ and the step size $a(t)$ is chosen to be a sufficiently small number, ensuring the

convergence of the update (14) to the optimal dual variables $\lambda^*$ [6]. For example, we can choose $a(t)$ as $\frac{1+q}{t+q}$, where $q$ is a nonnegative constant.

Now, let us summarize the primal-dual approach in Algorithm 1 that solves the problem **P1'** and hence **P1**. There are two major loops in Algorithm 1. The outmost loop applies the sub-gradient method to find the optimal dual variables. The inner loop (Lines 4–6) describes the Gauss-Seidel method of iteratively computing $\mathbf{x}^*$ given a set of dual variables. The outer loop converges to the optimal dual variables $\lambda^* = (\lambda_1^*, \lambda_2^*, \cdots, \lambda_K^*)$, based on which the inner loop solves the optimal primal variables $\mathbf{x}^* = (x_1^*, x^*, \cdots, x_N^*)$. Below, we formally establish the optimality of Algorithm 1.

**Theorem 4.** *As $\epsilon \to 0$ in Algorithm 1, Ideal yields the optimal solution to the problem **P1** defined in* (3)–(5).

*Proof. Due to the convexity of **P1'** and the equivalence between **P1'** and **P1**, it follows immediately that Algorithm 1 proposed based on a primal-dual approach solves **P1** [6].* ∎

## B    Proof of Theorem 1

For convenience, we restate the theorem.

**Theorem 1.** *The optimal core speeds that solve **P1** satisfy $0 < x_1^* \le x_2^* \le \cdots \le x_N^*$. If only the $L_1$ latency constraint is imposed, then $x_1^* = x_2^* = \cdots = x_N^*$.*

We prove Theorem 1 by contradiction. Suppose that there exists an $i \in \{1, 2, \cdots, N-1\}$ such that the optimal service speeds satisfy $x_{i+1}^* < x_i^*$. Based on the strong duality of **P1'**, the optimal core speeds can be obtained by minimizing Lagrangian given the optimal dual variables $\lambda^*$. In other words, the optimal speed of core $i$, $x_i^* = x_i^*(\lambda^*)$, satisfies,

$$e'(x_i^*)(x_i^*)^2 = \sum_{j=i}^{N} \left\{ \sum_{k=1}^{K} \left[ \lambda_k^* p_k \left( \sum_{m=1}^{j} \frac{w_m}{x_m^*} \right)^{p_k - 1} \right] \frac{f_j}{F_i^c} \right\} \tag{15}$$

where $F_i^c = 1 - F_{i-1}$. By assumption, we have $x_i^* > x_{i+1}^*$. For notational convenience, we define

$$z(i, k) = \sum_{j=i}^{N} \left[ \left( \sum_{m=1}^{j} \frac{w_m}{x_m^*} \right)^{p_k - 1} \cdot \frac{f_j}{F_i^c} \right]. \tag{16}$$

Thus, it follows immediately $e'(x_i^*)(x_i^*)^2 = \sum_{k=1}^{K} \lambda_k^* p_k z(i, k)$. Because $e'(x)$ is positive and increasing (which follows from the convexity of $e(x)$), we have $e'(x_i^*)(x_i^*)^2 > e'(x_{i+1}^*)(x_{i+1}^*)^2$, i.e., $\sum_{k=1}^{K} \lambda_k^* p_k z(i, k) > \sum_{k=1}^{K} \lambda_k^* p_k z(i+1, k)$. Therefore, there exists at lease one $k \in \{1, 2, \cdots, K\}$ such that $z(i, k) < z(i+1, k)$.

Considering such a $k \in \{1, 2, \cdots, K\}$, we can show the following equalities

$$
\begin{aligned}
& z(i+1, k) - z(i, k) \\
=& \sum_{j=i+1}^{N} \left[ \left( \sum_{m=1}^{j} \frac{w_m}{x_m^*} \right)^{p_k-1} \cdot \frac{f_j}{F_{i+1}^c} \right] - \sum_{j=i}^{N} \left[ \left( \sum_{m=1}^{j} \frac{w_m}{x_m^*} \right)^{p_k-1} \cdot \frac{f_j}{F_i^c} \right] \\
\geq& \frac{f_i}{F_{i+1}^c \cdot F_i^c} \left\{ \sum_{j=i+1}^{N} \left( \sum_{m=1}^{i} \frac{w_m}{x_m^*} \right)^{p_k-1} \cdot f_j - \sum_{m=1}^{i} \left( \frac{w_m}{x_m^*} \right)^{p_k-1} \cdot F_{i+1}^c \right\} \\
=& \frac{f_i}{F_{i+1}^c \cdot F_i^c} \left( \sum_{m=1}^{i} \frac{w_m}{x_m^*} \right)^{p_k-1} \left( \sum_{j=i+1}^{N} f_j - F_{i+1}^c \right) = 0 \ ,
\end{aligned}
$$

which contradicts with the inequality $z(i+1, k) < z(i, k)$ derived based on the assumption of $x_{i+1}^* < x_i^*$. Finally, we note that it can be easily shown from (13) that $e(x_i)x_i^2$ is the same for all $i = 1, 2, \cdots, N$ and hence, homogeneous speeds are optimal and can be easily derived from the $L_1$ latency constraint. Therefore, Proposition 1 is proved.

## C  Proof of Theorem 2

For convenience, we restate the theorem.

**Theorem 2.** *Given $e(x) = b \cdot x^{\alpha-1}$ and one $L_p$ latency constraint, then the degree of heterogeneity $\frac{x_N^*}{x_1^*}$ increases with increasing $p$ for $p \geq 1$.*

We first show the following lemma.

**Lemma 1.** *Suppose that $e(x) = b \cdot x^{\alpha-1}$, there is only one $L_p$ latency norm constraint, and $\mathbf{x}^* = (x_1^*, x_2^*, \cdots, x_N^*)$ are the optimal core speeds minimizing the average energy consumption subject to the latency constraint $\tilde{D}(p)$. If the latency constraint becomes $\theta \tilde{D}(p)$, then the optimal core speeds are $\frac{\mathbf{x}^*}{\theta} = (x_1^*, x_2^*, \cdots, x_N^*)/\theta$, for any $\theta > 0$.*

*Proof.* We prove Lemma 1 based on the equivalent problem **P1'**. The first-order derivative of the energy function $e(x) = bx^{\alpha-1}$ is $e'(x) = (\alpha-1)bx^{\alpha-2}$, where $b > 0$ and $\alpha \geq 2$. Given the (equivalent) latency constraint $D_p = [\tilde{D}(p)]^p$ where $p \geq 1$, we denote the unique optimal dual variable and primal variables as $\lambda^*$ and $\mathbf{x}^* = (x_1^*, x_2^*, \cdots, x_N^*)$, respectively. Next, we write the Karush Kuhn Tucker (KKT) conditions [9] as follows

$$
\left[ \sum_{i=1}^{N} \left( \sum_{j=1}^{i} \frac{w_j}{x_j^*} \right)^p \cdot f_i \right] - D_p \leq 0, \tag{17}
$$

$$
\lambda^* \geq 0, \tag{18}
$$

$$
\lambda^* \cdot \left\{ \left[ \sum_{i=1}^{N} \left( \sum_{j=1}^{i} \frac{w_j}{x_j^*} \right)^p \cdot f_i \right] - D_p \right\} = 0, \tag{19}
$$

$$
B \cdot (x_i^*)^{\alpha-2} - \frac{\lambda^* p}{(x_i^*)^2} \cdot \sum_{j=i}^{N} \left( \sum_{m=1}^{j} \frac{w_m}{x_m} \right)^{p-1} f_j = 0, \tag{20}
$$

$$
\forall \, i = 1, 2, \cdots, N,
$$

where (17) and (18) are primal and dual feasibility constraints, (19) is the complementary slackness condition, and (20) is the zero gradient condition in which $B = b(\alpha - 1)(1 - F_{i-1})$. In fact, (17) takes the equality due to $\lambda^* > 0$ and the complementary condition. Suppose now that we have a new latency constraint $\theta \tilde{D}(p)$, which results in a new equivalent latency constraint of $\theta^p D_p = [\theta \tilde{D}(p)]^p$. Next, we can construct a new set of primal variables as $\frac{\mathbf{x}^*}{\theta} = \frac{(x_1^*, x_2^*, \cdots, x_N^*)}{\theta}$ and a new dual variable as $\frac{\lambda^*}{\theta^{\alpha+p-1}}$. Then, by plugging into (17)–(20), it can be seen that the new primal variables and dual variables satisfy the KKT conditions. Therefore, based on the strict convexity of the problem **P1'**, it follows that $\frac{\mathbf{x}^*}{\theta} = \frac{(x_1^*, x_2^*, \cdots, x_N^*)}{\theta}$ solves **P1'** and hence **P1**, given the new latency constraint $\theta \tilde{D}(p)$, proving Lemma 1. ∎

Lemma 1 gives a somewhat unexpected result. If only one $L_p$ norm constraint is imposed, the degree of heterogeneity under the optimal schedule only depends on the value of $p$, irrespective of the actual delay norm constraint $\tilde{D}(p)$. The intuition underlying Lemma 1 is that when core speeds can be chosen arbitrarily, we can scale them to satisfy latency constraints and hence the degree of heterogeneity is constant once the value of $p$ is fixed. In practice, of course, the fastest core speed is not arbitrary and we assume that the software is designed to satisfy the demand of the longest job on the fastest processor.

Given only one latency constraint specified by a certain $L_p$ norm, we denote the optimal speed of core $i$ by $x_i^*$, for $i = 1, 2, \cdots, N$. Then, based on (15), we show the following equalities satisfied by $\mathbf{x}^* = (x_1^*, x_2^*, \cdots, x_N^*)$

$$F_i^c \cdot (x_i^*)^\alpha = \lambda^* p \sum_{j=i}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right], \ \forall i = 1, 2, \cdots, N,$$

where $F_i^c = 1 - F_{i-1}$ and $\lambda^*$ is the Lagrangian multiplier associated with the latency constraint. Then, with simple manipulations, we obtain the following equalities

$$\frac{F_i^c}{F_j^c} \left( \frac{x_i^*}{x_j^*} \right)^\alpha = \frac{\sum_{k=i}^N \left[ \left( \sum_{m=1}^k \frac{w_m}{x_m^*} \right)^{p-1} f_k \right]}{\sum_{k=j}^N \left[ \left( \sum_{m=1}^k \frac{w_m}{x_m^*} \right)^{p-1} f_k \right]}, \ \forall i, j = 1, 2, \cdots, N, \tag{21}$$

which are satisfied by the optimal core speeds $\mathbf{x}^* = (x_1^*, x_2^*, \cdots, x_N^*)$ that minimize the average energy consumption subject to a latency constraint specified by a $L_p$ norm.

Lemma 1 states that if the value of $p \geq 1$ is fixed, then the ratio of two optimal core speeds $\frac{x_j^*}{x_i^*}$, for $i, j = 1, 2, \cdots, N$, is constant regardless of the actual latency constraint. In other words, $\frac{x_j^*}{x_i^*}$ is solely determined by the value of $p$. We denote the ratio of $\frac{x_i^*}{x_{i-1}^*} = \beta_i(p)$, for $i = 2, 3, \cdots, N$, which stresses the dependency of the ratio on $p$. Hence, the degree of heterogeneity given a $L_p$ norm latency constraint is given by

$$\frac{x_N^*}{x_1^*} = \prod_{i=2}^N \beta_i(p), \tag{22}$$

which is independent of the actual latency constraint provided that $p$ is fixed.

Let us consider two types of latency constraints specified by $p_0$ and $p$ where $p_0 > p \geq 1$, and denote the corresponding optimal speeds of core $i$ as $x_i^*$ and $x_{i,0}^*$, respectively, for $i = 1, 2, \cdots, N$. In what follows, we shall prove $\prod_{i=2}^N \beta_i(p_0) > \prod_{i=2}^N \beta_i(p)$. First, we show $\beta_2(p_0) > \beta_2(p)$. Given $p$, we can choose a certain latency constraint $\tilde{D}_p$ such that the optimal core speed $x_1^* = w_1$ (i.e., time unit to process the first bin of CPU cycles). Thus, from (21), it follows that

$$
\begin{aligned}
F_2^c \left[ \beta_2(p) \right]^\alpha &= 1 - \frac{\left( \frac{w_1}{x_1^*} \right)^{p-1} f_1}{\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right]} \\
&= 1 - \frac{f_1}{\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right]}.
\end{aligned}
\tag{23}
$$

Similarly, for $p_0$, we can choose a certain latency constraint $\tilde{D}_{p_0}$ such that $x_{1,0}^* = w_1$, and the following equality holds

$$
F_2^c \left[ \beta_2(p_0) \right]^\alpha = 1 - \frac{f_1}{\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_{m,0}^*} \right)^{p_0-1} f_j \right]}.
\tag{24}
$$

Suppose that $\beta_2(p_0) \leq \beta_2(p)$, i.e., $x_{2,0}^* \leq x_2^*$ when $x_{1,0}^* = x_1^* = w_1$. Thus, by checking (23) and (24), we know that

$$
\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_{m,0}^*} \right)^{p_0-1} f_j \right] \leq \sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right].
\tag{25}
$$

Next, let us look at $\frac{x_3^*}{x_1^*}$ and $\frac{x_{3,0}^*}{x_{1,0}^*}$. Following (21), we know that

$$
F_3^c \left( \frac{x_3^*}{x_1^*} \right)^\alpha = 1 - \frac{f_1 + \left( 1 + \frac{w_2}{x_2^*} \right)^{p-1} \cdot f_2}{\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right]},
\tag{26}
$$

$$
F_3^c \left( \frac{x_{3,0}^*}{x_{1,0}^*} \right)^\alpha = 1 - \frac{f_1 + \left( 1 + \frac{w_2}{x_{2,0}^*} \right)^{p_0-1} \cdot f_2}{\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_{m,0}^*} \right)^{p_0-1} f_j \right]}.
\tag{27}
$$

Then, based on (25) and the assumption of $x_{2,0}^* \leq x_2^*$, we obtain $x_{3,0}^* \leq x_3^*$. Repeating this process, we can show $x_{i,0}^* < x_i^*$, for $i = 3, 4, \cdots, N$. Thus, we obtain

$$
\begin{aligned}
\sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_{m,0}^*} \right)^{p_0-1} f_j \right] &> \sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p_0-1} f_j \right] \\
&> \sum_{j=1}^N \left[ \left( \sum_{m=1}^j \frac{w_m}{x_m^*} \right)^{p-1} f_j \right],
\end{aligned}
\tag{28}
$$

which contradicts (25). Thus, $\beta_2(p_0) \leq \beta_2(p)$ does not hold, and $\beta_2(p_0)$ must be greater than $\beta_2(p)$ when $p_0 > p \geq 1$.

Next, we turn to $\beta_3(p)$ and $\beta_3(p_0)$ and suppose $\beta_3(p)\beta_2(p) \geq \beta_3(p_0)\beta_2(p_0)$. As it has been shown $\beta_2(p) < \beta_2(p_0)$, we can establish the inequality $\beta_3(p) > \beta_3(p_0)$. By appropriately choosing two latency constraints without changing the core speed ratios (as shown by Lemma 1), we can have two optimal speeds of core 2 are $x_2^*$ and $x_{2,0}^*$, for $p$ and $p_0$, respectively, such that $\frac{w_1}{x_1^*} + \frac{w_2}{x_2^*} = \frac{w_1}{x_{1,0}^*} + \frac{w_2}{x_{2,0}^*} = 1$. Equivalently, we have $\frac{w_1\beta_2(p)+w_2}{x_2^*} = \frac{w_1\beta_2(p_0)+w_2}{x_{2,0}^*}$, which leads to $\frac{x_2^*}{x_{2,0}^*} = \frac{w_1\beta_2(p)+w_2}{w_1\beta_2(p_0)+w_2}$. Then, we obtain $\frac{x_3^*}{x_{3,0}^*} = \frac{x_2^*\beta_3(p)}{x_{2,0}^*\beta_3(p_0)} = \frac{w_1\beta_2(p)\beta_3(p)+w_2\beta_3(p)}{w_1\beta_2(p_0)\beta_3(p_0)+w_2\beta_3(p_0)} > 1$, i.e., $x_3^* > x_{3,0}^*$. Similar with (23) and based on $\frac{w_1}{x_1^*} + \frac{w_2}{x_2^*} = \frac{w_1}{x_{1,0}^*} + \frac{w_2}{x_{2,0}^*} = 1$, we can show that

$$\frac{F_3^c}{F_2^c}[\beta_3(p)]^\alpha = 1 - \frac{f_2}{\sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_m^*}\right)^{p-1} f_j\right]} \tag{29}$$

$$\frac{F_3^c}{F_2^c}[\beta_3(p_0)]^\alpha = 1 - \frac{f_2}{\sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_{m,0}^*}\right)^{p_0-1} f_j\right]} \tag{30}$$

Then, from $\beta_3(p) > \beta_3(p_0)$, we have

$$\sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_{m,0}^*}\right)^{p_0-1} f_j\right] < \sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_m^*}\right)^{p-1} f_j\right]. \tag{31}$$

Next, let us look at $\frac{x_4^*}{x_2^*}$ and $\frac{x_{4,0}^*}{x_{2,0}^*}$. Following (21), we obtain

$$\frac{F_4^c}{F_2^c}\left(\frac{x_4^*}{x_2^*}\right)^\alpha = 1 - \frac{f_2 + \left(1 + \frac{w_3}{x_3^*}\right)^{p-1}\cdot f_3}{\sum_{j=1}^N\left[\left(\sum_{m=2}^j \frac{w_m}{x_m^*}\right)^{p-1} f_j\right]}, \tag{32}$$

$$\frac{F_4^c}{F_2^c}\left(\frac{x_{4,0}^*}{x_{2,0}^*}\right)^\alpha = 1 - \frac{f_2 + \left(1 + \frac{w_3}{x_{3,0}^*}\right)^{p_0-1}\cdot f_3}{\sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_{m,0}^*}\right)^{p_0-1} f_j\right]}. \tag{33}$$

Then, from $x_3^* > x_{3,0}^*$, we obtain the following equalities

$$\left(1 + \frac{w_3}{x_3^*}\right)^{p-1} \leq \left(1 + \frac{w_3}{x_3^*}\right)^{p_0-1} \leq \left(1 + \frac{w_3}{x_{3,0}^*}\right)^{p_0-1}. \tag{34}$$

Then, from (31)–(34), we can see that $x_{4,0}^* \leq x_4^*$. Repeating this process, we can show that $x_{i,0}^* < x_i^*$, for $i = 4, 5, \cdots, N$. Hence, by combining this result with the inequality of $x_3^* \geq x_{3,0}^*$, we obtain

$$\sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_{m,0}^*}\right)^{p_0-1} f_j\right] > \sum_{j=1}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_m^*}\right)^{p_0-1} f_j\right]$$
$$> \sum_{j=2}^N\left[\left(\sum_{m=1}^j \frac{w_m}{x_m^*}\right)^{p-1} f_j\right], \tag{35}$$

which contradicts (31). Thus, $\beta_3(p)\beta_2(p) \geq \beta_3(p_0)\beta_2(p_0)$ does not hold, and $\beta_3(p_0)\beta_2(p_0)$ must be greater than $\beta_3(p)\beta_2(p)$ when $p_0 > p \geq 1$.

By continuing the above procedure, we can prove that $\prod_{j=2}^{i} \beta_i(p_0) > \prod_{j=2}^{i} \beta_i(p) \geq 1$ when $p_0 > p \geq 1$, for all $i = 2, 3, \cdots, N$. Therefore, it follows immediately from (22) that the degree of heterogeneity increases as $p$ becomes larger, proving Theorem 2.

## D   Proof of Theorem 3

For convenience, we restate the theorem.

**Theorem 3.** *Given $e(x) = b \cdot x^{\alpha-1}$ and $K$ $L_p$ latency constraints specified by $1 \leq p_1 \leq p_2 \leq \cdots \leq p_K \leq \infty$, then the degree of heterogeneity $\frac{x_N^*}{x_1^*}$ satisfies:*

$$1 \leq \frac{x_N^*}{x_1^*} \leq \left(\frac{1}{f_N}\right)^{\frac{1}{\alpha}} \tag{36}$$

*where $f_N$ is the probability that a job has the maximum service demand of $\tilde{w}_N$.*

Denote the optimal dual variable and core speeds by $\lambda^* = (\lambda_1^*, \lambda_2^*, \cdots, \lambda_K^*)$ and $\mathbf{x}^* = (x_1^*, x_2^*, \cdots, x_N^*)$, respectively. Based on (15), we obtain

$$\left(\frac{x_N^*}{x_1^*}\right)^{\alpha} = \frac{\sum_{k=1}^{K} \left[\lambda_k^* p_k \left(\sum_{m=1}^{N} \frac{w_m}{x_m^*}\right)^{p_k-1}\right]}{\sum_{j=i}^{N} \left\{\sum_{k=1}^{K} \left[\lambda_k^* p_k \left(\sum_{m=1}^{j} \frac{w_m}{x_m^*}\right)^{p_k-1}\right] f_j\right\}} \tag{37}$$

$$\leq \frac{\lambda_K^* p_K \left(\sum_{m=1}^{N} \frac{w_m}{x_m^*}\right)^{p_K-1}}{\sum_{j=i}^{N} \lambda_K^* p_K \left(\sum_{m=1}^{j} \frac{w_m}{x_m^*}\right)^{p_K-1} \cdot f_j} \tag{38}$$

$$\leq \lim_{p_K \to \infty} \frac{\left(\sum_{m=1}^{N} \frac{w_m}{x_m^*}\right)^{p_K-1}}{\sum_{j=i}^{N} \left(\sum_{m=1}^{j} \frac{w_m}{x_m^*}\right)^{p_K-1} \cdot f_j} = \frac{1}{f_N} \tag{39}$$

where the inequalities (38) and (39) result from Theorem 2 which states that the degree of heterogeneity increases with $p$. Note that the inequalities (38) and (39) become equalities when the maximum latency constraint specified by the $L_\infty$ norm is dominant. Thus, after a simple mathematical manipulation, we obtain $\frac{x_N^*}{x_1^*} \leq \left(\frac{1}{f_N}\right)^{\frac{1}{\alpha}}$. The inequality of $\frac{x_N^*}{x_1^*} \geq 1$ directly follows Theorem 1, and it becomes an equality when the average latency constraint specified by the $L_1$ norm is dominant. Thus, Theorem 3 is proved.

## E   Algorithm to solve P2

Before presenting the branch-and-bound algorithm, we first consider a practical constraint that the core speeds are both upper and lower bounded and then

recast the problem of **P1** into the following

$$\mathbf{P3}: \ \min_{\mathbf{x}} \sum_{i=1}^{N} \{[1 - F_{i-1}] \cdot e(x_i) \cdot w_i\} \tag{40}$$

$$s.t., \ \text{Constraint (4) and } \mathbf{x}_{\min} \preceq \mathbf{x} \preceq \mathbf{x}_{\max}, \tag{41}$$

where $\mathbf{x}_{\min} = x_{\min} \cdot \mathbf{I}^{(1 \times N)}$ and $\mathbf{x}_{\max} = x_{\max} \cdot \mathbf{I}^{(1 \times N)}$ are minimum and maximum speed constraints, respectively, and $\mathbf{I}^{(1 \times N)}$ is a $1 \times N$ unit vector.

Similar as **P1**, **P3** is also a convex optimization problem. We can solve it using a primal-dual approach with a small modification on Algorithm 1 (shown in Appendix A). The only difference from Algorithm 1 is at Line 6, where we derive the optimal $\tilde{x}_i$ that minimizes the Lagrangian. Specifically, after solving $\tilde{x}_i$ based on (13), we need to add another step to apply the lower and upper bounds on $\tilde{x}_i$, i.e., $\tilde{x}_i = [\tilde{x}_i]_{x_{\min}}^{x_{\max}}$. We refer to this algorithm that solves **P3** as Algorithm 2, which is omitted for brevity. Next, we present the branch-and-bound algorithm in the following four steps.

*Decomposition* We first decompose **P2** into $M$ sub-problems, indexed by $\mathbf{P2}_1, \mathbf{P2}_2, \cdots, \mathbf{P2}_M$. Each sub-problem $\mathbf{P2}_m$ is expressed as follows:

$$\mathbf{P2}_m: \ \min_{\mathbf{x}} \sum_{i=1}^{N} \{[1 - F_{i-1}] \cdot e(x_i) \cdot w_i\} \tag{42}$$

$$s.t., \ \text{Constraint (4)} \tag{43}$$

$$x_i \in \{s_1, s_2, \cdots, s_M\}, i = 2, \cdots, N \tag{44}$$

$$x_1 = s_m, \tag{45}$$

where we fix $x_1 = s_m$ and minimize the average energy consumption over $\mathbf{x} \backslash \{x_1\}$.[3] After solving all the $M$ sub-problems, we can select one sub-problem (say, $\mathbf{P2}_m$) that yields the minimum average energy consumption and then, combined with $x_1 = s_m$, we obtain the optimal core speeds $\mathbf{x}^*$. Each sub-problem itself is an combinatorial problem and can be further decomposed into $M$ smaller problems by fixing another core speed. Thus, the original problem can be solved recursively, which serves as the basis for applying the branch-and-bound technique.

*Lower and upper bounds* By replacing (43) with $x_i \in [s_1, s_M]$ for $i = 1, 2, \cdots, N$ and solving the problem **P3**, the constraint (43) is relaxed and the resulting average energy consumption is a lower bound on that of problem **P2** (i.e., the minimum average energy consumption in **P2** is no less than that in **P3**).

To find an upper bound on the minimum average energy consumption in **P2**, we propose a greedy algorithm, as described in Algorithm 3, which never outperforms the optimal solution to **P2** in terms of the average energy consumption. In the greedy algorithm, all the core speeds are chosen to be the minimum value $s_1$ initially. If not all the latency norm constraints are satisfied, we greedily increase the core speed such that the average energy increase is minimum (i.e., Lines 3–6 in Algorithm 3). Repeat this process until all the core speeds have reached their maximum $s_M$ or all the $L_p$ latency norm constraints are satisfied.

---

[3] We can also fix any other core speed other than $x_1$.

---

**Algorithm 3** Greedy

---

1: Initialize $\mathbf{x} = s_1 \cdot \mathbf{I}^{(1 \times N)}$
2: **while** $\mathbf{x} \neq s_M \cdot \mathbf{I}^{(1 \times N)}$ **and** Constraint (5) is not satisfied **do**
3:     $\Omega \leftarrow \{i \,|\, i = 1, 2, \cdots, N, x_i \neq s_M\}$
4:     $\Delta x_i \leftarrow \arg \min_{s \in \{s_1, s_2, \cdots, s_M\}} (s > x_i), \forall i \in \Omega$
5:     $i = \arg \min_{i \in \Omega} \{[e(\Delta x_i) - e(x_i)] \, (1 - F_{i-1}) \cdot w_i\}$
6:     $x_i \leftarrow \Delta x_i$
7: **end while**
8: **return** $\mathbf{x}^* = \mathbf{x}$

---

Next, we define the following notations that facilitate the description of our branch-and-bound algorithm.

**Definition 1.** *$LB(\mathcal{X})$ is the minimum average energy consumption obtained by solving **P3** with the constraint $\mathcal{X}$ as its additional input. $UB(\mathcal{X})$ is the minimum average energy consumption obtained by using the proposed greedy algorithm (i.e., Algorithm 3) with the additional constraint $\mathcal{X}$ as its additional input.*

We explain Definition 1 using an example. If $\mathcal{X} = \{x_1 = s_m\}$, we compute $LB(\mathcal{X})$ by solving **P3** with an additional constraint of $x_1 = s_m$. The variation of **P3** with additional constraints specified by $\mathcal{X}$ is still convex and can be efficiently solved using our proposed primal-dual algorithm. Similarly, we compute $UB(\mathcal{X})$ using the proposed greedy algorithm with an additional constraint of $x_1 = s_m$. We use $LB(\varnothing)$ and $UB(\varnothing)$ to represent the minimum average energy consumptions obtained by solving the original problem **P3** and by using the greedy algorithm without additional constraints, respectively.

*Fixing rule* A core component of branch-and-bound algorithms is the "fixing" rule, which determines the next decision variable to be fixed. We define the "fixing" rule as follows.

**Definition 2.** *$next(\mathcal{X})$ is the index of the core speed that the proposed greedy algorithm uses given the constraint set $\mathcal{X}$ as the input.*

In essence, we select and fix the core speed which, if increased to the next faster speed out of the available speeds $\{s_1, s_2, \cdots, s_M\}$, results in the minimum increase in average energy.

*Algorithm* We describe our branch-and-bound algorithm in Algorithm 4. The parameter *IterateMax* is the maximum number of iterations selected based on the desired accuracy and the problem scale. The algorithm generates an $M$-ary tree, where each node represents a constraint set and all the leaf nodes are stored in the set $\mathcal{S}$. The algorithm begins with an empty constraint set $\mathcal{X}_0 = \varnothing$ as the parent node of the tree. In each iteration, we choose a leaf node and split it into $M$ new leaf nodes, each of which represents a new constraint set with an additional core speed fixed to be one of the permissible values in $\{s_1, s_2, \cdots, s_M\}$. In the splitting process (Lines 4–7), we split the node that corresponds to the constraint

---
**Algorithm 4** SEM-M
---
1: Initialize: $i \leftarrow 0$, $\mathcal{X}_0 \leftarrow \varnothing$, set of leaves of a single-node tree $\mathcal{S} \leftarrow \mathcal{X}_0$
2: Compute lower and upper bounds: $L_0 = LB(\mathcal{X}_0)$ and $U_0 = UB(\mathcal{X}_0)$
3: **while** $U_i - L_i > \epsilon$ **or** $i < IterateMax$ **do**
4:     Choose the splitting node: $\mathcal{X}^* = \arg\min_{\mathcal{X} \in \mathcal{S}} LB(\mathcal{X})$
5:     Choose the core speed to fix: $n = next(\mathcal{X}^*)$
6:     Generate $M$ new constraint sets:
       $\mathcal{X}_{i+1}^1 = \mathcal{X}^* \cup \{x_n = s_1\}, \cdots, \mathcal{X}_{i+1}^M = \mathcal{X}^* \cup \{x_n = s_M\}$
7:     Update the set of leaves:
       $\mathcal{S} \leftarrow (\mathcal{S} \backslash \{\mathcal{X}^*\}) \cup \{\mathcal{X}_{i+1}^1\} \cup \cdots \cup \{\mathcal{X}_{i+1}^M\}$
8:     Compute upper and lower bounds for $M$ new constraint sets:
       $LB(\mathcal{X}_{i+1}^1), \cdots, LB(\mathcal{X}_{i+1}^M), UB(\mathcal{X}_{i+1}^1), \cdots, UB(\mathcal{X}_{i+1}^M)$
9:     Update global upper and lower bounds:
       $L_{i+1} = \min_{\mathcal{X} \in \mathcal{S}} LB(\mathcal{X})$ and $U_{i+1} = \max_{\mathcal{X} \in \mathcal{S}} UB(\mathcal{X})$
10:     $i \leftarrow i + 1$
11: **end while**
12: Choose the best constraint set thus far:
    $\bar{\mathcal{X}} = \min_{\mathcal{X} \in \mathcal{S}} UB(\mathcal{X})$
13: **return** $\mathbf{x}^*$ achieved by the greedy algorithm (i.e., Algorithm 3) with $\bar{\mathcal{X}}$ as the constraint
---

set resulting in the minimum average energy consumption (obtained by solving **P3** with an additional constraint specified by the node to be split). The reason behind our splitting process is that after splitting this node, the global lower bound will likely be increased, while splitting any other node keeps the global lower bound unchanged and hence the algorithm does not shrink the gap between the global upper and lower bounds. Besides the maximum number of iterations, another termination criterion is the difference between the global upper and lower bounds. Specifically, if $U_i - L_i$ is no greater than a sufficiently small positive number $\epsilon$, it is guaranteed that the solution obtained using the greedy algorithm with an appropriate constraint set is close-to-optimal. Therefore, by increasing *IterateMax* and using a sufficiently small positive number $\epsilon$, SEM yields an arbitrarily close-to-optimal solution, while the global optimality is achieved at the expense of increasing the computation cost (which, however, is not a major concern in our context, as we only recompute the optimal solution when the service demand distribution or latency requirement changes).