

Scalably Scheduling Power-Heterogeneous Processors

Anupam Gupta^{1*}, Ravishankar Krishnaswamy^{*1}, and Kirk Pruhs^{2**}

¹ Computer Science Dept., Carnegie Mellon University, Pittsburgh, PA 15213, USA.

² Computer Science Dept., University of Pittsburgh, Pittsburgh, PA 15260, USA.

Abstract. We show that a natural online algorithm for scheduling jobs on a heterogeneous multiprocessor, with arbitrary power functions, is scalable for the objective function of weighted flow plus energy.

1 Introduction

Many prominent computer architects believe that architectures consisting of heterogeneous processors/cores, such as the STI Cell processor, will be the dominant architectural design in the future [8, 13, 12, 17, 18]. The main advantage of a heterogeneous architecture, relative to an architecture of identical processors, is that it allows for the inclusion of processors whose design is specialized for particular types of jobs, and for jobs to be assigned to a processor best suited for that job. Most notably, it is envisioned that these heterogeneous architectures will consist of a small number of high-power high-performance processors for critical jobs, and a larger number of lower-power lower-performance processors for less critical jobs. Naturally, the lower-power processors would be more energy efficient in terms of the computation performed per unit of energy expended, and would generate less heat per unit of computation. For a given area and power budget, heterogeneous designs can give significantly better performance for standard workloads [8, 17]; Emulations in [12] suggest a figure of 40% better performance, and emulations in [18] suggest a figure of 67% better performance. Moreover, even processors that were designed to be homogeneous, are increasingly likely to be heterogeneous at run time [8]: the dominant underlying cause is the increasing variability in the fabrication process as the feature size is scaled down (although run time faults will also play a role). Since manufacturing yields would be unacceptably low if every processor/core was required to be perfect, and since there would be significant performance loss from derating the entire chip to the functioning of the least functional processor (which is what would be required in order to attain processor homogeneity), some processor heterogeneity seems inevitable in chips with many processors/cores.

* Supported in part by NSF award CCF-0729022 and an Alfred P. Sloan Fellowship.

** Supported in part by NSF grants CNS-0325353, IIS-0534531, and CCF-0830558, and an IBM Faculty Award.

The position paper [8] identifies three fundamental challenges in scheduling heterogeneous multiprocessors: (1) the OS must discover the status of each processor, (2) the OS must discover the resource demand of each job, and (3) given this information about processors and jobs, the OS must match jobs to processors as well as possible. In this paper, we address this third fundamental challenge. In particular, we assume that different jobs are of differing importance, and we study how to assign these jobs to processors of varying power and varying energy efficiency, so as to achieve the best possible trade-off between energy and performance.

Formally, we assume that a collection of jobs arrive in an online fashion over time. When a job j arrives in the system, the system is able to discover a *size* $p_j \in \mathbb{R}_{>0}$, as well as a *importance/weight* $w_j \in \mathbb{R}_{>0}$, for that job. The importance w_j specifies an upper bound on the amount of energy that the system is allowed to invest in running j to reduce j 's flow by one unit of time (assuming that this energy investment in j doesn't decrease the flow of other jobs)—hence jobs with high weight are more important, since higher investments of energy are permissible to justify a fixed reduction in flow. Furthermore, we assume that the system knows the allowable speeds for each processor, and the system also knows the power used when each processor is run at its set of allowable speeds. We make no real restrictions on the allowable speeds, or on the power used for these speeds.³ The online scheduler has three component policies:

Job Selection: Determines which job to run on each processor at any time.

Speed Scaling: Determines the speed of each processor at each time.

Assignment: When a new job arrives, it determines the processor to which this new job is assigned.

The objective we consider is that of *weighted flow plus energy*. The rationale for this objective function is that the optimal schedule under this objective gives the best possible weighted flow for the energy invested, and increasing the energy investment will not lead to a corresponding reduction in weighted flow (intuitively, it is not possible to speed up a collection of jobs with an investment of energy proportional to these jobs' importance).

We consider the following natural online algorithm that essentially adopts the job selection and speed scaling algorithms from the uniprocessor algorithm in [5], and then greedily assigns the jobs based on these policies.

Job Selection: Highest Density First (HDF)

Speed Scaling: The speed is set so that the power is the fractional weight of the unfinished jobs.

Assignment: A new job is assigned to the processor that results in the least increase in the projected future weighted flow, assuming the adopted speed scaling and job selection policies, and ignoring the possibility of jobs arriving in the future.

³ So the processors may or may not be speed scalable, the speeds may be continuous or discrete or a mixture, the static power may or may not be negligible, the dynamic power may or may not satisfy the cube root rule, etc.

Our main result is then:

Theorem 1. *This online algorithm is scalable for scheduling jobs on a heterogeneous multiprocessor with arbitrary power functions to minimize the objective function of weighted flow plus energy.*

In this context, *scalable* means that if the adversary can run processor i at speed s and power $P(s)$, the online algorithm is allowed to run the processor at speed $(1 + \epsilon)s$ and power $P(s)$, and then for all inputs, the online cost is bounded by $O(f(\epsilon))$ times the optimal cost. Intuitively, a scalable algorithm can handle almost the same load as optimal; for further elaboration, see [20, 19]. Theorem 1 extends theorems showing similar results for weighted flow plus energy on a uniprocessor [5, 2], and for weighted flow on a multiprocessor without power considerations [9]. As scheduling on identical processors with the objective of total flow, and scheduling on a uniprocessor with the objective of weighted flow, are special cases of our problem, constant competitiveness is not possible without some sort of resource augmentation [16, 3].

Our analysis is an amortized local-competitiveness argument. As is usually the case with such arguments, the main technical hurdle is to discover the “right” potential function. The most natural strawman potential function to try is the sum over all processors of the single processor potential function used in [5]. While one can prove constant competitiveness with this potential in some special cases (e.g. where for each processor the allowable speeds are the non-negative reals, and the power satisfies the cube-root rule), one can not prove constant competitiveness for general power functions with this potential function. The reason for this is that the uniprocessor potential function from [5] is not sufficiently accurate. Specifically, one can construct configurations where the adversary has finished all jobs, and where the potential is much higher than the remaining online cost. This did not mess up the analysis in [5] because to finish all these jobs by this time the adversary would have had to run very fast in the past, wasting a lot of energy, which could then be used to pay for this unnecessarily high potential. But since we consider multiple processors, the adversary may have no jobs left on a particular processor simply because it assigned these jobs to a different processor, and there may not be a corresponding unnecessarily high adversarial cost that can be used to pay for this unnecessarily high potential.

Thus, the main technical contribution in this paper is a seemingly more accurate potential function expressing the additional cost required to finish one collection of jobs compared to another collection of jobs. Our potential function is arguably more transparent than the one used in [5], and we expect that this potential function will find future application in the analysis of other power management algorithms.

In section 3, we show that a similar online algorithm is $O(1/\epsilon)$ -competitive with $(1 + \epsilon)$ -speedup for *unweighted* flow plus energy. We also remark that when the power functions $P_i(s)$ are restricted to be of the form s^{α_i} , our algorithms give a $O(\alpha^2)$ -competitive algorithm (with no resource augmentation needed) for the problem of minimizing weighted flow plus energy, and an $O(\alpha)$ -competitive algorithm for minimizing the unweighted flow plus energy, where $\alpha = \max_i \alpha_i$.

1.1 Related Results

Let us first consider previous work for the case of a single processor, with unbounded speed, and a polynomially bounded power function $P(s) = s^\alpha$. [21] gave an efficient offline algorithm to find the schedule that minimizes average flow subject to a constraint on the amount of energy used, in the case that jobs have unit work. [1] introduced the objective of flow plus energy and gave a constant competitive algorithm for this objective in the case of unit work jobs. [6] gave a constant competitive algorithm for the objective of weighted flow plus energy. The competitive ratio was improved by [15] for the unweighted case using a potential function specifically tailored to integer flow. [4] extended the results of [6] to the bounded speed model, and [10] gave a *nonclairvoyant* algorithm that is $O(1)$ -competitive.

Still for a single processor, dropping the assumptions of unbounded speed and polynomially-bounded power functions, [5] gave a 3-competitive algorithm for the objective of unweighted flow plus energy, and a 2-competitive algorithm for fractional weighted flow plus energy, both in the uniprocessor case for a large class of power functions. The former analysis was subsequently improved by [2] to show 2-competitiveness, along with a matching lower bound on the competitive ratio.

Now for multiple processors: [14] considered the setting of multiple *homogeneous* processors, where the allowable speeds range between zero and some upper bound, and the power function is polynomial in this range. They gave an algorithm that uses a variant of round-robin for the assignment policy, and job selection and speed scaling policies from [6], and showed that this algorithm is scalable for the objective of (unweighted) flow plus energy. Subsequently, [11] showed that a randomized machine selection algorithm is scalable for weighted flow plus energy (and even more general objective functions) in the setting of polynomial power functions. Both these algorithms provide non-migratory schedules and compare their costs with optimal solutions which could even be migratory. In comparison, as mentioned above, for the case of polynomial power functions, our techniques can give a deterministic constant-competitive online algorithm for non-migratory weighted flow time plus energy. (Details appear in the final version.)

In non-power-aware settings, the paper most relevant to this work is that of [9], which gives a scalable online algorithm for minimizing weighted flow on unrelated processors. Their setting is even more demanding, since they allow the processing requirement of the job to be processor dependent (which captures a type of heterogeneity that is orthogonal to the performance energy-efficiency heterogeneity that we consider in this paper). Our algorithm is based on the same general intuition as theirs: they assign each new job to the processor that would result in the least increment in future weighted flow (assuming HDF is used for job selection), and show that this online algorithm is scalable using an amortized local competitiveness argument. However, it is unclear how to directly extend their potential function to our power-aware setting; we had success only

in the case that each processor had allowable speed-power combinations lying in $\{(0, 0), (s_i, P_i)\}$.

1.2 Preliminaries

Scheduling Basics. We consider only non-migratory schedules, which means that no job can ever run on one processor, and later run on some other processor. In general, migration is undesirable as the overhead can be significant. We assume that preemption is allowed, that is, that jobs may be suspended, and restarted later from the point of suspension. It is clear that if preemption is not allowed, bounded competitiveness is not obtainable. The speed is the rate at which work is completed; a job j with size p_j run at a constant speed s completes in $\frac{p_j}{s}$ seconds. A job is completed when all of its work has been processed. The flow of a job is the completion time of the job minus the release time of the job. The weighted flow of a job is the weight of the job times the flow of the job. For a $t \geq r_j$, let $p_j(t)$ be the remaining unprocessed work on job j at time t . The fractional weight of job j at this time is $w_j \frac{p_j(t)}{p_j}$. The fractional weighted flow of a job is the integral over times between the job's release time and its completion time of its fractional weight at that time. The density of a job is its weight divided by its size. The job selection policy Highest Density First (HDF) always runs the job of highest density. The inverse density of a job is its size divided by its weight.

Power Functions. The power function for processor i is denoted by $P_i(s)$, and specifies the power used when processor is run at speed s . We essentially allow any reasonable power function. However, we do require the following minimal conditions on each power function, which we adopt from [5]. We assume that the allowable speeds are a countable collection of disjoint subintervals of $[0, \infty)$. We assume that all the intervals, except possibly the rightmost interval, are closed on both ends. The rightmost interval may be open on the right if the power $P_i(s)$ approaches infinity as the speed s approaches the rightmost endpoint of that interval. We assume that P_i is non-negative, and P_i is continuous and differentiable on all but countably many points. We assume that either there is a maximum allowable speed T , or that the limit inferior of $P_i(s)/s$ as s approaches infinity is not zero (if this condition doesn't hold then, then the optimal speed scaling policy is to run at infinite speed). Using transformations specified in [5], we may assume without loss of generality that the power functions satisfy the following properties: P is continuous and differentiable, $P(0) = 0$, P is strictly increasing, P is strictly convex, and P is unbounded. We use Q_i to denote P_i^{-1} ; i.e., $Q_i(y)$ gives us the speed that we can run processor i at, if we specify a limit of y .

Local Competitiveness and Potential Functions. Finally, let us quickly review amortized local competitiveness analysis on a single processor. Consider

an objective G . Let $G_A(t)$ be the increase in the objective in the schedule for algorithm A at time t . So when G is fractional weighted flow plus energy, $G_A(t)$ is $P_A^t + w_A^t$, where P_A^t is the power for A at time t and w_A^t is the fractional weight of the unfinished jobs for A at time t . Let OPT be the offline adversary that optimizes G . A is locally c -competitive if for all times t , if $G_A(t) \leq c \cdot G_{OPT}(t)$. To prove A is $(c+d)$ -competitive using an amortized local competitiveness argument, it suffices to give a potential function $\Phi(t)$ such that the following conditions hold (see for example [19]).

Boundary condition: Φ is zero before any job is released and Φ is non-negative after all jobs are finished.

Completion condition: Φ does not increase due to completions by either A or OPT.

Arrival condition: Φ does not increase more than $d \cdot OPT$ due to job arrivals.

Running condition: At any time t when no job arrives or is completed,

$$G_A(t) + \frac{d\Phi(t)}{dt} \leq c \cdot G_{OPT}(t) \quad (1)$$

The sufficiency of these conditions for proving $(c + d)$ -competitiveness follows from integrating them over time.

2 Weighted Flow

Our goal in this section is to prove Theorem 1. We first show that the on-line algorithm is $(1 + \epsilon)$ -speed $O(\frac{1}{\epsilon})$ -competitive for the objective of *fractional* weighted flow plus energy. Theorem 1 then follows since HDF is $(1 + \epsilon)$ -speed $O(\frac{1}{\epsilon})$ -competitive for fixed processor speeds [7] for the objective of (integer) weighted flow.

Let OPT be some optimal schedule minimizing fractional weighted flow. Let $w_{a,i}^t(q)$ denote the total fractional weight of jobs in processor i 's queue that have an inverse density of at least q . Let $w_{a,i}^t := w_{a,i}^t(0)$ be the total fractional weight of unfinished jobs in the queue. Let $w_a^t := \sum_i w_{a,i}^t$ be the total fractional weight of unfinished jobs in all queues. Let $w_{o,i}^t(q)$, $w_{o,i}^t$, and w_o^t be similarly defined for OPT. When the time instant being considered is clear, we drop the superscript of t from all variables.

We assume that once OPT has assigned a job to some processor, it runs the BCP algorithm [5] for job selection and speed scaling—i.e., it sets the speed of the i^{th} processor to $Q_i(w_{o,i})$, and hence the i^{th} processor uses power $W_{o,i}$, and uses HDF for job selection. We can make such an assumption because the results of [5] show that the fractional weighted flow plus energy of the schedule output by this algorithm is within a factor of two of optimal. Therefore, the only real difference between OPT and the online algorithm is the assignment policy.

2.1 The Assignment Policy

To better understand the online algorithm's assignment policy, define the “shadow potential” for processor i at time t to be

$$\hat{\Phi}_{a,i}(t) = \int_{q=0}^{\infty} \int_{x=0}^{w_{a,i}^t(q)} \frac{x}{Q_i(x)} dx dq \quad (2)$$

The shadow potential captures (up to a constant factor) the total fractional weighted flow to serve the current set of jobs if no jobs arrive in the future. Based on this, the online algorithm's assignment policy can alternatively be described as follows:

Assignment Policy. When a new job with size p_j and weight w_j arrives at time t , the assignment policy assigns it to a processor which would cause the smallest increase in the shadow potential; i.e. a processor minimizing

$$\begin{aligned} & \int_{q=0}^{d_j} \int_{x=0}^{w_{a,i}^t(q)+w_j} \frac{x}{Q_i(x)} dx dq - \int_{q=0}^{d_j} \int_{x=0}^{w_{a,i}^t(q)} \frac{x}{Q_i(x)} dx dq \\ &= \int_{q=0}^{d_j} \int_{x=w_{a,i}^t(q)}^{w_{a,i}^t(q)+w_j} \frac{x}{Q_i(x)} dx dq \end{aligned}$$

2.2 Amortized Local Competitiveness Analysis

We apply a local competitiveness argument as described in subsection 1.2. Because the online algorithm is using the BCP algorithm on each processor, the power for the online algorithm is $\sum_i P_i(Q_i(w_{a,i})) = w_a$. Thus $G_A = 2w_a$. Similarly, since OPT is using BCP on each processor $G_{\text{OPT}} = 2w_o$.

Defining the potential function For processor i , define the potential

$$\Phi_i(t) = \frac{2}{\epsilon} \int_{q=0}^{\infty} \int_{x=0}^{(w_{a,i}^t(q)-w_{o,i}^t(q))_+} \frac{x}{Q_i(x)} dx dq \quad (3)$$

Here $(\cdot)_+ = \max(\cdot, 0)$. The global potential is then defined to be $\Phi(t) = \sum_i \Phi_i(t)$. Firstly, we observe that the function $x/Q_i(x)$ is increasing and subadditive. Then, the following lemma will be useful subsequently, the proof of which will appear in the full version of the paper.

Lemma 1. *Let g be any increasing subadditive function with $g(0) \geq 0$, and $w_a, w_o, w_j \in \mathbb{R}_{\geq 0}$. Then,*

$$\int_{x=w_a}^{w_a+w_j} g(x) dx - \int_{x=(w_a-w_o-w_j)_+}^{(w_a-w_o)_+} g(x) dx \leq 2 \int_{x=0}^{w_j} g(w_o+x) dx$$

That the boundary and completion conditions are satisfied are obvious. In Lemma 2 we prove that the arrival condition holds, and in Lemma 3 we prove that the running condition holds.

Lemma 2. *The arrival condition holds with $d = \frac{4}{\epsilon}$.*

Proof. Consider a new job j with processing time p_j , weight w_j and inverse density $d_j = p_j/w_j$, which the algorithm assigns to processor 1 while the optimal solution assigns it to processor 2. Observe that $\int_{q=0}^{d_j} \int_{x=w_{o,2}(q)}^{w_{o,2}(q)+w_j} \frac{x}{Q_2(x)} dx dq$ is the increase in OPT's fractional weighted flow due to this new job j . Thus our goal is to prove that the increase in the potential due to job j 's arrival is at most this amount. The change in the potential $\Delta\Phi$ is:

$$\frac{2}{\epsilon} \int_{q=0}^{d_j} \left(\int_{x=(w_{a,1}(q)-w_{o,1}(q))+}^{(w_{a,1}(q)-w_{o,1}(q))+w_j} \frac{x}{Q_1(x)} dx - \int_{x=(w_{a,2}(q)-w_{o,2}(q)-w_j)+}^{(w_{a,2}(q)-w_{o,2}(q))+} \frac{x}{Q_2(x)} dx \right) dq$$

Now, since $x/Q_1(x)$ is an increasing function we have that

$$\int_{x=(w_{a,1}(q)-w_{o,1}(q))+}^{(w_{a,1}(q)-w_{o,1}(q))+w_j} \frac{x}{Q_1(x)} dx \leq \int_{x=w_{a,1}(q)}^{w_{a,1}(q)+w_j} \frac{x}{Q_1(x)} dx$$

and hence the change of potential can be bounded by

$$\frac{2}{\epsilon} \int_{q=0}^{d_j} \left(\int_{x=w_{a,1}(q)}^{w_{a,1}(q)+w_j} \frac{x}{Q_1(x)} dx - \int_{x=(w_{a,2}(q)-w_{o,2}(q)-w_j)+}^{(w_{a,2}(q)-w_{o,2}(q))+} \frac{x}{Q_2(x)} dx \right) dq$$

Since we assigned the job to processor 1, we know that

$$\int_{q=0}^{d_j} \int_{x=w_{a,1}(q)}^{w_{a,1}(q)+w_j} \frac{x}{Q_1(x)} dx dq \leq \int_{q=0}^{d_j} \int_{x=w_{a,2}(q)}^{w_{a,2}(q)+w_j} \frac{x}{Q_2(x)} dx dq$$

Therefore, the change in potential is at most

$$\Delta\Phi \leq \frac{2}{\epsilon} \int_{q=0}^{d_j} \left(\int_{x=w_{a,2}(q)}^{w_{a,2}(q)+w_j} \frac{x}{Q_2(x)} dx - \int_{x=(w_{a,2}(q)-w_{o,2}(q)-w_j)+}^{(w_{a,2}(q)-w_{o,2}(q))+} \frac{x}{Q_2(x)} dx \right) dq$$

Applying Lemma 1, we get:

$$\Delta\Phi \leq \left(2 \cdot \frac{2}{\epsilon}\right) \int_{q=0}^{d_j} \int_{x=w_{o,2}(q)}^{w_{o,2}(q)+w_j} \frac{x}{Q_2(x)} dx dq$$

Lemma 3. *The running condition holds with constant $c = 1 + \frac{1}{\epsilon}$.*

Proof. Let us consider an infinitesimally small interval $[t, t+dt)$ during which no jobs arrive and analyze the change in the potential $\Phi(t)$. Since $\Phi(t) = \sum_i \Phi_i(t)$, we can do this on a per-processor basis. Fix a single processor i , and time t . Let $w_i(q) := (w_{a,i}(q) - w_{o,i}(q))_+$, and $w_i := (w_{a,i} - w_{o,i})_+$. Let q_a and q_o denote

the inverse densities of the jobs being executed on processor i by the algorithm and optimal solution respectively (which are the densest jobs in their respective queues, since both run HDF). Define $s_a = Q_i(w_{a,i})$ and $s_o = Q_i(w_{o,i})$. Since we assumed that OPT uses the BCP algorithm on each processor, OPT runs processor i at speed s_o . Since the online algorithm is also using BCP, but has $(1 + \epsilon)$ -speed augmentation, the online algorithm runs the processor at speed $(1 + \epsilon)s_a$. Hence the fractional weight of the job the online algorithm works on decreases at a rate of $s_a(1 + \epsilon)/q_a$. Therefore, the quantity $w_{a,i}(q)$ drops by $s_a dt(1 + \epsilon)/q_a$ for $q \in [0, q_a]$. Likewise, $w_{o,i}(q)$ drops by $s_o dt/q_o$ for $q \in [0, q_o]$ due to the optimal algorithm working on its densest job. We consider several different cases based on the values of $q_o, q_a, w_{o,i}$, and $w_{a,i}$ and establish bounds on $d\Phi_i(t)/dt$; Recall the definition of $\Phi_i(t)$ from equation (3):

$$\Phi_i(t) = \frac{2}{\epsilon} \int_{q=0}^{\infty} \int_{x=0}^{(w_{a,i}^t(q) - w_{o,i}^t(q))^+} \frac{x}{Q_i(x)} dx dq$$

Case (1): $w_{a,i} < w_{o,i}$: The only possible increase in potential function occurs due to the decrease in $w_{o,i}(q)$, which happens for values of $q \in [0, q_o]$. But for q 's in this range, $w_{a,i}(q) \leq w_{a,i}$ and $w_{o,i}(q) = w_{o,i}$. Thus the inner integral is empty, resulting in no increase in potential. The running condition then holds since $w_{a,i} < w_{o,i}$.

Case (2): $w_{a,i} > w_{o,i}$: To quantify the change in potential due to the online algorithm working, observe that for any $q \in [0, q_a]$, the inner integral of Φ_i decreases by

$$\int_{x=0}^{w_i(q)} \frac{x}{Q_i(x)} dx - \int_{x=0}^{w_i(q) - (1+\epsilon)\frac{s_a dt}{q_a}} \frac{x}{Q_i(x)} dx = \frac{w_i(q)}{Q_i(w_i(q))} (1 + \epsilon) \frac{s_a dt}{q_a}$$

Here, we have used the fact that dt is infinitesimally small to get the above equality. Hence, the total drop in Φ_i due to the online algorithm's processing is

$$\begin{aligned} \frac{2}{\epsilon} \int_{q=0}^{q_a} \frac{w_i(q)}{Q_i(w_i(q))} (1 + \epsilon) \frac{s_a dt}{q_a} dq &\geq \frac{2}{\epsilon} \int_{q=0}^{q_a} \frac{w_i}{Q_i(w_i)} (1 + \epsilon) \frac{s_a dt}{q_a} dq \\ &= \frac{2}{\epsilon} \frac{w_i}{Q_i(w_i)} (1 + \epsilon) s_a dt \end{aligned}$$

Here, the first inequality holds because $x/Q_i(x)$ is a non-decreasing function, and for all $q \in [0, q_a]$, we have $w_{a,i}(q) = w_{a,i}$ and $w_{o,i}(q) \leq w_{o,i}$ and hence $w_i(q) \geq w_i$.

Now to quantify the increase in the potential due to the optimal algorithm working: observe that for $q \in [0, q_o]$, the inner integral of Φ_i increases by at most

$$\int_{x=w_i(q)}^{w_i(q) + \frac{s_o dt}{q_o}} \frac{x}{Q_i(x)} dx = \frac{w_i(q)}{Q_i(w_i(q))} \frac{s_o dt}{q_o}$$

Again notice that we have used that fact that here dt is an infinitesimal period of time that in the limit is zero. Hence the total increase in Φ_i due to the optimal algorithm's processing is at most

$$\frac{2}{\epsilon} \int_{q=0}^{q_o} \frac{w_i(q)}{Q_i(w_i(q))} \frac{s_o dt}{q_o} dq \leq \frac{2}{\epsilon} \int_{q=0}^{q_o} \frac{w_i}{Q_i(w_i)} \frac{s_o dt}{q_o} dq = \frac{2}{\epsilon} \frac{w_i}{Q_i(w_i)} s_o dt.$$

Again here, the first inequality holds because $x/Q_i(x)$ is a non-decreasing function, and for all $q \in [0, q_o]$, we have $w_{a,i}(q) \leq w_{a,i}$ and $w_{o,i}(q) = w_{o,i}$ and hence $w_i(q) \leq w_i$.

Putting the two together, the overall increase in $\Phi_i(t)$ can be bounded by

$$\begin{aligned} \frac{d\Phi_i(t)}{dt} &\leq \frac{2}{\epsilon} \frac{w_{a,i} - w_{o,i}}{Q_i(w_{a,i} - w_{o,i})} [-(1 + \epsilon)s_a + s_o] \\ &= \frac{2}{\epsilon} (w_{a,i} - w_{o,i}) \frac{[-(1 + \epsilon)Q_i(w_{a,i}) + Q_i(w_{o,i})]}{Q_i(w_{a,i} - w_{o,i})} \\ &\leq -\frac{2}{\epsilon} \epsilon (w_{a,i} - w_{o,i}) = -2(w_{a,i} - w_{o,i}) \end{aligned}$$

It is now easy to verify that by plugging this bound on $\frac{d\Phi_i(t)}{dt}$ into the running condition that one gets a valid inequality.

Case (3): $w_{a,i} = w_{o,i}$: In this case, let us just consider the increase due to OPT working. The inner integral in the potential function starts off from zero (since $w_{a,i} - w_{o,i} = 0$) and potentially (in the worst case) could increase to

$$\int_0^{\frac{s_o dt}{q_o}} \frac{x}{Q_i(x)} dx$$

(since $w_{o,i}$ drops by $s_o dt/q_o$ and $w_{a,i}$ cannot increase). However, since $x/Q_i(x)$ is a monotone non-decreasing function, this is at most

$$\int_0^{\frac{s_o dt}{q_o}} \frac{w_{o,i}}{Q_i(w_{o,i})} dx = \frac{s_o dt}{q_o} \frac{w_{o,i}}{Q_i(w_{o,i})}$$

Therefore, the total increase in the potential $\Phi_i(t)$ can be bounded by

$$\frac{2}{\epsilon} \int_{q=0}^{q_o} \frac{w_{o,i}}{Q_i(w_{o,i})} \frac{s_o dt}{q_o} dq = \frac{2}{\epsilon} s_o dt \frac{w_{o,i}}{Q_i(w_{o,i})} = \frac{2}{\epsilon} w_{o,i} dt$$

It is now easy to verify that by plugging this bound on $\frac{d\Phi_i(t)}{dt}$ into the running condition, and using the fact that $w_{a,i} = w_{o,i}$, one gets a valid inequality.

3 The Algorithm for Unweighted Flow

In this section, we give an immediate assignment based scheduling policy and the potential function that can be used to show (using basically the same line

of reasoning as in the last section) that it is $O(1/\epsilon)$ -competitive against a non-migratory adversary for the objective of *unweighted* flow plus energy, assuming the online algorithm has resource augmentation of $(1 + \epsilon)$ in speed. Note that this result has a better competitiveness than the result for weighted flow from Section 2, but holds only for the unweighted case. Once again we can assume that the adversary is running the BCP algorithm on each processor. Just like the weighted case, the crux of our algorithm would be in designing a good assignment policy.

Our algorithm works as follows: Each processor maintains a queue of jobs that have currently been assigned to it. At some time instant t , for any processor i , let $n_{a,i}^t(q)$ denote the number of jobs in processor i 's queue that have a remaining processing time of at least q . Let $n_{a,i}^t$ denote the total number of unfinished jobs in the queue. Now our algorithm is the following:

When a new job arrives, the assignment policy assigns it to a processor which would cause the smallest increase in the “shadow potential”; i.e., a processor minimizing

$$\int_{q=0}^p \sum_{j=1}^{n_{a,i}^t(q)+1} \frac{j}{Q_i(j+1)} dq - \int_{q=0}^p \sum_{j=1}^{n_{a,i}^t(q)} \frac{j}{Q_i(j+1)} dq = \int_{q=0}^p \frac{(n_{a,i}^t(q) + 1)}{Q_i(n_{a,i}^t(q) + 2)} dq$$

The job selection on each processor is SRPT (Shortest Remaining Processing Time), and we set the power of processor i at time t to $n_{a,i}^t + 2$ if $n_{a,i}^t \neq 0$ (and power zero otherwise). Once the job is assigned to a processor, it is never migrated.

Note that, even without accounting for resource augmentation, the online algorithm runs processor i at a speed of $Q_i(n_{a,i}^t + 2) \cdot \mathbf{1}_{(n_{a,i}^t > 0)}$, instead of $Q_i(n_{a,i}^t + 1) \cdot \mathbf{1}_{(n_{a,i}^t > 0)}$ as in [5]; since OPT uses BCP without any changes, it runs processor i at speed $Q_i(n_{o,i}^t + 1) \cdot \mathbf{1}_{(n_{o,i}^t > 0)}$.

We now describe our potential function Φ . For time t and processor i . Define $n_{o,i}^t$ as the number of unfinished jobs assigned to processor i by the optimal solution at time t , and $n_{o,i}^t(q)$ to be the number of these jobs with remaining processing time at least q . The global potential function is $\Phi(t) = \sum_i \Phi_i(t)$, where $\Phi_i(t)$ is the potential for processor i defined as:

$$\Phi_i(t) = \frac{4}{\epsilon} \int_{q=0}^{\infty} \sum_{j=1}^{(n_{a,i}^t(q) - n_{o,i}^t(q))_+} j / Q_i(j+1) dq \quad (4)$$

Recall that $(x)_+ = \max(x, 0)$, and $Q_i = P_i^{-1}$.

Acknowledgments: We thank Sangyeun Cho and Bruce Childers for helpful discussions about heterogeneous multicore processors, and also Srivatsan Narayanan for several useful discussions.

References

1. Albers, S., Fujiwara, H.: Energy-efficient algorithms for flow time minimization. *ACM Transactions on Algorithms* 3(4) (2007)
2. Andrew, L.L., Wierman, A., Tang, A.: Optimal speed scaling under arbitrary power functions. *SIGMETRICS Performance Evaluation Review* 37(2), 39–41 (2009)
3. Bansal, N., Chan, H.L.: Weighted flow time does not admit $o(1)$ -competitive algorithms. In: *SODA*. pp. 1238–1244 (2009)
4. Bansal, N., Chan, H.L., Lam, T.W., Lee, L.K.: Scheduling for speed bounded processors. In: *ICALP* (1). pp. 409–420 (2008)
5. Bansal, N., Chan, H.L., Pruhs, K.: Speed scaling with an arbitrary power function. In: *SODA*. pp. 693–701 (2009)
6. Bansal, N., Pruhs, K., Stein, C.: Speed scaling for weighted flow time. *SIAM Journal on Computing* 39(4) (2009)
7. Becchetti, L., Leonardi, S., Marchetti-Spaccamela, A., Pruhs, K.: Online weighted flow time and deadline scheduling. *J. Discrete Algorithms* 4(3), 339–352 (2006)
8. Bower, F.A., Sorin, D.J., Cox, L.P.: The impact of dynamically heterogeneous multicore processors on thread scheduling. *IEEE Micro* 28(3), 17–25 (2008)
9. Chadha, J.S., Garg, N., Kumar, A., Muralidhara, V.N.: A competitive algorithm for minimizing weighted flow time on unrelated machines with speed augmentation. In: *STOC*. pp. 679–684 (2009)
10. Chan, H.L., Edmonds, J., Lam, T.W., Lee, L.K., Marchetti-Spaccamela, A., Pruhs, K.: Nonclairvoyant speed scaling for flow and energy. In: *STACS*. pp. 255–264 (2009)
11. Greiner, G., Nonner, T., Souza, A.: The bell is ringing in speed-scaled multiprocessor scheduling. In: *SPAA '09: Proceedings of the twenty-first annual symposium on Parallelism in algorithms and architectures*. pp. 11–18. ACM, New York, NY, USA (2009)
12. Kumar, R., Tullsen, D.M., Jouppi, N.P.: Core architecture optimization for heterogeneous chip multiprocessors. In: *International conference on parallel architectures and compilation techniques*. pp. 23–32. ACM (2006)
13. Kumar, R., Tullsen, D.M., Ranganathan, P., Jouppi, N.P., Farkas, K.I.: Single-isa heterogeneous multi-core architectures for multithreaded workload performance. *SIGARCH Computer Architecture News* 32(2), 64 (2004)
14. Lam, T.W., Lee, L.K., To, I.K.K., Wong, P.W.H.: Competitive non-migratory scheduling for flow time and energy. In: *SPAA*. pp. 256–264 (2008)
15. Lam, T.W., Lee, L.K., To, I.K.K., Wong, P.W.H.: Speed scaling functions for flow time scheduling based on active job count. In: *European Symposium on Algorithms*. pp. 647–659 (2008)
16. Leonardi, S., Raz, D.: Approximating total flow time on parallel machines. *Journal of Computer and Systems Sciences* 73(6), 875–891 (2007)
17. Merritt, R.: CPU designers debate multi-core future. *EE Times* (February 2008)
18. Morad, T.Y., Weiser, U.C., Kolodny, A., Valero, M., Ayguade, E.: Performance, power efficiency and scalability of asymmetric cluster chip multiprocessors. *IEEE Computer Architecture Letters* 5(1), 4 (2006)
19. Pruhs, K.: Competitive online scheduling for server systems. *SIGMETRICS Performance Evaluation Review* 34(4), 52–58 (2007)
20. Pruhs, K., Sgall, J., Torng, E.: Online scheduling. In: *Handbook on Scheduling*. CRC Press (2004)
21. Pruhs, K., Uthaisombut, P., Woeginger, G.J.: Getting the best response for your erg. *ACM Transactions on Algorithms* 4(3) (2008)