

## Scheduling with Equipartition

2000; Edmonds

JEFF EDMONDS

York University, Toronto, ON, Canada

### Synonyms

Round Robin and Equi-partition are the same algorithm. Average Response time and Flow are basically the same measure.

### Problem Definition

The task is to schedule a set of  $n$  on-line jobs on  $p$  processors. The jobs are  $J = \{J_1, \dots, J_n\}$  where job  $J_i$  has a release/arrival time  $r_i$  and a sequence of phases  $\langle J_i^1, J_i^2, \dots, J_i^{q_i} \rangle$ . Each phase is represented by  $\langle w_i^q, \Gamma_i^q \rangle$ , where  $w_i^q$  denotes the amount of work and  $\Gamma_i^q$  is the speedup function specifying the rate  $\Gamma_i^q(\beta)$  at which this work is executed when given  $\beta$  processors.

A phase of a job is said to be *fully parallelizable*  $\checkmark$  if its speedup function is  $\Gamma(\beta) = \beta$ . It is said to be *sequential*  $\perp$  if its speedup function is  $\Gamma(\beta) = 1$ .<sup>1</sup> A speedup function  $\Gamma$  is *nondecreasing* iff  $\Gamma(\beta_1) \leq \Gamma(\beta_2)$  whenever  $\beta_1 \leq \beta_2$ ,<sup>2</sup> is *sublinear*  $\prec$  iff  $\Gamma(\beta_1)/\beta_1 \geq \Gamma(\beta_2)/\beta_2$ ,<sup>3</sup> and is *strictly-sublinear* by  $\alpha$  iff  $\Gamma(\beta_2)/\Gamma(\beta_1) \leq (\beta_2/\beta_1)^{1-\alpha}$ .

An *s-speed scheduling algorithm*  $S_s(J)$  allocates  $s \times p$  processors each point in time to the jobs  $J$  in a way such that all the work completes.<sup>4</sup> More formally, it constructs a function  $\mathcal{S}(i, t)$  from  $\{1, \dots, n\} \times [0, \infty)$  to  $[0, sp]$  giving the number of processors allocated to job  $J_i$  at time  $t$ . (A job is allowed to be allocated a non-integral number of processors.) Requiring that for all  $t$ ,  $\sum_{i=1}^n \mathcal{S}(i, t) \leq sp$  ensures that at most  $sp$  processors are allocated at any given time. Requiring that for all  $i$ , there exist  $r_i = c_i^0 < c_i^1 < \dots <$

<sup>1</sup>Note that an odd feature of this definition is that a sequential job completes work at a rate of 1 even when absolutely no processors are allocated to it. This assumption makes things easier for the adversary and harder for any non-clairvoyant algorithm. Hence, it only makes these results stronger.

<sup>2</sup>A job phase with a nondecreasing speedup function executes no slower if it is allocated more processors.

<sup>3</sup>A measure of how efficient a job utilizes its processors is  $\Gamma(\beta)/\beta$ , which is the work completed by the job per time unit per processor. A sublinear speedup function is one whose efficiency does not increase with more processors. This is a reasonable assumption if in practice  $\beta_1$  processors can simulate the execution of  $\beta_2$  processors in a factor of at most  $\beta_2/\beta_1$  more time.

<sup>4</sup> $S^s(J)$  is defined to be the scheduler with  $p$  processors of speed  $s$ .  $S_s$  and  $S^s$  are equivalent on fully parallelizable jobs and  $S^s$  is  $s$  times faster than  $S_s$  on sequential jobs.

$c_i^q$  such that for all  $1 \leq q \leq q_i$ ,  $\int_{c_i^{q-1}}^{c_i^q} \Gamma_i^q(\mathcal{S}(i, t)) dt = w_i^q$  ensures that before a phase of a job begins, the job must have been released and all of the previous phases of the job must have completed. The completion time of a job  $J_i$ , denoted  $c_i$ , is the completion time of the last phase of the job.

The goal of a scheduling algorithm is to minimize the average response time,  $\frac{1}{n} \sum_{i \in J} (c_i - r_i)$ , of the jobs or equivalently its flow time  $S_s(J) = \sum_{i \in J} (c_i - r_i)$ . An alternative formalization is to integrate over time the number of jobs  $n_t$  alive at time  $t$ ,  $S_s(J) = \sum_{i \in J} \int_0^\infty (J_i \text{ is alive a time } t) \delta t = \int_0^\infty n_t \delta t$ .

A scheduling algorithm is said to be *on-line* if it lacks knowledge of which jobs will arrive in the future. It is said to be *non-clairvoyant* if it also lacks all knowledge about the jobs that are currently in the system, except for knowing when a job arrives and knowing when it completes.

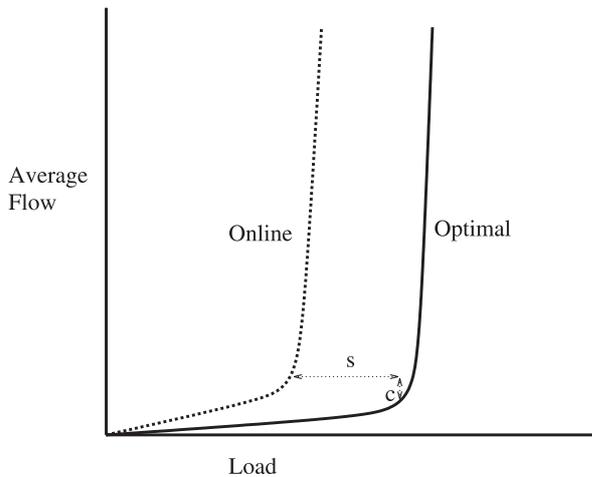
The two examples of *non-clairvoyant* schedulers that are often used in practice are *Equipartition* (also called *Round Robin*) and *Balance*.  $EQUI_s$  is defined to be the scheduler that allocates an equal number of processors to each job that is currently alive. That is, for all  $i$  and  $t$ , if job  $J_i$  is alive at time  $t$ , then  $\mathcal{EQUI}(i, t) = sp/n_t$ , where  $n_t$  is the number of jobs that are alive at time  $t$ . The schedule  $BAL_s$  is defined in [8] to be the schedule that allocates all of its processors to the job that has been allocated processors for the shortest length of time. (Though no one implements *Balance* directly, Unix uses a multi-level feedback (MLF) queue algorithm which in a way approximates *Balance*).

The most obvious worst-case measure of the goodness of an online non-clairvoyant scheduling algorithm  $S$  is its *competitive ratio*. This compares the perform of the algorithm to that of the optimal scheduler. However, in many cases, the limited algorithm is unable to compete against an all knowing all powerful optimal scheduler. To compensate the algorithm  $S_s$ , it is given extra speed  $s$ . An online scheduling algorithm  $S$  is said to be *s-speed c-competitive* if:  $\max_J S_s(J)/Opt(J) \leq c$ . For example, being *s-speed 2-competitive* means that the cost  $S_s(J)$  of scheduler  $S$  with  $s \times p$  processors on any instance  $J$  is at most twice the optimal cost for the same jobs when only given  $p$  processors.

### Key Results

If all jobs arrive at time zero (batch), then the flow time of  $EQUI$  is 2-competitive on fully parallelizable jobs [10] and  $(2 + \sqrt{3})$ -competitive on jobs with nondecreasing sublinear speedup functions [3]. (The time until the last job completes (makespan) on fully parallelizable jobs is the same for  $EQUI$  and  $OPT$ , but can be a factor of  $\Theta(\log n / \log \log n)$  worse for  $EQUI$  if the jobs can also have sequential phases [11].) Table 1 summarizes all the results.

Please note that the pagination is not final; in the print version an entry will in general not start on a new page.



Scheduling with Equipartition, Figure 1

To understand the motivation for this *resource augmentation analysis* [8], note that it is common for the quality of service of a system to have a *threshold property* with respect to the load that it is given. In this example, it seems that the online scheduling algorithm  $S$  performs reasonably well in comparison to the optimal scheduling algorithm. Despite this, one can see that the competitive ratio of  $S$  is huge by looking at the vertical gap between the curves when the load is near capacity. To explain why these curves are close, one must also measure the horizontal gap between curves.  $S$  performs at most  $c$  times worse than optimal, when either the load is decreased or equivalently the speed is increased by factor of  $s$

80 When the jobs have arbitrary arrival times and are  
 81 fully parallelizable, the optimal schedule simply allocates  
 82 all the processors to the jobs with least remaining work.  
 83 This, however, requires the scheduler to know the amount  
 84 of work per job. Without this knowledge,  $EQUI$  and  $BAL$   
 85 are unable to compete with the optimal and hence can do  
 86 a factor of  $\Omega(n/\log n)$  and  $\Omega(n)$  respectively worse and no  
 87 non-clairvoyant scheduler has a better competitive ratio  
 88 than  $\Omega(n^{1/3})$  [9,10]. Randomness improves the competi-  
 89 tive ratio of  $BAL$  to  $\Theta(\log n \log \log n)$  [7]. Having more (or  
 90 faster) processors also helps.  $BAL_s$  achieves a  $s = 1 + \epsilon$  speed  
 91 competitive ratio of  $\frac{s}{s-1} = 1 + \frac{1}{\epsilon}$  [8].

92 If some of the jobs are fully parallelizable and some are  
 93 sequential jobs, it is hard to believe that any non-clairvoyant  
 94 scheduler, even with  $sp$  processors, can perform well.  
 95 Not knowing which jobs are which, it waists too many pro-  
 96 cessors on the sequential jobs. Being starved, the fully par-  
 97 allelizable jobs fall further and further behind and then  
 98 other fully parallelizable jobs arrive which fall behind as  
 99 well. For example, even the randomized version of  $BAL$  can  
 100 have an arbitrarily bad competitive ratio, even when given  
 101 arbitrarily fast processors.

$EQUI$ , however, does amazingly well.  $EQUI_s$  achieves  
 a  $s = 2 + \epsilon$  speed competitive ratio of  $2 + \frac{4}{\epsilon}$  [1]. This was  
 later improved to  $1 + \mathcal{O}(\frac{\sqrt{s}}{s-2})$ , which is better for large  $s$  [1].  
 The intuition is that  $EQUI_s$  is able to automatically “self  
 adjust” the number of processors wasted on the sequential  
 jobs. As it falls behind, it has more uncompleted jobs in  
 the system and hence allocates fewer processors to each  
 job and hence each job utilizes the processors that it is  
 given more efficiently. The extra processors are enough  
 to compensate for the fact that some processors are still  
 wasted on sequential jobs. For example, suppose the job  
 set is such that  $OPT$  has  $\ell_t$  sequential jobs and at most  
 one fully parallelizable job alive at any point in time  $t$ .  
 (The proof starts by proving that this is the worst case.)  
 It may take a while for the system under  $EQUI_s$  to reach  
 a “steady state”, but when it does,  $m_t$ , which denotes the  
 number of fully parallelizable jobs it has alive at time  $t$ ,  
 converges to  $\frac{\ell_t}{s-1}$ . At this time,  $EQUI_s$  has  $\ell_t + m_t$  jobs  
 alive and  $OPT$  has  $\ell_t + 1$ . Hence, the competitive ratio  
 is  $EQUI_s(J)/OPT(J) = (\ell_t + \frac{\ell_t}{s-1})/(\ell_t + 1) \leq \frac{s}{s-1}$  **CE2**,  
 which is  $1 + \frac{1}{\epsilon}$  for  $s = 1 + \epsilon$ . This intuition makes it appear  
 that speed  $s = 1 + \epsilon$  is sufficient. However, unless the speed  
 is at least 2 then the competitive ratio can be bad during  
 the time until it reaches this steady state, [8].

More surprisingly if all the jobs are *strictly sublinear*,  
 i. e., are not fully parallel, then  $EQUI$  performs competi-  
 tively with no extra processors [1]. More specifically, it is  
 shown that if all the speedup functions are no more fully  
 parallelizable than  $\Gamma(\beta) = \beta^{1-\alpha}$  than the competitive ra-  
 tio is at most  $2^{\frac{1}{\alpha}}$ . For intuition, suppose the adversary al-  
 locates  $\frac{p}{n}$  processors to each of  $n$  jobs and  $EQUI$  falls be-  
 hind enough so that it has  $2^{\frac{1}{\alpha}}n$  uncompleted jobs. Then it  
 allocates  $p/(2^{\frac{1}{\alpha}}n)$  processors to each, completing work at  
 an overall rate of  $(2^{\frac{1}{\alpha}}n)\Gamma(p/(2^{\frac{1}{\alpha}}n)) = 2 \cdot n\Gamma(p/n)$ . This is  
 a factor of 2 more than that by the adversary. Hence, as in  
 the previous result,  $EQUI$  has twice the speed and so per-  
 forms competitively.

The results for  $EQUI_s$  can be extended further. There  
 is a competitive  $s = (8 + \epsilon)$ -speed non-clairvoyant sched-  
 uler that only preempts when the number of jobs in the  
 system goes up or down by a factor of two (in some sense  
 $\log n$  times). There is  $s = (4 + \epsilon)$ -speed one that includes  
 both sublinear  $\llcorner$  and superlinear  $\lrcorner$  jobs. Finally, there is  
 a  $s = \mathcal{O}(\log p)$  speed one that includes both nondecreasing  
 $\llcorner$  and gradual  $\triangleleft$  jobs.

The proof of these results for  $EQUI_s$  require techniques  
 that are completely new. For example, the previous results  
 prove that their algorithm is competitive by proving that at  
 every point in time, the number of jobs alive under their al-

**CE2** Unbalanced parantheses. Please check.

**Scheduling with Equipartition, Table 1**

Each row represents a specific scheduler and a class  $\mathcal{J}$  of job sets. Here  $EQUI_s$  denotes the Equipartition scheduler with  $s$  times as many processors and  $EQUI^s$  the one with processors that are  $s$  times as fast. The graphs give examples of speedup functions from the class of those considered. The columns are for different extra resources ratios  $s$ . Each entry gives the corresponding ratio between the given scheduler and the optimal

	$s = 1$	$s = 1 + \epsilon$	$s = 2 + \epsilon$	$s = 4 + 2\epsilon$	$s = \mathcal{O}(\log p)$
Batch	[2.71, 3.74]				
Det. Non-clair	$\Omega(n^{\frac{1}{3}})$	–			
Rand. Non-clair	$\tilde{\Theta}(\log n)$	–			
Rand. Non-clair	$\Omega(n^{\frac{1}{2}})$	$\Omega(\frac{1}{\epsilon})$			
$BAL_s$	$\Omega(n)$	$1 + \frac{1}{\epsilon}$		$\frac{2}{s}$	
$BAL_s$	$\Omega(s^{-1/\alpha}n)$				
$EQUI_s$	$\Omega(\frac{n}{\log n})$	$\Omega(n^{1-\epsilon})$	$[1 + \frac{1}{\epsilon}, 2 + \frac{4}{\epsilon}]$		$\geq 1$
$EQUI^s$	$\Omega(\frac{n}{\log n})$	$\Omega(n^{1-\epsilon})$	$[\frac{2}{3}(1 + \frac{1}{\epsilon}), 2 + \frac{4}{\epsilon}]$		$[\frac{2}{3}, \frac{16}{3}]$
$EQUI$	$[1.48^{1/\alpha}, 2^{1/\alpha}]$				
$EQUI'_s$ Few Preempts			$\Omega(n^{1-\epsilon})$	$\Theta(1)$	
$HEQUI_s$			$\Omega(n^{1-\epsilon})$	$\Theta(1)$	
$HEQUI'_s$			$\Omega(n)$		$\Theta(1)$

151 gorithm is within a constant fraction of that under the opti- 179  
 152 mal schedule. This, however, is simply not true with this less 180  
 153 restricted model. There are job sets such that for a period of 181  
 154 time the ratio between the numbers of alive jobs under the 182  
 155 two schedules is unbounded. Instead, a potential function 183  
 156 is used to prove that this can only happen for a relatively 184  
 157 short period of time.

158 The proof first transforms each possible input into 185  
 159 a canonical input that as described above only has paral- 186  
 160 lelizable or sequential phases. Having the number of fully 187  
 161 parallelizable jobs alive under  $EQUI_s$  at time  $t$  be much big- 188  
 162 ger than the number of sequential jobs alive at this same 189  
 163 time is bad for  $EQUI_s$  because it then has many more jobs 190  
 164 alive than  $OPT$  and hence is currently incurring much 191  
 165 higher costs. On the other hand, this same situation is also 192  
 166 good for  $EQUI_s$  because it means that it is allocating a larger 193  
 167 fraction of its processors to the fully parallelizable jobs and 194  
 168 hence is catching up to  $OPT$ . Both of these aspects of the 195  
 169 current situation is carefully measured in a potential func- 196  
 170 tion  $\Phi(t)$ . It is proven that at each point in time, the oc- 197  
 171 curred cost to  $EQUI_s$  plus the gain  $\frac{d\Phi(t)}{dt}$  in this potential 198  
 172 function is at most  $c$  times the costs occurred by  $OPT$ . As- 199  
 173 suming that the potential function begins and ends at zero, 200  
 174 the result follows.

175 More formally, the potential function is  $\Phi(t) = F(t) +$  201  
 176  $Q(t)$  where  $Q(t)$  is total sequential work finished by  $EQUI_s$   
 177 by time  $t$  minus the total sequential work finished by  
 178 the adversary by time  $t$ . To define  $F(t)$  requires some

179 preliminary definitions. For  $u \geq t$ , define  $m_u(t)$  ( $\ell_u(t)$ )  
 180 to be number of fully parallelizable (sequential) phases  
 181 executing under  $EQUI_s$  at time  $u$ , for which  $EQUI_s$  at  
 182 time  $u$  has still not processed as much work as the ad-  
 183 versary processed at time  $t$ . Let  $n_u(t) = m_u(t) + \ell_u(t)$ .  
 184 Then  $F(t) = \int_t^\infty f_u(m_u(t), \ell_u(t)) du$ , where  $f_u(m, \ell) =$   
 185  $\frac{s}{s-2} \frac{(m-\ell)(m+\ell)}{n_u}$ . As the definition of the potential function  
 186 suggests, the analysis is quite complicated.

**Applications** 187

188 In addition to being interesting results on their own, they  
 189 have been powerful tools for the theoretical analysis of  
 190 other on-line algorithms. For example, in [2,4] TCP was  
 191 reduced to this problem and in [5], the online broadcast  
 192 scheduling problem was reduced to this problem.

**Open Problems** 193

194 An open question is whether there is an algorithm that is  
 195 competitive when given processors of speed  $s = 1 + \epsilon$  (as  
 196 opposed to  $s = 2 + \epsilon$ ). There is a candidate algorithm that  
 197 is part way between  $EQUI_s$  and  $BAL_s$ .

**Cross References** 198

- ▶ Flow Time Minimization 199
- ▶ List Scheduling 200
- ▶ Minimum Flow 201

- 202 ▶ [Minimum Weighted Completion Time](#)
- 203 ▶ [Online List Update](#)
- 204 ▶ [Online Load Balancing](#)
- 205 ▶ [Schedulers for Optimistic Rate Based Flow Control](#)
- 206 ▶ [Shortest Elapsed Time First Scheduling](#)

### 207 Recommended Reading

- 208 1. Edmonds, J.: Scheduling in the dark. Improved results:  
209 manuscript 2001. In: Theor. Comput. Sci. **235**, 109–141 (2000).  
210 In: 31st Ann. ACM Symp. on Theory of Computing, 1999
- 211 2. Edmonds, J.: On the Competitiveness of AIMD-TCP within  
212 a General Network. In: **CE3** LATIN, Latin American Theoretical  
213 Informatics, vol. 2976, pp. 577–588 (2004). Submitted to  
214 Journal Theoretical Computer Science and/or Lecture Notes in  
215 Computer Science **CE4**
- 216 3. Edmonds, J., Chinn, D., Brecht, T., Deng, X.: Non-clairvoyant  
217 Multiprocessor Scheduling of Jobs with Changing Execution  
218 Characteristics. In: 29th Ann. ACM Symp. on Theory of Com-  
219 puting, 1997, pp. 120–129. Submitted to SIAM J. Comput. **CE5**
- 220 4. Edmonds, J., Datta, S., Dymond, P.: TCP is Competitive Against  
221 a Limited Adversary. In: SPAA, ACM Symp. of Parallelism in Al-  
222 gorithms and Architectures, 2003, pp. 174–183
- 223 5. Edmonds, J., Pruhs, K.: Multicast pull scheduling: when fairness  
224 is fine. *Algorithmica* **36**, 315–330 (2003)
- 225 6. Edmonds, J., Pruhs, K.: A maiden analysis of longest wait first.  
226 In: Proc. 15th Symp. on Discrete Algorithms (SODA) **CE6**
- 227 7. Kalyanasundaram, B., Pruhs, K.: Minimizing flow time nonclair-  
228 voyantly. In: Proceedings of the 38th Symposium on Founda-  
229 tions of Computer Science, October 1997
- 230 8. Kalyanasundaram, B., Pruhs, K.: Speed is as powerful as clair-  
231 voyance. In: Proceedings of the 36th Symposium on Founda-  
232 tions of Computer Science, October 1995, pp. 214–221
- 233 9. Matsumoto, **CE7**: Competitive Analysis of the Round Robin Al-  
234 gorithm. in: 3rd International Symposium on Algorithms and  
235 Computation, 1992, pp. 71–77
- 236 10. Motwani, R., Phillips, S., Torng, E.: Non-clairvoyant scheduling.  
237 Theor. Comput. Sci. **130** (Special Issue on Dynamic and On-  
238 Line Algorithms), 17–47 (1994). Preliminary Version in: Pro-  
239 ceedings of the 4th Annual ACM-SIAM Symposium on Discrete  
240 Algorithms, 1993, pp. 422–431
- 241 11. Robert, J., Schabanel, N.: Non-Clairvoyant Batch Sets Schedul-  
242 ing: Fairness is Fair enough. Personal Correspondence (2007)

**CE3** Please provide editor.

**CE4** Please clarify and update.

**CE5** Please update.

**CE6** Please provide location and date of the symposium.

**CE7** Please provide initials of “Matsumoto”