

A Brief Survey of Multi-Processor Scheduling For Hard Real-Time Systems

Xin Lin^a, Xiaorong Zhu^a, Lijia Liu^a

^a*Department of Computer Science, The University of Texas at Austin*

Abstract

In class, both of scheduling algorithms [1] and priority inheritance protocols [2] in the context of a single processor were examined in details. Nevertheless, the emergence and popularity of distributed computing system gave rise to the need to solve multi-processor scheduling and priority inheritance problems. As the supplementary study, this paper surveys existing scheduling algorithms in the context of multiple processors. The very first section outlines the background of multi-processor scheduling problems, as well as system models, terminology, and the metrics of scheduling algorithms. After that, partitioned scheduling and global scheduling, as the primary objects of our research, will be fully explored. Moreover, we will also give brief sketch to the hybrid approaches of partitioned scheduling and global scheduling.

Keywords: System, Scheduling Algorithm, Task Management

1. Introduction

Real-time systems has become prevalent in diverse application areas, such as traffic controls, embedded automotive electronics, networked multimedia etc. Typically, a real-time system comprises a real-time computer system and a controlled object, which are connected through sensors, actuators and some other input-output interfaces. The controlling real-time computer reacts/responds to its controlled object based on the currently known information about the object. For example, a real-time computer controls a device according to the data read by sensors at periodic intervals and responds by sending signals to actuators within a time bound. (Figure?)

The purpose of real-time systems is to meet various timing constraints imposed on it by its external world during the operation of the system. The

13 instant at which a response from the controlling system must be delivered to
14 its controlled system is called a deadline. A deadline can be classified into
15 three categories[3]:

- 16 • Soft Deadline: the response still has effectiveness even after the deadline
17 has passed. However the later the execution of a task after the deadline,
18 the more penalty it pays.
- 19 • Firm Deadline: the response has no utilize if the deadline has passed.
20 Moreover the earlier the execution of a task before a firm deadline, the
21 more rewards it gains.
- 22 • Hard Deadline: a task must complete its computation by the deadline,
23 otherwise a catastrophe could happen.

24 Commands and Control systems, Air traffic control systems are examples
25 for hard real-time systems. On-line transaction systems, airline reservation
26 systems are soft real-time systems[4]. Such a system is called embedded
27 real-time system if it is a component of a larger system.

28 Companies are motivated to build embedded real-time systems by their
29 advanced effectiveness and flexibility. In order to beat the competition in
30 industries, they need to satisfy customers needs and lower associated costs
31 to the maximum extent possible. Such aspiration has resulted in a rapid
32 progress in the area of embedded real-time technology. For instance, sili-
33 con vendors used to increase processor performance by concentrating on the
34 miniaturization needed. However this approach has led to problems with
35 both high power consumption and excessive heat dissipation[4]. Thus using
36 multiprocessor platforms for high-end real-time applications is an increasing
37 trend toward instead[4].

38 Multiprocessor systems can be classified into three categories from the
39 perspectives of scheduling[4]:

- 40 • Heterogeneous: processors in the system are different. Not all tasks
41 may be able to execute on all processors and the execution rate of a
42 task depends on both the processor and the task.
- 43 • Homogeneous: processors in the system are identical. The execution
44 rate of all tasks is the same on all processors. Thus it depends only on
45 the task.

- Uniform: the execution rate of a task depends only on the speed of processors. For example a processor with speed 2 will execute all tasks at exactly twice the rate of a processor with speed 1.

In this paper we are concerned with scheduling algorithms of **homogeneous** multiprocessor for **hard** real-time system. The paper is organized as follows: section 2 is problem statement, section 3 is key concepts, section 4 is performance evaluation. we will describe the partitioned scheduling and global scheduling in detail in section 5 and section 6, respectively. A conclusion is given at the end of this paper (need to revise).

2. Preliminaries

This section reviews the fundamental terminologies and notations used in multiprocessor scheduling algorithms.

2.1. Multiprocessor Scheduling

In computer science, multiprocessor scheduling is an NP-hard optimization problem[wiki]. The problem can be described as "given a set J of jobs where job j_i has length l_i and a number of processors m , what is the minimum possible time required to schedule all jobs in J on m processors such that none overlap?"[5]. Clearly multiprocessor schedulers should concern with the allocation of the resources and the satisfaction of timing constraints[4]:

- Resource Allocation: on which processor a task executes.
- Propriety Determination: when and in what order, with respect to jobs of other tasks, each job executes so that all jobs can be finished before their deadlines.

Then real-time scheduling algorithms for multiprocessor systems can be classified from perspectives of allocation changes and priority changes (referred to as migration-based and priority-based classifications [6]).

2.2. Classification of Real-time Scheduling Algorithms

Real-Time scheduling can be categorized into hard scheduling and soft scheduling based on the types of deadlines. The hard real-time scheduling can be further classified into two types: static and dynamic. A static scheduling generates its scheduling decisions off-line according to the prior

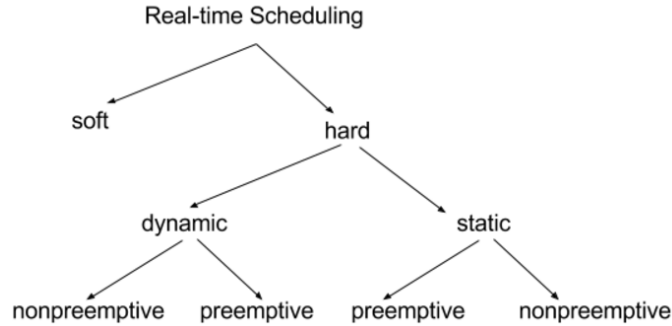


Figure 1: Taxonomy of Real-time Scheduling Algorithms

information of the jobs in a task. The decisions are made at compile time, thus its run-time overhead is small. By contrast, the scheduling decisions are made at run time in dynamic scheduling system, which leads to significant overheads. Furthermore, both dynamic scheduling and static scheduling can have preemptive or nonpreemptive scheduling:

- Preemptive: the currently executing task can be preempted by a higher priority task at anytime
- Nonpreemptive: the currently executing task can not be preempted until completion

In this paper, we are mainly concerned with dynamic preemptive scheduling algorithms. Figure 1 illustrates a rough taxonomy of real-time scheduling algorithms.

More specifically, the dynamic scheduling algorithm can be classified by its priority policy[4]:

- Fixed task-level priority: the jobs of a task has a single fixed priority, such as Rate Monotonic RM algorithm
- Fixed job-level priority: the jobs of a task may have different priorities, but each job has a single static priority, such earliest deadline first (EDF) scheduling.
- Dynamic priority: A single job may have different priorities at different times, such as least laxity first (LLF) scheduling.

98 From the respective of allocation policy, the real-time scheduling algo-
99 rithms for multiprocessor systems can be also classified as follows[4]:

- 100 • Non-migration: the jobs of a task must execute on one single processor
101 and migration is not permitted.
- 102 • Task-level migration: the jobs of a task may be allocated to different
103 processors, however, each job can only execute on one single processor.
- 104 • Job-level migration: a single job can be executed on different proces-
105 sors, however, no parallel execution of a job is permitted.

106 Scheduling algorithms without migrations are referred partitioned algo-
107 rithms while those that permits migrations are referred as global algorithms.
108 These two algorithms have been discussed in detail in section 3 and section
109 4, respectively.

110 2.3. Classifications of Task Models

111 Two simple task models have been defined for multiprocessor real-time
112 scheduling[4]:

- 113 • periodic task model: the jobs of a task arrive strictly periodically, sep-
114 arated by a fixed time interval.
- 115 • sporadic task Model: the interval time of two continuous jobs from a
116 task must be larger than a fixed time interval.

117 Each task π_i is characterized by: the deadline D_i , the period/minimum
118 separation time P_i , the worst-case computation time C_i and the utilization
119 U_i , where $U_i = C_i/P_i$. The response time R_i of a task is defined as the
120 interval between the task request time and the task satisfaction time. The
121 hyperperiod $H(\pi)$ is defined as the least common multiple (lcm) of all the
122 periodic tasks.

123 Some other important terminologies realted to taskset are[4]:

- 124 • schedulable task: its worst-case response time under a scheduling algo-
125 rithm is less than or equal to its deadline.
- 126 • schedulable taskset: all of its tasks are schedulable under a scheduling
127 algorithm.

- feasible taskset: for each possible sequence of jobs that may be generated by a taskset in given a system, there exists some scheduling algorithm to satisfy it without missing any deadlines.

2.4. Notations

The notations that are used frequently in multiprocessing real-time scheduling and in this paper are as follows[4]:

- π_i : task i at priority level i
- B_i : blocking time at priority level i
- C_i : worst-case execution time of task π_i
- D_i : relative deadline of task π_i
- δ_i : density of task π_i , $\delta_i = C_i / \min(D_i, T_i)$
- δ_{max} : maximum density of any task in the taskset
- δ_{sum} : taskset density (sum of task densities)
- f_A : speedup factor (resource augmentation factor) for scheduling algorithm A
- $h(t)$: processor demand in the interval $[0, t)$
- $H(\pi)$: hyperperiod of the taskset
- $load(\pi)$: processor load of taskset π
- $load(\pi, k)$: processor load of taskset π , due to tasks of priority higher than or equal to k
- m : number of processors
- $M_A(\pi)$: minimum number of processors needed to schedule taskset π using scheduling algorithm A
- n : number of tasks
- N : number of jobs (typically in the hyperperiod of the taskset)

- 153 • R_i : worst-case response time of task π
- 154 • R_A : approximation ratio for scheduling algorithm A
- 155 • t : time
- 156 • T_i : minimum interarrival time of task π
- 157 • u_i : utilization of task π
- 158 • u_{max} : maximum utilization of any task in the taskset
- 159 • u_{sum} : taskset utilization
- 160 • U_A : utilization upper bound for scheduling algorithm A

161 3. Partitioned Scheduling

162 In this section, we will review some partitioned approaches to multipro-
 163 cessor real-time scheduling, and then compare their performance.

164 3.1. *Characteristic of Partitioned Scheduling*

165 Partitioned scheduling provides capability for performing parallel pro-
 166 cessing and also for automation of batch execution of multiple processes. It
 167 can improve performance by breaking down a process that works on a large
 168 data set, to multiple parallel processes that work on smaller data sets, each
 169 contains part of the original data set. The general solution involves two algo-
 170 rithms: one to assign tasks to processors, known as the allocation, the other
 171 to schedule tasks that are assigned to each individual processor. The parti-
 172 tioning strategy also requires that all occurrences of a task to be executed
 173 on the same processor.

174 One main advantage of partitioned scheduling is that, after the allocation
 175 of tasks to processors, we can apply the optimal real-time scheduling tech-
 176 niques and analyses for uni-processor systems to each individual processor.

177 Comparing to the global scheduling, partitioned scheduling also has some
 178 advantages.

- 179 • A task that overruns its worst-case execution time will not affect tasks
 180 on other processors.
- 181 • There is no migration cost as each task only runs on a single processor.

- 182 • For each processor, the run-queue is much smaller comparing to the
 183 single global queue in the global scheduling, thus the overheads of man-
 184 aging the run-queue can be neglected.

185 However, the main disadvantage of the partitioned approach is that the
 186 allocation problem is analogous to a NP-Hard problem.

187 3.2. Rate Monotonic Next Fit Scheduling (RMNF)

188 The rate-monotonic-next-fit scheduling [7] is a partitioned scheduling al-
 189 gorithm for multiprocessors systems. According to this algorithm, tasks are
 190 first sorted in the descending order of their request rates. The scheduling
 191 scheme is that:

- 192 1. Starts from the first task i , and the first processor j .
- 193 2. Assign the task i to the processor j if j meets the requirement that
 194 together with all tasks that have been assigned to the processor j , tasks
 195 can be feasibly scheduled on processor j according to the rate-monotonic
 196 scheduling algorithm for a single processor. Otherwise, assign i to $j +$
 197 1 and increase j by 1.
- 198 3. Get the next task and repeat the previous two steps, unless there is no
 199 task left.

200 Let N be the number of processors required to be feasibly schedule a set of
 201 task by the RMNF algorithm, and N_0 be the minimum number of processors
 202 required to feasibly schedule the same set of tasks. Then as N_0 approaches
 203 infinity, we can get that [7]

$$2.4 \leq \frac{N}{N_0} \leq 2.67 \quad (1)$$

204 3.3. Rate Monotonic First-Fit Scheduling (RMFF)

205 The rate-monotonic-first-fit scheduling algorithm [7] is similar to the
 206 RMNF. According to the RMFF, the tasks are first sorted in the descending
 207 order of their request rates. We will use a counter N to be the number of
 208 processors required for scheduling the given tasks. N is first initiated to 1.
 209 The scheduling scheme is that:

- 210 1. Starts from the first task i .

2. Starts from the lowest-indexed processor. Assign the task i to the first processor that meets the requirement that together with all tasks that have been assigned to the processor, tasks can be feasibly scheduled on the processor according to the rate-monotonic scheduling algorithm for a single processor. If no processor with index less than N meets the requirement, increase N by 1.
3. Get the next task and repeat the previous step, unless there is no task left.

Similarly, the bound of the rate of the number of processors N required to be feasibly schedule the task set by this algorithm, and the minimum number of processors N_0 required to be feasibly schedule the same set of tasks can be obtained when N_0 approaches infinity:

$$2 \leq \lim_{N_0 \rightarrow \infty} \frac{N}{N_0} \leq \frac{4 \times 2^{1/3}}{1 + 2^{1/3}} \approx 2.33 \quad (2)$$

3.4. Rate Monotonic Best-Fit Scheduling (RMBF)

The rate-monotonic-best-fit scheduling [8] algorithm is based on the bin-packing heuristic. Best-Fit scheme chooses to assign tasks on a processor that can maximize the utilization of that processor.

Similar to RMNF and RMFF, we first sort all given tasks according to non-decreasing periods. Keep a counter N to be the number of processors required for scheduling the given set of tasks. N is set to 1 at the beginning. The scheduling algorithm works like this:

1. Starts from the first task i .
2. Starts from the lowest-indexed processor j . Let k_j and U_j denote the number of tasks already assigned to the processor j and the total utilization of the k_j tasks. Let u_i denote the utilization of the task i . Find the smallest-indexed processor j such that task i together with all k_j tasks can be feasibly scheduled by the rate-monotonic scheduling algorithm, and $2(1 + \frac{U_j}{k_j})^{-k_j} - 1$ be as small as possible, then assign task i to j , and set $k_j = k_j + 1$, $U_j = u_i + U_j$. If $j < m$ then set $m = j$.
3. Get the next task and repeat the previous step, unless there is no task left.

Similarly, the bound of rate of N and N_0 (with the same definitions in RMFF) can be obtained when N_0 approaches infinity:

$$2.3 \leq \lim_{N_0 \rightarrow \infty} \frac{N}{N_0} \leq 2 + \frac{3 - 2^{3/2}}{2(2^{1/3} - 1)} \approx 2.33 \quad (3)$$

Table 1: Approximation Ratios

Algorithm	Approximation Ratio $\frac{N}{N_0}$
RMNF	2.67
RMFF	2.33
RMBF	2.33

243 3.5. Comparison and Summary

244 In this section, we had introduced the three partitioned scheduling algo-
 245 rithms RMNF, RMFF and RMBF. All these three algorithms require that
 246 tasks are sorted according to their non-decreasing periods/request rates.
 247 They are only applicable to be applied to the situations that the periods
 248 of incoming tasks are fixed and static. The performance of these three algo-
 249 rithms are evaluated in the approximation ratio, which is defined to be the
 250 rate of N and N_0 , the number of processors required to be feasibly sched-
 251 ule a set of task by the algorithm and the minimum number of processors
 252 required to feasibly schedule the same set of tasks. As shown in the Table
 253 1, RMNF has a approximation ratio as 2.67, while RMFF and RMBF have
 254 better performance.

255 4. Global Scheduling

256 In this section, we will dive into another branch of scheduling strategies
257 – global approaches to multiprocessor real-time scheduling.

258 4.1. Overview

259 Global scheduling algorithms, as its name suggests, globally schedules
260 any feasible periodic task set. In contrast to partitioned scheduling, global
261 scheduling schedules jobs and tasks in one single shared queue instead of
262 multiple local, dedicated queues. By this means, the automatic load balanc-
263 ing and lower average response time can be achieved by global approaches.
264 In addition to this, global approaches also take advantages of the simpler
265 implementations and the existence of optimal schedulers.

266 Global scheduling strategies include Global Fixed-Job-Priority Scheduling,
267 Global Fixed-Task-Priority Scheduling, and Global Dynamic Priority
268 Scheduling. Although there are various categories of global scheduling algo-
269 rithms, the focus of this paper is on the Global Dynamic Priority Scheduling.
270 In the following subsection, it will be characterized in details.

271 4.2. Global Dynamic Priority Scheduling

272 In this subsection, we will present our in-depth exploration to the track
273 of global dynamic priority scheduling algorithm. To the best of our knowl-
274 edge, a number of global dynamic priority scheduling algorithms are optimal
275 for periodic tasksets with explicit or implicit deadlines. For example, Pro-
276 portionate Fairness algorithm and its variants including PD, PD², ERFair,
277 BF, SA [9], and LLREF [10] as well, are all optimal for offline environment.
278 Nevertheless, no algorithms until now are optimal to cope with online pre-
279 emptive scheduling problem, where tasksets are sporadic and multi-processor
280 environments are enforced. On the other hand, despite of its optimality and
281 dominance in theory, the usage of global dynamic priority algorithms are
282 limited in practice. This is because the existence of frequent preemption and
283 migration between tasks gives rise to excessive overheads in potential.

284 The following part of this subsection will provide brief summary of three
285 classic global dynamic priority scheduling algorithms. They are respectively
286 Proportionate Fairness Algorithm (PFair), and Largest Local Remaining Ex-
287 ecution First (LLREF).

4.2.1. PFair

Baruah et al [11] introduced Proportionate Fairness Algorithm. The Pfair class of algorithms that allow full migration and fully dynamic priorities have been shown to be theoretically optimal – i.e., they achieve a schedulable utilization bound (below which all tasks meet their deadlines) that equals the total capacity of all processors. Here are some fundamental principles of Proportionate Fairness algorithms:

1. Timeline is divided into equal length slots.
2. Task period and execution time are multiples of the slot size.
3. Each task receives amount of slots proportional to its task utilization.

The essential part of PFair algorithms is the quantum-based optimization defined over the lag of each task $lag(\tau_i, t)$, with the goal of minimizing the maximum lags of all tasks $\max_t |lag(\tau_i, t)|$.

$$\underbrace{lag(\tau_i, t)}_{error} = \underbrace{t \cdot \left(\frac{c_i}{T_i}\right)}_{fluid\ execution\ in\ [0, t)} - \underbrace{allocated(\tau_i, t)}_{real\ execution\ in\ [0, t)} \quad (4)$$

The generation of an optimal schedule is based on the above definition of lag . PFair algorithm does execute all urgent tasks with $lag(\tau_i, t) > 0$ and $lag(\tau_i, t + 1) \geq 0$ if τ_i executes. On top of that, PFair algorithm does not execute tnegru tasks, for which $lag(\tau_i, t) < 0$ and $lag(\tau_i, t + 1) \leq 0$ if τ_i does not execute. Besides, for other tasks, only those that have the least t such that $lag(\tau_i, t) > 0$ are executed.

The PFair algorithm will assign priorities to tasks at every time slot, which indicated itself as one of job-level dynamic priority scheduling policies. However, this characteristic gives rise to some issues, for example, frequent preemptions and frequent migrations.

Proportionate Fairness algorithm, as a strong candidate for solving resource allocation problems, has wide variety of interesting applications and powerful theoretical support. For instances, Kelly et al [12] presented its application on the problem of rate control for communication networks. Application of PFair algorithm on LANs and hoc networks was fulfilled by Jiang et al [13]. Besides, Bonald’s paper [14] provides in-depth queueing analysis over proportionate fairness as well as max-min fairness and balanced fairness.

318 4.2.2. *LLREF*

319 LLREF was firstly introduced by Cho et al [10]. Similar to PFair class
320 of algorithms, LLREF is also based on the fluid scheduling model, where
321 each task executes at a constant rate at all times. The principal idea of
322 LLREF is that given M processors, M largest local remaining execution
323 time tasks are selected first for every secondary event. This is also called
324 the LLREF scheduling policy. To reason over task execution behavior on
325 multiprocessors, a novel abstraction called Time and Local Execution Time
326 Domain Plane (T-L Plane) was developed.

327 This algorithm divides the schedule into Time and Local execution time
328 planes (TL-planes), which are determined by task deadlines. The algorithm
329 schedules tasks by creating smaller local jobs within each TL-plane. The
330 only parameters considered by the algorithm during a TL-plane are the pa-
331 rameters of the local jobs. When a TL-plane completes, the next TL-plane
332 is started. The duration of each TL-plane is the amount of time between
333 consecutive deadlines.

334 4.3. *Summary*

335 The global scheduling paradigm has advantages over the partitioned ap-
336 proach. First of all, if tasks can join and leave the system at run-time,
337 then it may be necessary to reallocate tasks to processors in the partitioned
338 approach. In addition, the partitioned approach cannot produce optimal
339 real-time schedules – one that meets all task deadlines when task utiliza-
340 tion demand does not exceed the total processor capacity – for periodic task
341 sets, since the partitioning problem is analogous to the bin-packing problem
342 which is known to be NP-hard in the strong sense. On top of that, in some
343 embedded processor architectures with no cache and simpler structures, the
344 overhead of migration has a lower impact on the performance. Finally, global
345 scheduling can theoretically contribute to an increased understanding of the
346 properties and behaviors of real-time scheduling algorithms for multiproces-
347 sors.

348 However, the global scheduling paradigm has also several disadvantages.
349 Firstly, global scheduling strategies are much more complicated to implement
350 than partitioned scheduling. In other words, for the partitioned approach,
351 once a set of tasks are allocated to processors, the multiprocessor real-time
352 scheduling problem becomes a collection of single processor real-time schedul-
353 ing problems. The ease of programming partitioned scheduling is obvious
354 since the single processor scheduling problem has already been well-studied

355 and optimal algorithms with easy implementations already exist. Secondly,
356 migrating tasks at run-time means more runtime overhead in that migrating
357 tasks may suffer cache misses on the newly assigned processor. If the task set
358 is fixed and known in advanced, it is obvious that the partitioned approach
359 provides more appropriate solutions.

360 **5. Hybrid Approaches**

361 **6. Conclusions**

References

- [1] C. L. Liu, J. W. Layland, Scheduling algorithms for multiprogramming in a hard-real-time environment, *Journal of the ACM (JACM)* 20 (1973) 46–61.
- [2] L. Sha, R. Rajkumar, J. P. Lehoczky, Priority inheritance protocols: An approach to real-time synchronization, *Computers, IEEE Transactions on* 39 (1990) 1175–1185.
- [3] A. Mohammadi, S. G. Akl, Scheduling algorithms for real-time systems, School of Computing, Queens University (2005).
- [4] R. I. Davis, A. Burns, A survey of hard real-time scheduling for multiprocessor systems, *ACM Computing Surveys (CSUR)* 43 (2011) 35.
- [5] M. R. Garey, D. S. Johnson, *Computers and intractability: a guide to the theory of np-completeness*. 1979, San Francisco, LA: Freeman (1979).
- [6] J. Carpenter, S. Funk, P. Holman, A. Srinivasan, J. Anderson, S. Baruah, A categorization of real-time multiprocessor scheduling problems and algorithms, *Handbook on scheduling algorithms, methods, and models* (2004) 30–1.
- [7] S. K. Dhall, C. Liu, On a real-time scheduling problem, *Operations research* 26 (1978) 127–140.
- [8] Y. Oh, S. H. Son, Tight performance bounds of heuristics for a real-time scheduling problem (1993).
- [9] A. Khemka, R. Shyamasundar, An optimal multiprocessor real-time scheduling algorithm, *Journal of parallel and distributed computing* 43 (1997) 37–45.
- [10] H. Cho, B. Ravindran, E. D. Jensen, An optimal real-time scheduling algorithm for multiprocessors, in: *Real-Time Systems Symposium, 2006. RTSS’06. 27th IEEE International, IEEE*, pp. 101–110.
- [11] S. K. Baruah, N. K. Cohen, C. G. Plaxton, D. A. Varvel, Proportionate progress: A notion of fairness in resource allocation, *Algorithmica* 15 (1996) 600–625.

- 393 [12] F. P. Kelly, A. K. Maulloo, D. K. Tan, Rate control for communication
394 networks: shadow prices, proportional fairness and stability, Journal of
395 the Operational Research society (1998) 237–252.
- 396 [13] L. B. Jiang, S. C. Liew, Proportional fairness in wireless lans and ad hoc
397 networks, in: Wireless Communications and Networking Conference,
398 2005 IEEE, volume 3, IEEE, pp. 1551–1556.
- 399 [14] T. Bonald, L. Massoulié, A. Proutiere, J. Virtamo, A queueing anal-
400 ysis of max-min fairness, proportional fairness and balanced fairness,
401 Queueing systems 53 (2006) 65–84.