

Synthesis of Safe, QoS Extendible, Application Specific Schedulers for Heterogeneous Real-Time Systems *

Christos KLOUKINAS

Sergio YOVINE

VERIMAG, Centre Équation, 2 avenue de Vignate, 38610 GIÈRES, France

E-mail: Christos.Kloukinas@imag.fr & Sergio.Yovine@imag.fr

Abstract

We present a new scheduler architecture, which permits adding QoS policies to the scheduling decisions. We also present a new scheduling synthesis method which allows a designer to obtain a safe scheduler for a particular application. Our scheduler architecture and scheduler synthesis method can be used for heterogeneous applications where the tasks communicate through various synchronisation primitives. We present a prototype implementation of this scheduler architecture and related mechanisms on top of an open-source OS for embedded systems.

1. Introduction

Safety & mission-critical systems need to be of extremely high quality, due to the great dangers and the high cost of their potential failure. For this reason, when they are multi-threaded they must be guaranteed to be free of deadlocks and all threads must be guaranteed to meet their deadlines under all circumstances.

The current practice for avoiding deadlocks is to use the *immediate priority ceiling protocol* (IPCP) [8] for the sharing of non-preemptable resources. This approach has a certain number of disadvantages though. All the priority inheritance family of protocols are pessimistic in nature and, therefore, can refuse access to a shared resource even when there is no real danger of a deadlock at the current situation. Additionally, in order to apply the IPCP, the designer of the system must choose a set of priorities for all the tasks in the system. Furthermore, for each shared resource the designer must identify the threads which use it, in order to assign a priority to that resource (*i.e.*, the *ceiling* priority of the resource). More importantly, however, the IPCP cannot on its own support tasks which synchronise using monitors and communicate using condition variables, as for example is done in JAVA [3]. In such cases one should split tasks, which means that the complexity of what the designers have

to do in order to use IPCP is quite high. Worse yet, not all cases of communication through condition variables can be translated according to these rules, since there can be tasks which wait on a condition variable while still holding some resource locked. These tasks cannot be split since the underlying assumptions of the IPCP clearly demand that tasks finish executing without holding any resources. Therefore, it is not straightforward how one can use the IPCP with a language such as *R-T JAVA* [9].

Besides, the methods currently used do not allow designers to easily extend them for incorporating QoS to the scheduler decisions. Being able to extend a scheduler with QoS characteristics could allow us to experiment with ways to minimise energy consumption, or further increase the speed of the system, by minimising, for example, the number of context switches. An example of this can be found in [2] where the authors present a dynamic scheduling method which also treats the QoS aspect of the system. However, in this work the authors consider a fixed task model where all tasks are periodic and do not consider deadlock situations or the communication aspect of the system.

In the following, we present a method for synthesising *QoS extendible* and *safe* schedulers following the controller synthesis paradigm and continuing previous work at Verimag [1]. We start in section 2 by presenting the overall architecture of the scheduler. Then, in section 3 we present the model of the systems we consider and in section 4 the particular method we use for the synthesis of the scheduler. In section 5 we present a prototype implementation of our scheduler on top of a real operating system and we conclude in section 6.

2. Scheduler Architecture

We consider a set of threads synchronising through monitors and communicating through condition variables. All threads and shared objects are created at the initialisation phase. We only consider applications executing on a single processor for the moment. The architecture of the schedulers we synthesise consists of two three-layered stacks, as

*Partially funded by the French RNTL project Expresso.

shown in Figure 1.

The left stack is responsible for selecting an application thread for execution. The right stack is responsible for selecting an application thread for the reception of a notification. This stack allows the scheduler to control the communication aspect of the system. Being able to control which thread will be notified for a particular event is something that other scheduling policies like the PCP do not offer, since they concentrate only on the selection of threads for execution.

After one of the scheduler stacks is finished, it passes control to an underlying *R-T OS* which provides low-level kernel mechanisms. Such mechanisms include the ability to create, suspend and resume an application thread, as well as the ability to create, set and disable alarms for future events (e.g., arrival of next period or the timeout of a **waitTimed**).

2.1. Controlling the Currently Executing Thread

The left stack takes control of the system when the application calls one of **monitorEnter**, **monitorExit**, **waitForPeriod**, **wait** and **waitTimed**, or when an alarm expires. In these cases, it must choose one of the available threads as the thread which should next be run on the processor. It does this in three steps, each one performed at a different layer.

In the first layer, referred to as the *Ready-Exec* scheduler, it calculates the set of threads \mathcal{R}_{exec} which are ready to execute without directly blocking due to mutual exclusion. That is, it examines whether an application thread will try to enter a monitor which is already occupied by another thread, if it is chosen as the next thread to execute. Having calculated the set \mathcal{R}_{exec} , this layer passes it to the next layer.

The *Safe-Exec* scheduler layer is responsible for calculating the subset \mathcal{S}_{exec} of \mathcal{R}_{exec} , consisting of those threads that can *safely* execute. Safety here refers both to deadlock freedom (i.e., entering a monitor would not cause a deadlock later on), as well as, to meeting the timing constraints of the different threads (i.e., choosing a thread for execution will not delay another thread enough to make it miss its deadline).

The *Safe-Exec* layer passes the set \mathcal{S}_{exec} to the third layer *QoS-Exec*, which calculates the set $\mathcal{Q}_{exec} \subseteq \mathcal{S}_{exec}$, consisting of the *safe* threads which respect the QoS requirements.

2.2. Controlling the Notified Thread

The right stack is passed control when the application calls **notify** or **notifyAll**. The reason for this is that the threads which will be notified (if any) cannot ever be selected for execution. This is because they will immediately try to re-enter the monitor after being notified and thus get blocked by the notifier (which is already in the monitor). Nevertheless, when a thread notifies a condition variable, then we can control which among the threads waiting

for the notification should receive the event. Indeed, languages (like JAVA [3]) which provide the monitor construct, or thread libraries (like POSIX threads [6]) which offer it to languages which do not provide it, leave this point *unspecified*, allowing each implementation to choose the thread to be notified as it convenes it the best.

The top layer *Ready-Notif* calculates the set \mathcal{R}_{notif} of threads waiting on the condition variable being notified. The apparent cases where we cannot effect any control on the system are three. First, the case where no thread waits on the condition variable being notified. Second, the case where only one thread waits on the condition variable. Third, the case where the notifying thread does a **notifyAll**, in which case we are obliged to notify all the threads waiting on the condition variable. In these cases, the right stack of the scheduler does not make any control decision, but simply *changes the PC* of all threads waiting on the current notification, that is, the threads belonging to the set \mathcal{R}_{notif} , to mark them as notified. Otherwise, the top layer passes the set \mathcal{R}_{notif} to the middle layer *Safe-Notif*, which calculates the subset \mathcal{S}_{notif} of \mathcal{R}_{notif} consisting of those threads which, if notified, will not cause the system to enter into a deadlock state or cause some other thread to miss its deadline (i.e., they are safe). Finally, the *Safe-Notif* layer passes the \mathcal{S}_{notif} set to the bottom *QoS-Notif* layer, which calculates the subset \mathcal{Q}_{notif} of \mathcal{S}_{notif} , consisting of the threads which we can safely notify and also respect the QoS properties of the application. The *QoS-Notif* layer is also responsible for choosing one of the threads in \mathcal{Q}_{notif} as the recipient of the current notification and marking it as notified by changing its *PC*.

2.3. QoS Policies

By incorporating a layer for QoS in the scheduler, we offer an *extendible* mechanism for providing additional properties to the system. We allow the QoS policies to use the same information which is available to the scheduler. That is, the QoS layer has access to the *program counters (PC)* of the threads, the *currently executing thread* (T_{Exec}), as well as the value of the *system clock* (C_{System}).

The complexity of the QoS layer is controlled by the application designer. In choosing a QoS policy (or policies, since these are composable) the designer can balance between the execution time and extra memory space needed by the policy and the gains to the overall system quality the particular policy can offer. A QoS policy is, for example, the *local minimisation of context switches* in order to speed-up the execution. This policy can be implemented quite easily, since all one needs to examine is whether the currently executing thread T_{Exec} is in the set \mathcal{S}_{exec} of threads which are safe to execute next. If this is the case, then we can let it continue its execution, by setting the set \mathcal{Q}_{exec} equal to the singleton $\{T_{Exec}\}$. This particular policy has another ad-

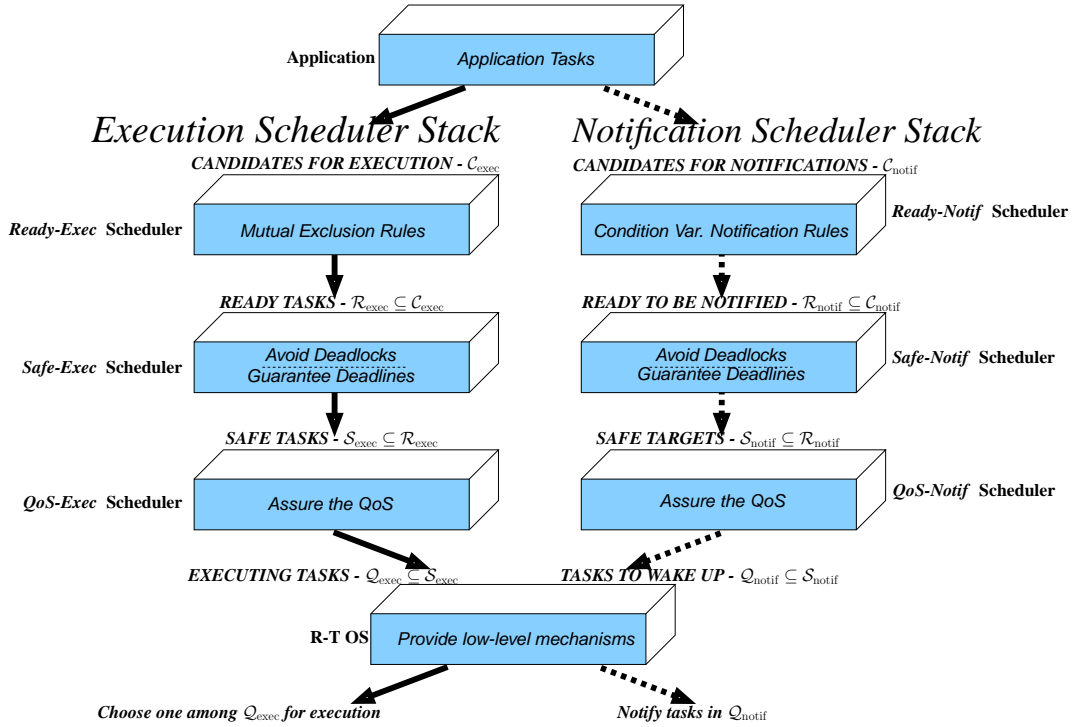


Figure 1. A three-layered scheduler architecture

vantage: by decreasing the number of context switches, we also decrease the cache misses of the application, since now there are fewer points in the execution where the threads compete for the cache, potentially flushing each other's data out of it. This can help *decrease the energy consumption* of the system, since a cache miss can lead to two main memory accesses, which are known to be quite demanding with respect to energy [7]. In fact, since a cache reads and flushes one *cache line* at a time (*i.e.*, multiple consecutive memory addresses) the benefits can be even greater, both with respect to energy consumption and execution speed.

3. System Model

The model of the system we construct is the parallel composition of an automaton which is responsible for advancing time and firing the alarms, one automaton for each of the application threads and two more automata, for the *QoS-Exec* and the *QoS-Notif* scheduler layers respectively. The automaton of time and the automata of the application threads perform a finite number of actions and then block, letting the scheduler automata respond. The actions of the time and application automata being *uncontrollable*, the only *controllable* actions are those of the two scheduler automata. Thus, our model can be seen as a two player game with the scheduler automata on one side (*i.e.*, the controller) and the time and application thread automata on the other (*i.e.*, the plant). In this game, the automata related

to the application simulate the locking and unlocking of resources, as well as, the waiting and notification on condition variables. The computations performed by the application threads are simulated just by their minimum and maximum execution times. Each statement s of an application thread (where s is one of **monitorEnter**, **monitorExit**, **wait**, **waitTimed**, **waitForPeriod**, a conditional or a computation) is modelled by a separate automaton state and a transition from it to the next statement position ($@s'$) which is taken when the statement s can be executed. The only exception to this rule is the case of the **wait** and **waitTimed** statements. These statements are effectively modelled by two states; the first one models the release of the mutex associated with the condition variable on which we wait and the second one models the attempt of the thread to reacquire the mutex, once it has been notified. The advancement of time is the responsibility of a single automaton (see Figure 2-a) which, in addition, enables transitions in the application automata which correspond to timeouts, such as in the case of a **waitTimed** or a **waitForPeriod**. This automaton is also responsible for advancing the local thread clocks, that is, the clocks which model the time spent by threads in computations (and in waiting when doing a **waitTimed**). These clocks are set to zero at the beginning of a computation by a thread and are incremented alongside with the global time, until the duration of the computation is over (or the timeout of the **waitTimed** has expired). The two automata for the *QoS-Exec* and *QoS-Notif* scheduler layers

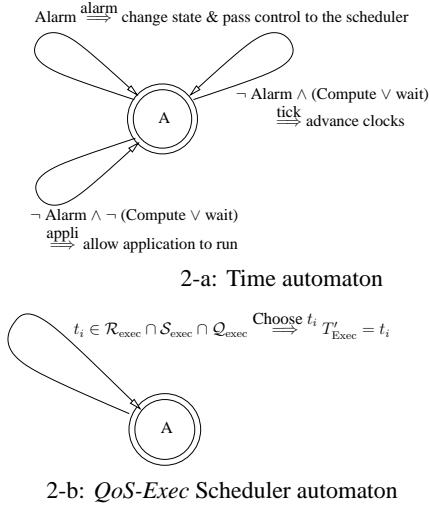


Figure 2. Time & Scheduler automata

are passed control as described in sub-sections 2.1 and 2.2 and decide which of the application automata should be allowed to execute next or be notified of an event. These automata are comprised of a single state and $n + 1$ transitions, where n is the number of application threads (the additional transition being for the *idle* thread). A transition of these automata selects one of the threads for execution (resp. for notification). It is guarded by a predicate which asserts that the corresponding thread belongs to the Q_{exec} (resp. Q_{notif}) set (see Figure 2-b). The transition choosing the idle thread asserts that the rest of the threads are not safe to execute (resp. no thread waits to be notified on the current event).

To summarise, the state in our model comprises of: (i) a program counter (PC_i) for each of the application threads, (ii) a local clock (C_i) for each thread which is used for their computations and the timeouts if they execute a **waitTimed**, (iii) a global clock (C_{System_i}) for modelling the periods of each periodic thread, (iv) a variable (T_{Exec}) holding the currently executing thread, (v) two boolean variables ($Exec_Sched_Enabled$ & $Notif_Sched_Enabled$) for controlling whether it is one of the scheduler automata (and which of them), or the time (when they are both true) or the time and application automata (when they are both false) which should execute, and (vi) the *boolean* variables of the application threads used in conditionals which are *significant*¹ with respect to the use of resources and communication of events. Our model goes through three different modes of execution, as shown in Figure 3. In the “Time Only” mode (where $Exec_Sched_Enabled = Notif_Sched_Enabled = \mathbf{true}$) the time automaton is the sole automaton enabled in the system and it can fire one or more alarms, if any is enabled. If an alarm is fired then the execution mode changes to “Schedulers Only” (where

¹We perform *slicing* of the original code to identify these variables.

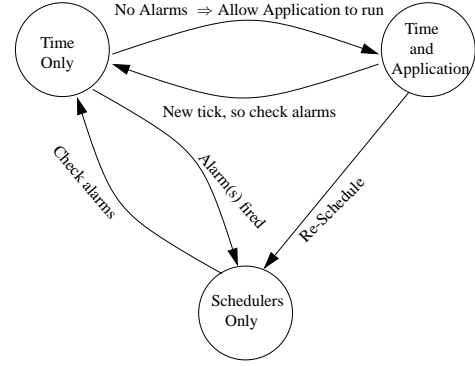


Figure 3. Model execution modes

$Exec_Sched_Enabled = \neg Notif_Sched_Enabled$), so that our scheduler can treat the alarm. If there is no alarm to be fired then the execution mode changes to “Time and Application” (where $Exec_Sched_Enabled = Notif_Sched_Enabled = \mathbf{false}$). At this mode, both the time automaton and the automata of the application are enabled. If the time automaton gets to execute first, then a tick (*i.e.*, a time step) is performed and we pass back to the “Time Only” mode, so as to check if an alarm is now enabled. If it is one of the application automata which gets to execute first, then it executes until it needs to perform an action which causes re-scheduling, in which case it passes control to the schedulers (*i.e.*, the mode now becomes “Schedulers Only”). If the application automaton needs to execute a time guarded action (*i.e.*, a computation), then it blocks, allowing time to advance.

As an example, let us consider the model shown in Figure 4. Here, the application consists of three threads, one of which is a *periodic* one (the *User*) and two *aperiodic* ones (the *Writer* and the *Refresher*). One should note that the *Writer* and the *Refresher* are continually enabled aperiodic threads and do not have any deadlines directly associated with them.

4. Scheduler Synthesis

In order to synthesise the *Safe-Exec* and *Safe-Notif* scheduler layers, we first construct the set of reachable states and, thus, identify the deadlocks. These are the states where the application threads are deadlocked, or the states where some thread has missed its deadline (since in that case we block the system explicitly). The existence of these states indicates that the predicate we are currently using to describe the set *Safe-Exec* (resp. *Safe-Notif*) needs to be constrained even further. This predicate starts with the value of **true**, thus accepting initially all threads in the set *Ready-Exec* (resp. *Ready-Notif*) as safe. Having obtained the deadlocked states, we do a backwards traversal of the whole state space starting from the deadlocked states, until we reach a

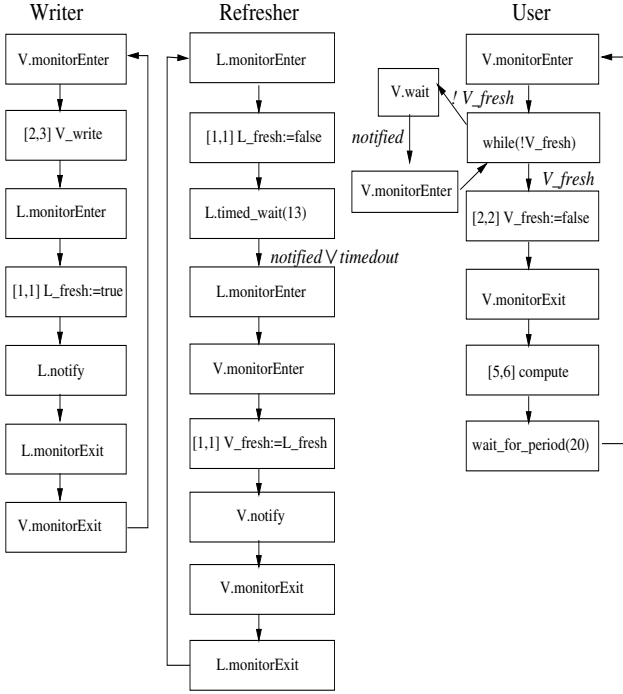


Figure 4. Application automata

(Clock variables are omitted for readability – each computation is annotated with its duration interval.)

state which corresponds to a choice of one of the scheduler automata. There, we identify the choice $T_{Exec} = t_i$ which allowed the path leading to a deadlock state(s) and create a new constraint for the layer *Safe-Exec* (resp. *Safe-Notif*). This constraint is constructed by changing the set \mathcal{S}_{exec} (resp. \mathcal{S}_{notif}) to be:

$$\mathcal{S}'_{exec}(\overrightarrow{state}) = \mathcal{S}_{exec}(\overrightarrow{state}) \setminus \{t_i\}$$

If at some point we find that $\mathcal{S}'_{exec}(\overrightarrow{state})$ is equal to the empty set, then we add the current state to the set of deadlocked states and continue the synthesis procedure.

4.1. State Space Reduction & Application Analysis

Even though the basic idea of synthesising the *Safe-Exec* and *Safe-Notif* scheduler layers is simple, it is evident that in practice it suffers from the state explosion problem. Therefore, it is imperative that we use techniques to minimise the size of the state space. Altisen *et al.* [1] proposed to constrain the system with a high level policy (*i.e.*, FIFO scheduling, EDF scheduling, *etc.*) the idea being that this will constrain enough the system to allow us to construct the entire (constrained) state space of it. However, this method can sometimes over-constrain the system and remove all possible paths which would allow us to avoid the unsafe states. Besides, it is not always clear how one can apply policies such as EDF, RMA, *etc.* when the application consists of heterogeneous threads. Our method

consists of synthesising schedulers for successively more detailed models, adding thus complexity to a model only when we have already calculated how we can constrain the more abstract one. The scheduler synthesis is performed in two major steps. In the first one, we examine the *untimed* model of the system and search for constraints which can guarantee the *absence of deadlocks*. In the second one, we re-introduce time into the model and after constraining it for avoiding deadlocks, we search for those constraints which can guarantee that *all threads meet their deadlines*.

4.2. Deadlock Avoidance

Examining the untimed model of the system first, has the disadvantage that some of the deadlocks we identify are not possible in reality, due to the existing timing relations. On the other hand, adding time to a model significantly increases its size and thus renders the analysis and synthesis a lot more difficult. Thus, searching for deadlocks in the untimed model allows us to examine a much smaller search space and thus allows us to attack larger systems. For instance, the untimed model of the application shown in Figure 4 is 97% smaller than the timed one and it allowed us to discover 8 constraints which can help us avoid all the 10 deadlocks caused by the use of shared resources. More importantly, finding and removing *all* deadlocks in the untimed model means that the application is *logically* correct and allows a designer to experiment with different underlying platforms and algorithms for implementing the application computations.

4.3. Guaranteeing Deadlines

Having found all the potential deadlocks in the system, we add the synthesised \mathcal{S}_{exec} and \mathcal{S}_{notif} scheduler sets obtained so far to the timed version of the initial model, in order to search for the *timeliness* constraints, which can guarantee that all threads will meet their deadlines. In order to make the problem more tractable, our synthesis algorithm works in three steps.

Safety-preserving abstraction Not all time instances are visible when the scheduler automata take control during the execution of the system. That is, there are certain states of a system where the only event allowed is the advancement of time. In such cases, we do not really gain anything by explicitly constructing the complete sequence of all the different time steps. Instead, we can jump directly to the last state of this sequence, where the time has advanced enough to allow some other event to occur. This state space reduction can be effectively obtained through the *branching bisimulation equivalence (bbe) reduction* [11], which eliminates “unobservable” actions (in our case the Tick action) but only when doing so preserves the branching structure of processes. The preservation of the branching structure of

the application is crucial for us, since the synthesis of the scheduler depends on it for calculating the states where a controllable action can help avoid taking a path which leads to an undesired state. Given a set of equivalent states under the *bbe* reduction, we elect as a representative of this set the state which has the maximum global clock value. In other words, in the *bbe* reduced system, our scheduler takes its decisions at the latest moment possible. For the application shown in Figure 4 we obtain a 74% reduction of the state-space.

Non Preemptable Threads We then consider that the application threads cannot be *preempted* while they are computing. The non-preemption hypothesis reduces the state space, since it removes all the cases where the execution of a thread is suspended by an interrupt (*e.g.*, for starting a new period of some other thread). In the application of Figure 4 the reduction obtained by not allowing preemption is 40% when applied on the initial timed model and 80% when applied to the *bbe* reduced timed model. Once we can indeed safely schedule the system under the hypothesis that threads are never preempted, then we can use the constraints obtained during this step to *reduce even further* the state space that we have to construct and analyse when we do allow threads to be preempted. Indeed, for the application of Figure 4, we can reduce the state space by an additional 10% with the 30 constraints we construct during this step.

The non-preemption of threads is easily added to our models *through the use of a QoS policy* that forbids the schedulers from choosing a thread for execution, when another thread is already in a state where it is performing a computation :

$$\mathcal{Q}_{\text{exec}}(\overrightarrow{\text{state}}) = \{t . t \in \mathcal{S}_{\text{exec}}(\overrightarrow{\text{state}}) \wedge \neg \exists t' \neq t . \text{computes}(t')\}$$

However, we cannot safely schedule all systems when we do not allow threads to be preempted. This means that for these systems we will not obtain any scheduling constraints and, therefore, will be obliged to examine the large, unconstrained state space of the timed model.

Allowing Preemption Having performed the scheduler synthesis for the case where threads are not preemptable, we add the additional constraints we synthesised (if any) to the model and perform the scheduler synthesis once more, this time allowing threads to be preempted. This is the final step of the scheduler synthesis, which provides us with the whole set of constraints that we must impose on the application in order to guarantee that it will be deadlock free and that it will meet all the deadlines of the threads. For the application of Figure 4, this last synthesis step produces

an additional 18 constraints and thus we can safely schedule this application with a total of 56 constraints, avoiding both deadlocks and missed deadlines. These 56 constraints are all part of the *Safe-Exec* layer, since in this application there is always at most one thread waiting to be notified on a particular condition variable and thus we cannot control the communication aspect of the application. This final *safe* timed model has a state space which is 96% smaller than the original, unsafe one.

4.4. Not Observing Clocks

Having synthesised a safe scheduler for an application does not necessarily mean that we can implement it easily on a usual *R-T OS* though. The difficulty of implementing it as is, arises from the fact that the constraints we produce during the synthesis use the state of the system to decide what are the safe choices at each point during the execution and, therefore, also make reference to the values of the local clocks of the threads. However, these clocks do not really exist in the application but were only introduced as a way to model the computations of the threads. Introducing them means that we will have to add for each thread an additional timer object and reset and activate (*resp.* reset and deactivate) the timer before (*resp.* after) each computation and read its value when making a scheduling decision. As using timers may substantially increase the execution time of the scheduler, we investigate the possibility of synthesising a clock-free one, which only examines the *PCs* of the threads. This will make the scheduler itself faster to execute, since in order to make a scheduling decision it now only needs to examine the n values of the different *PCs* and not the $2n + 1$ values of the *PCs*, the local thread clocks and the global clock.

On the other hand, removing the clocks from the constraints can introduce states where the scheduler will take the wrong decision and cause a thread to miss its deadline. These states are those where a scheduler gets called at the same configuration of thread *PCs* but at different time instances. Since the time instance (and therefore the clock values) are different, the safe sets $\mathcal{S}_{\text{exec}}$ of these states can be different themselves as well. When we decide to not observe the clock values while scheduling, we are effectively unable to differentiate among these different sets and all these states become *equivalent*, as far as our scheduler is concerned. Therefore, if we wish our scheduler to always make a decision which is *safe*, then the $\mathcal{S}_{\text{exec}}$ set of this *equivalence class* of states should be the *intersection* of the $\mathcal{S}_{\text{exec}}$ sets of the states which belong to the same equivalence class.

$$\mathcal{S}_{\text{exec}}(\overrightarrow{\text{class}}_j) = \bigcap_{\text{state}_i \in \text{class}_j} \mathcal{S}_{\text{exec}}(\overrightarrow{\text{state}}_i)$$

Sometimes, the $\mathcal{S}_{\text{exec}}(\overrightarrow{\text{class}}_j)$ set will be empty, if the

scheduler decisions at the members of this equivalence class were conflicting. When encountering such an equivalence class whose $\mathcal{S}_{\text{exec}}$ set is the empty set, we need to add its members to the set of deadlocked states and continue the synthesis algorithm, until we find a set of constraints which helps us to avoid the whole class.

The example of Figure 4 is indeed an application which can be scheduled without observing the clock values. However, in some cases it may be impossible to safely schedule the application without taking into account the values of the clocks. Thus, there is a trade-off between schedulability and execution cost of the scheduler.

4.5. Other QoS Policies

Once we have synthesised a safe scheduler, we can compose it with other QoS policies to choose among the safe threads those which better realise the QoS requirements of the system. However, it can always be the case that some of these policies can cause the application to miss a deadline, by choosing no thread for execution (*i.e.*, setting the $\mathcal{Q}_{\text{exec}}$ set to be the empty set). For this reason we must verify them and change them if they can indeed cause the application to miss a deadline. Since the verification is performed on the *safely schedulable* application, the size of the state space we must explore is quite small. For the example application of Figure 4, we verified two additional QoS policies. The first is a fixed priority QoS policy, where $\text{User} > \text{Refresher} > \text{Writer}$ when they are *safe* and it is effectively verified in a model 98.4% smaller than the original timed one. The second is a QoS which (locally) minimises the number of context switches and it is verified on a model which is slightly even smaller (a 98.8% reduction).

5. Scheduler Implementation

In this section we present the implementation of our scheduler over an *OS* for embedded systems. Our current implementation does not make use of priorities, mutexes or condition variables of the underlying system. It rather uses *thread suspension and resumption* to simulate these mechanisms, which allows us to provide an implementation of our approach even on operating systems which have a very limited number of priorities or have no priorities whatsoever.

The actions which activate our scheduler in order to elect a new thread to execute (**monitorEnter**, **monitorExit**, **waitForPeriod**, **wait** and **waitTimed**, and the expiration of a timeout), are implemented similarly. First, the *OS* scheduler is locked so as to avoid interrupt handlers from changing the system state. Then our scheduler examines the current state of the system, *i.e.*, the *PCs* of all the threads, and decides which should be the next thread to execute, by using the synthesised sets $\mathcal{R}_{\text{exec}}$ & $\mathcal{S}_{\text{exec}}$ and the user-provided set $\mathcal{Q}_{\text{exec}}$. Finally, our scheduler suspends all threads, resumes the one it has chosen for execution and returns after

having unlocked the *OS* scheduler, thus re-allowing interrupts to occur. The only case which is treated differently, is the case where an interrupt arrives to signal a timeout (either for a **waitTimed** or a **waitForPeriod**). In this case, the associated interrupt handler changes the *PC* of the respective thread to mark it as no longer waiting and, then, suspends the currently executing thread (if any) and resumes the thread that was waiting. Once this thread starts to execute, it calls our scheduler in its turn, to assure that it can safely continue. If this is the case, our scheduler will allow it to execute, otherwise it will suspend it and resume another thread. It is easy to see that if more than one interrupts arrive at the same time, our scheduler will be called consecutively more than once. However, we proved that there can be no more than 3 consecutive calls to our scheduler ever.

The actions which do not cause our scheduler to elect a new thread for execution (**notify** and **notifyAll**) are implemented as follows. These actions make use of its right stack, which deals with the communication aspects of the system. When a **notify** occurs, the scheduler is called and it checks whether there are any threads waiting on the event notified. If so, it selects one of them *using the $\mathcal{Q}_{\text{notif}}$ set*, marks it as notified and gives back the execution to the thread which did the **notify**.

We have implemented our architecture on top of *eCos* [5], an open-source *OS* for embedded systems. We run our test-bed application on a 330 MHz Pentium II machine using *synthetic-Linux* as the execution platform of *eCos*, which means that *eCos* and the application are running as a single Linux process. Our experiments showed that the application did indeed honour its deadlines as expected. Besides, the measurements showed an average execution time of our scheduler of the order of 0.66 microseconds. Given the fact that our prototype implementation is not particularly optimised for speed, this is a rather small execution cost.

6. Conclusions

We have presented a new methodology for building application-driven schedulers for heterogeneous systems and a prototype implementation of our scheduler using an open-source *OS* for embedded systems. Our work continues the one described in [1], where a thorough discussion about related approaches is presented. Applications comprising of heterogeneous thread types have also been considered in [4], without taking into account thread interdependencies due to the sharing of non-preemptable resources.

The advantages of our synthesis method is that we can handle larger models than if we would have tried to attack the original timed version of the model at once. In addition, following our method a designer is better able to understand the behaviour of a system, since we successively drive him through: (i) the states which cause a deadlock later on, (ii)

Table 1. Experimental results

Model kind	States	Red.	Trans.	Red.	Dead.	Red.	Constraints
<i>Model Abstractions & Optimisations</i>							
T <i>original (i.e., Preemption)</i>	45470	0.00%	48786	0.00%	367	0.00%	—
U	1352	97.03%	1645	96.63%	10	97.28%	—
T <i>No Preemption</i>	27266	40.04%	29118	40.31%	134	63.49%	—
T <i>Preemption, bbe reduction</i>	11437	74.85%	13648	72.02%	1	99.73%	—
T <i>No Preemption, bbe reduction</i>	8648	80.98%	10038	79.42%	1	99.73%	—
<i>Synthesis Steps</i>							
U	1352	97.03%	1645	96.63%	10	97.28%	0
U, No Deadlocks	1200	97.36%	1451	97.03%	0	100.00%	8
T <i>No Preemption, bbe reduction, No Deadlocks</i>	8642	80.99%	10027	79.45%	1	99.73%	8
T <i>No Preemption, bbe reduction, Safe</i>	1542	96.61%	1668	96.58%	0	100.00%	38
T <i>Preemption, bbe reduction, No Clocks</i>	4640	89.80%	5532	88.66%	1	99.73%	38
T <i>Preemption, bbe reduction, No Clocks, Safe</i>	1593	96.50%	1740	96.43%	0	100.00%	56
<i>QoS Policies (reduced with bbe)</i>							
T Safe, Fixed Priorities	728	98.40%	750	98.46%	0	100.00%	56
T Safe, Locally Min. Context Switches	549	98.79%	573	98.83%	0	100.00%	56

the states where a system is overloaded (and, therefore, he needs to allow preemption of threads), and finally, (iii) the states where the scheduler also needs to observe the values of the local clocks measuring the duration of each computation of the threads. Our method can be applied to applications comprising of any mix of periodic, aperiodic, *etc.* threads, which share non-preemptable system resources and communicate through condition variables. Finally, the robustness of the synthesised scheduler with respect to the assumed execution times of the computations can be verified by enlarging the corresponding intervals. For example, this allows checking for the case where computations take less time than specified in the model.

A disadvantage of our method is that we must build the entire state space before we can synthesise a scheduler for an application. We plan to address this problem in future versions of our tools, which will perform the synthesis in an on-the-fly manner while constructing the state-space, as for example was done in [10]. We also plan to study ways to perform the synthesis symbolically, without explicitly constructing the state space graph.

In this article we focused on models instead of a particular programming language. Such models can be extracted from programs using static analysis techniques. We indeed plan to develop such a model extraction for Java, so as to be able to schedule real-time Java programs.

References

- [1] K. Altisen, G. Göbner, and J. Sifakis. Scheduler modeling based on the controller synthesis paradigm. *Real-Time Systems*, 23(1):55–84, July 2002.
- [2] G. C. Buttazzo, G. Lipari, M. Caccamo, and L. Abeni. Elastic scheduling for flexible workload management. *IEEE Trans. Computers*, 51(3):289–302, Mar. 2002.
- [3] J. Gosling, B. Joy, and G. Steele. *The Java Language Specification*. Addison-Wesley, Reading, MA, USA, 1996.
- [4] D. Isović and G. Fohler. Efficient scheduling of sporadic, aperiodic, and periodic tasks with complex constraints. In *IEEE RTSS'00*, Orlando, Florida, USA, Nov. 2000.

- [5] A. J. Massa. *Embedded Software Development with eCos*. Prentice Hall, 2002.
- [6] The Open Group. *The Single UNIX Specification, Version 2: Threads*, 1997. (www.unix-systems.org/single_unix_specification_v2/xsh/threads.html)
- [7] F. Parain, M. Banâtre, G. Cabillie, T. Higuera, V. Issarny, and J.-P. Lesot. Techniques de réduction de la consommation dans les systèmes embarqués temps-réel. TR-1332, IRISA, France, May 2000. (In French).
- [8] R. Rajkumar, L. Sha, and J. P. Lehoczky. An experimental investigation of synchronisation protocols. In *IEEE Workshop on Real-Time Operating Systems & Software*, 1989.
- [9] Real-Time for Java Expert Group. The real-time specification for Java. Tech. Report, RTJ.org, Dec. 2001.
- [10] S. Tripakis and K. Altisen. On-the-fly controller synthesis for discrete and dense-time systems. In *FM'99*, volume 1708 of *LNCS*, Toulouse, France, Sept. 1999. Springer-Verlag.
- [11] R. J. van Glabbeek and W. P. Weijland. Branching time and abstraction in bisimulation semantics. *JACM*, 43(3), 1996.

A. Experimental Results

Table 1 shows the results of the experiments with the case study. The state-space construction and reduction have been carried out using the CADP tools². Section “*Model Abstractions & Optimisations*”, presents the state-space reductions obtained at each step. *T* indicates a *Timed* and *U* an *Untimed* model. Section “*Synthesis Steps*” shows how the model sizes change during the scheduler synthesis process. The attribute “No Clocks” refers to the fact that the synthesised scheduler does not observe the values of the clocks. The column “Constraints” reports the number of constraints *applied* to the model, that is the number of constraints synthesised in the *previous* step(s). Section “*QoS Policies*” shows the size of the *safe* model, when we apply to it different QoS policies. The first one uses the fixed priority User > Refresher > Writer. The second one locally minimises the number of thread context switches, by selecting the previously executing thread if it is still safe to do so.

²<http://www.inrialpes.fr/vasy/cadp/>