# CISTER

# Conference Paper

## Hard real-time multiprocessor scheduling resilient to core failures

**Borislav Nikolic**

**Konstantinos Bletsas and Stefan M. Petters**

# Hard real-time multiprocessor scheduling resilient to core failures

Borislav Nikolic, Konstantinos Bletsas and Stefan M. Petters

CISTER Research Center

Polytechnic Institute of Porto (ISEP-IPP)

Rua Dr. António Bernardino de Almeida, 431

4200-072 Porto

Portugal

Tel.: +351.22.8340509, Fax: +351.22.8321159

E-mail:

http://www.cister.isep.ipp.pt

## Abstract

Most multiprocessor scheduling theory overlooksthe possibility of hardware failures that entirely nullify thecomputation carried out by a task instance, and potentiallyalso make the respective processor henceforth unusable. Yet,such failures may occur, causing the system to fail. Motivatedby this reality, we introduce a new concept of hard real-timeschedulability guarantees for critical multiprocessor systems andanalysis for their derivation. Namely, all deadlines must be met,even in the event of a core failure. A scheduling approach,based on global fixed priorities, and accompanying analysis, forachieving such guarantees are then formulated

# Hard real-time multiprocessor scheduling resilient to core failures

Borislav Nikolić, Konstantinos Bletsas and Stefan M. Petters

CISTER/INESC-TEC, ISEP, IPP, Porto, Portugal

Email: {borni, ksbs, smp}@isep.ipp.pt

*Abstract*—Most multiprocessor scheduling theory overlooks the possibility of hardware failures that entirely nullify the computation carried out by a task instance, and potentially also make the respective processor henceforth unusable. Yet, such failures may occur, causing the system to fail. Motivated by this reality, we introduce a new concept of hard real-time schedulability guarantees for critical multiprocessor systems and analysis for their derivation. Namely, all deadlines must be met, even in the event of a core failure. A scheduling approach, based on global fixed priorities, and accompanying analysis, for achieving such guarantees are then formulated.

## I. Introduction

Hard-real-time multiprocessor scheduling requires offline guarantees of schedulability. Its critical systems subdomain though, requires even stricter schedulability guarantees, in the presence of unlikely events like execution overruns and hardware faults.

For dealing with the former kind of events (execution overruns), and in the context of mixed-criticality systems, a comprehensive toolset of analysis and design techniques is gradually being assembled [8]. However, another kind of rare event endangering a system involves a core failing, temporarily or permanently. This causes whichever task instance (job) was executing on the respective processor to be terminated short of completion as a result, and all the computation performed by it since its arrival to be wasted. Additionally, if this failure is permanent, rather than just temporary, it leaves the system with one usable core less. Even if such a failure is detected immediately, and the terminated task restarted, the schedulability of the system may be compromised because (i) the job restarted from the beginning still has to meet the absolute deadline of the *original* terminated job, which may be too close, and, if the failure is permanent, (ii) the same workload as before has to be scheduled on *fewer* cores.

*Motivational example 1*: A multiprocessor system with 100 cores has a single task with an execution time of 6 time units and a deadline of 10. It arrives at time $t = 0$ but at $t = 5$ the task is killed due to the core failure. And even if it is restarted immediately, it will miss its deadline at $t = 10$.

Adding processors, in this example, does not help meet deadlines. What *would* help, would be to speculatively launch a redundant copy of the task, in case the original one is killed. But such duplication wastes processing capacity in the general case and even so, does not always bring the desired resilience, as demonstrated with the following example:

*Motivational example 2*: Consider a multiprocessor system with three cores and two sporadic tasks $\tau_1 = (10, 10, 10)$ and $\tau_2 = (2\epsilon, 10, 10)$, where $\tau_1$ has a higher priority, and the values in brackets denote the execution time, the deadline and the period, respectively. This workload is schedulable (e.g. with global fixed-priorities) even if two copies of each task are released at each arrival. But if both tasks arrive at $t = 0$ and the core where $\tau_2$ runs fails permanently immediately afterwards, a deadline is missed (Figure 1).
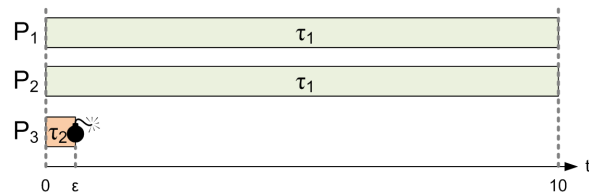


Fig. 1: Accompanying graphic for Motivational example 2.

So ideally, we seek both minimal execution redundancy *and* also offline-provable resilience in case of a core failure. To that end, we propose an approach for scheduling hard-real time tasks resiliently with respect to core failures, based on global fixed priorities. Global scheduling is a natural choice because the single run-queue tends to balance the load on the available processors, smoothing out any transient load due to the failure. It also makes it irrelevant, which is the core that fails, since cores are interchangeable and pooled together – unlike (semi- or full) partitioning. As for fixed-priority scheduling, it is a well-understood widely supported policy.

**Core failure semantics:** We disregard other kinds of faults (e.g., corruption of main memory) and only consider hardware failures specific to each core (i.e. its ALUs, registers, private caches) and failures in its core-specific software layers (OS, workload). We assume a hardware facility exists for immediate fault detection e.g., the contents of a write-back L1 cache cannot be trusted or the ALU output written to a register is suspect. How this is implemented is beyond the scope of this work, but engineering reality offers examples[1].

---

[1] For example, parity bits in L1 caches. Also, in critical systems, pairs of cores are sometimes set to operate in lockstep, with the same inputs [11], effectively being used as a single processor; mismatching core outputs then indicate a fault in one of the two.

Regardless of the failure type (i.e. a transient or a permanent one), a fault indicates that the state of whichever task was executing on the affected processor at detection time, is corrupted. Hence, the task should be aborted, but the original deadline associated with it must still be met – either by a newly launched instance of the task or by a redundant instance already launched speculatively (to cover for a potential failure, and prevent the situation described in Motivational example 1). Conversely, we assume that the state of all other tasks in the system can be trusted, either because they currently execute on non-faulty cores or (if not currently running) because their state resides entirely in main memory, which can be trusted. So all other tasks, except the one currently executing on the processor that failed, are unaffected.

Depending on the criticality of the system, as well as the nature of the fault, it may be reasonable to treat the fault as transient and reuse the core after the task in question has been killed, or alternatively, immediately take the faulty core permanently offline, which leaves the system with one core less. In this work, we consider both possibilities. In any case though, if the system is critical at all, we assume that the operator will seek to take it offline at the first convenient opportunity for doing so in an orderly manner. Therefore, although resilience guarantees under *multiple* core failures are theoretically interesting, in this work we only consider a single failure and believe that this captures most practical scenarios, even if not all. To illustrate, if the mean-time-before-failure (MTBF) for a typical core used in critical systems is hundreds of thousands of hours, the probability of a second core failing within the few minutes needed for safely shutting down the system in a controlled manner, would be too low to be of practical concern in most contexts.

In order for the above semantics and process model to be applicable to a real system, some requirements exist both for the platform and O/S and also in the application design. We already mentioned the need for a facility for detecting failures as soon as they occur. A software facility for starting (and stopping) job copies, as needed, is also required. However, the mere fact that our approach may force a task to, effectively, execute in parallel with itself, may introduce, in the absence of design safeguards against this, synchronisation hazards (i.e., when performing I/O or when accessing variables/manipulating application state) where there were previously none. These and other practical aspects and challenges are discussed in more detail in our technical report [16].

## II. RELATED WORK

As already mentioned, the focus of this work is on multiprocessors with global fixed priorities. For such a model, Lundberg [15] proposed the response-time analysis. Subsequently, Bertogna and Cirinei [6] proposed an analysis which improves the above result. The improvement is twofold. First, the authors derived a tighter upper-bound on the workload that higher-priority tasks can generate within the analysed time interval. Second, they observed that if an interfering task is "too large", not all its workload will necessarily cause interference, as some parts of it might be executed in parallel with the analysed task.

Guan et al. [13] further improve the above result. Their analysis is inspired by the aforementioned method of Bertogna

and Cirinei [6], upon which they apply concepts similar to the window analysis framework proposed by Baruah [4]. Specifically, this approach allows to derive even tighter upper-bounds on the workload generated by higher-priority tasks, and consequently derive tighter response-time estimates. This is achieved by proving an upper-bound on the number of interfering tasks that can have *carry-in workload* (defined as workload from jobs released earlier than the start of the time window in consideration) in the worst-case scenario. Recently, Sun et al. [19] proposed an improvement over the analysis of Guan et al. [13], which is applicable to task sets with arbitrary deadlines.

Until now, core failures and their implications on the schedulability analysis have received very little attention from the real-time community. In the uniprocessor domain, one notable work addressing transient core failures is that of Pathan [18]. Specifically, in order to recover from task errors, caused by transient hardware or software faults, the author proposes the release of *backup* jobs. Understandably, permanent processor failures are not covered therein, simply because they are impossible to recover from on a uniprocessor. However, on multiprocessors, this is a possibility. Ghosh et al. [12] proposed a method which allows for multiple core failures, but it is applicable only to task-sets with low-utilised tasks, whereas in our work such a restriction does not exist. Moreover, the works of Pathan [17] and Cirinei et al. [9] consider only transient faults in the multiprocessor setup.

Note that the schedulability analysis with core failures can also be studied from the perspective of *mode changes* (e.g. [14]), whereas the functioning of the system, before and after a core failure, can be perceived as distinctive system modes. Also note that any potentially different workload requirements, before and after a failure, can be analysed with the *mixed-criticality* computational model (e.g. [20]). The state-of-the-art methods from the aforementioned areas indeed efficiently take into account potential workload variations, associated to different system states. However, none of the existing works allows the possibility of variations in available hardware resources, which is an essential requirement for the integration of core failure events in the schedulability analysis. Partly, an exception is the work of Baruah and Guo [5], in the context of mixed-criticality scheduling, which allows degradation in the processor speed. Still, that problem differs from ours in the sense that Baruah and Guo [5] focus on uniprocessors, while in this work we consider multiprocessors. Additionally, in the aforementioned study, the workload is comprised of a finite sequence of independent jobs, whereas we target sporadically recurring tasks.

## III. OVERVIEW

### A. System model

We assume a multiprocessor system with $m$ identical cores and a task-set $\tau$ with $n$ sporadic tasks $\{\tau_1, \tau_2, \ldots, \tau_{n-1}, \tau_n\}$. Each task $\tau_i$ has a worst-case execution time (WCET) $C_i$, a relative deadline $D_i$ and a minimum inter-arrival time $T_i$. We assume constrained deadlines ($D_i \leq T_i$). Tasks are independent, share no resources and are to be scheduled according to a fully preemptive global fixed-priority scheme. Each task has an associated unique priority, assigned at design time, i.e. not given *a priori* as a part of the problem instance.

We use the terms "processor" and "core" interchangeably. We refer to the currently executing job affected by a core failure, and aborted with no salvageable state, as the *wasted* job. A job launched after the detection of a failure, as a replacement instance of the wasted job is termed its *copy job*, with the wasted job being referred to as its corresponding *main job*. For a task that requires redundancy, in order to meet its deadline, even in the case of a core failure, the respective copy jobs may need to be launched speculatively, just in case the corresponding main job is wasted too close to its deadline; then the two may concurrently co-exist in the system. We refer to that co-existence as "overlap" and to tasks for which overlap is possible at run-time as "overlapping". Once a main job by an overlapping task completes, its corresponding copy job is immediately terminated, as it no longer serves any purpose.

Whether a task is overlapping or not, is known offline. It is not part of the input to the design process though; it is simply an arrangement reached at the design stage itself, in order to ensure schedulability.

### B. Problem refinement

As mentioned, the system is restricted to global fixed task priority scheduling, with the actual priority assignment originally unspecified. The objective is for all tasks to meet their deadlines, assuming at most one core failure event (whether temporary or permanent). From an abstracted scheduling perspective, the only difference between these two failure types is the number of processors $m'$ that are available after a failure, which is $m$ after a temporary failure, and $m - 1$ after a permanent one. Resilience to a core failure is achieved via the use of copy jobs, as described earlier. However, the fact that we only require the system to survive a single core failure event means that, from the perspective of providing such a guarantee, it is no longer beneficial to use copy jobs after "surviving" that single failure. Thenceforth, only main jobs are released, as in conventional scheduling, and no more copy jobs are ever released. For the same reason, all active copy jobs of unaffected tasks, are also dropped, when the core fails.

Under these semantics, the problem amounts to specifying (i) a task priority assignment and (ii) the conditions that trigger the release of copy jobs at run-time, such that the system is provably schedulable even in the event of a core failure.

These decisions are intertwined, since, for example, whether some task needs to be overlapping or not may depend on the set of its higher-priority tasks. To simplify both the decision problem and the dispatching at run-time, we adopt the following arrangement:
– For a non-overlapping task $\tau_i$, a copy job is released at the time of failure of its corresponding main job.
– If $\tau_i$ is overlapping, then its copy jobs are released at a fixed (designer-set) offset $O_i$ relative to the arrival of their corresponding main jobs – unless the main job already completes before that time. For tie-breaking, a copy job always has a lower priority than its main job.

Effectively then, in the context of a particular priority assignment, we need to:
(i) Identify the tasks that require speculative copy job execution (i.e. the overlapping tasks) and
(ii) pick appropriate copy job offsets ($O_i$) for them.

As will be demonstrated with experimental evaluations (Section V), offsets entail a tradeoff: A small $O_i$ gives the copy job a greater deadline ($D_i - O_i$) within which to complete, hence it facilitates the schedulability of $\tau_i$. However, it increases the interference from speculative redundant computation onto lower-priority tasks.

### IV. SCHEDULABILITY ANALYSIS, OFFSET ASSIGNMENT AND PRIORITY ASSIGNMEMT

This section describes our proposed approach for scheduling with core-failure-resilience guarantees. First, we introduce an appropriate schedulability analysis, assuming a given priority assignment and a given offset assignment. This analysis is then used to pick the tasks that execute speculatively and their copy job offsets. It can also guide the priority assignment.

Note that the following analysis is general enough to cover either permanent or transitional core fault semantics. The only distinguishing aspect for the two scenarios, from the perspective of the analysis, is the number of available cores post-failure.

### A. System mode coverage in analysis

The system operation can be viewed in terms of two modes, S and D. The system starts in Mode S (standard) and immediately switches to Mode D (degraded) when a core failure occurs. Recall that if the failure affected some main job, upon the transition to Mode D, all copy jobs, except the one "backing up" the affected main job, are dropped. Afterwards, no further copy jobs are released. Conversely, if the failure affected some copy job, all copy jobs are dropped.

The schedulability of each task $\tau_i$ must be guaranteed in both modes, also covering the possibility that a mode switch occurs *during* the busy period of $\tau_i$. Therefore, $\tau_i$ must be provably schedulable under all the following scenarios, depending on what mode the system is in at the time of the task's absolute deadline:

- **Case 1:** Before any core fails (i.e. with the system in Mode S throughout the entire activation of $\tau_i$).

- **Case 2:** In Mode D, triggered by a core failure affecting some higher priority task $\tau_k$.

- **Case 3:** In Mode D, triggered by a core failure affecting $\tau_i$ itself.

- **Case 4:** In Mode D, triggered by a core failure affecting some lower priority task.

Case 1 can be covered by adapting the state-of-the-art analysis for global fixed-priority scheduling [13] and our model to each other, as we will describe. In Cases 2-4, we need to additionally account for the mode change (potentially also involving a reduction in the number of cores, if the core failure is permanent). But we note that Case 4 is dominated by Case 2, since lower-priority tasks cannot interfere with $\tau_i$; and the reduction in the number of cores (in the case of permanent failure) is already captured by Case 2. In other words, all other things being equal, a failure of a higher-priority task $\tau_k$ (Case 2) causes both an immediate reduction in the number of cores *and* an increase in the workload by $\tau_k$, which interferes with

$\tau_i$. By comparison, when a lower-priority task fails (Case 4), this only reduces the number of cores but it does not increase the higher-priority workload interfering with $\tau_i$.

Hence, we only consider Cases 1-3 and use the superscripts $\varnothing$, $\not{k}$ and $\not{l}$ to differentiate among them in our notation. But first we briefly present the state-of-the-art analysis [13] for standard global fixed-priority scheduling, that we build upon.

*B. The state-of-the art analysis by Guan et al. [13] for standard global fixed-priority scheduling*

Under that analysis [13], $\tau_i$ is deemed schedulable if $R_i \leq D_i$, where $R_i$ is an upper bound on its worst-case response time (WCRT), computed as

$$R_i = \begin{cases} C_i, & \text{if } |hp(i)| < m \\ C_i + I_i(R_i), & \text{otherwise} \end{cases} \quad (1)$$

where $I_i(R_i)$ is an upper bound on the interference that $\tau_i$ suffers from higher priority tasks during any interval of length $R_i$. This is computed as follows:

If a higher-priority task $\tau_j \in hp(i)$ has no carry-in workload within the busy-period of $\tau_i$, its workload over an interval of length $t$ is bounded by the expression

$$W_j^{NC}(t) = \left\lfloor \frac{t}{T_j} \right\rfloor \cdot C_j + [\![t \bmod T_j]\!]^{C_j} \quad (2)$$

where the term $t \bmod T_j \overset{\text{def}}{=} t - \left\lfloor \frac{t}{T_j} \right\rfloor$ (intuitively, the remainder of the division of $t$ by $T_j$) corresponds to the length of the carry-out interval and the operator $[\![\cdot]\!]$ is defined as

$$[\![x]\!]_{\min}^{\max} \overset{\text{def}}{=} \begin{cases} \min & \text{if } x < \min \\ x & \text{if } \min \leq x \leq \max \\ \max & \text{if } x > \max \end{cases}$$

with the arguments $\min$ and $\max$ being optional and defaulting to $-\infty$ and $+\infty$ respectively, if omitted (Figure 2).
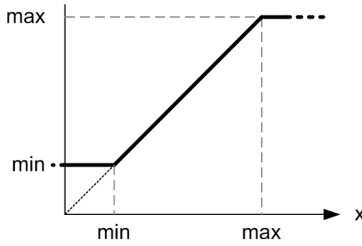


Fig. 2: Plot of $[\![x]\!]_{\min}^{\max}$ as a function of x.

However, in the general case, $\tau_j$ may also have carry-in workload. Then, its workload has the following upper-bound:

$$W_j^{CI}(t) = \overbrace{\left\lfloor \frac{[\![t - C_j]\!]_0}{T_j} \right\rfloor \cdot C_j}^{\text{body}} + \overbrace{C_j}^{\text{carry-out}}$$
$$+ \underbrace{[\![[\![t - C_j]\!]_0 \bmod T_j - (T_j - R_j)]\!]_0^{C_j-1}}_{\text{carry-in}} \quad (3)$$

In our technical report [16], we provide some intuition on how Guan et al. [13] derived these equations.

It is always the case that $W_j^{CI}(t) \geq W_j^{NC}(t)$. However, as Guan et al. [13] showed, improving on the work of Bertogna and Cirinei [6], at most $m - 1$ tasks can have carry-in. Accordingly, an upper bound on the interference suffered by $\tau_i$ within the time interval $R_i$ can then be computed as

$$I_i(R_i) = \left\lfloor \frac{1}{m} \Big( \sum_{j \in hp^{CI}(i)} [\![W_j^{CI}(R_i)]\!]^{R_i - C_i + 1} \right.$$
$$\left. + \sum_{j \in hp^{NC}(i)} [\![W_j^{NC}(R_i)]\!]^{R_i - C_i + 1} \Big) \right\rfloor \quad (4)$$

where $hp^{CI}(i)$ is the subset, of cardinality $m - 1$, of $hp(i)$ for which

$$\sum_{j \in hp^{CI}(i)} \left( [\![W_j^{CI}(R_i)]\!]^{R_i - C_i + 1} - [\![W_j^{NC}(R_i)]\!]^{R_i - C_i + 1} \right)$$

is maximised; and $hp^{NC}(i) \overset{\text{def}}{=} hp(i) \setminus hp^{CI}(i)$.

Due to $R_i$ appearing in both sides, Equation 1 is solved via a recurrence relation, as in uniprocessor WCRT analysis [3]. Note however that $hp^{CI}(i)$ and $hp^{NC}(i)$ need to be computed anew at every iteration.

*C. Analysis for the failure-resilient model*

*1) Schedulability analysis for Mode S (Case 1):* Our model differs from that of Guan et al. [13] mainly because of the copy jobs which, in the case of overlapping tasks, may concurrently exist with the corresponding main jobs. However, via a few key observations and transformations, we will adapt our model and the equations of Guan et al. [13] to each other.

Let $R_i^\varnothing$ denote an upper bound on the WCRT of the main job of $\tau_i$ under our model, under Case 1 (i.e. when the system is in Mode S).

*Lemma 1:* Under core-failure-resilient global fixed-priority scheduling with copy jobs, any copy job of task $\tau_j$, whose corresponding main job is not directly affected by a core failure, executes for at most:
– zero time units, if $\tau_j$ is non-overlapping;
– no more than $\min(C_j, R_j^\varnothing - O_j)$, if $\tau_j$ is overlapping.

*Proof:* If $\tau_j$ is non-overlapping, it never even releases a copy job *unless* its corresponding main job has failed – which would contradict the initial assumption.

If $\tau_j$ is overlapping, no copy job by it can execute for more than $C_j$ time units, by definition. But it cannot execute for more than $R_j^\varnothing - O_j$ time units either, because, according to the initial assumption, at most $R_j^\varnothing - O_j$ time units after the release of the copy job, its corresponding main job will have completed. And the copy job is terminated early at the same instant that its corresponding main job completes. ∎

*Remark 1:* The copy jobs by an overlapping task $\tau_j$ have an inter-arrival time of $T_j$.

*Proof:* Follows directly from the fact that copy jobs are released at a fixed offset $O_j$, relative to the respective main jobs (which have an inter-arrival time of $T_j$), or not at all. ∎

Lemma 1 and Remark 1 allow us to conveniently model, in Mode S, each overlapping task $\tau_j$ as two distinct tasks (main and copy), conforming to the semantics of classical scheduling (a single job per activation) and released at an offset $O_j$, with the main task $\tau_j$ having a WCET of $C_i$ and the copy task $\tau'_j$ a WCET of $C'_j \stackrel{\text{def}}{=} \min(C_j, R^\varnothing_j - O_j)$ time units[2]. Figure 3 illustrates the relationship between copy offset $O_j$ and $C'_j$.
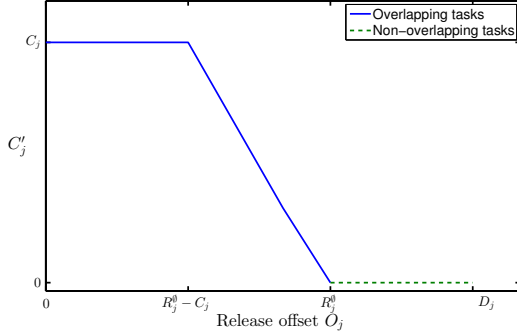


Fig. 3: The relationship between $O_j$ and $C'_j$.

By treating the main and the copy jobs as two independent tasks, we can then safely apply the existing analysis of Guan et al. [13] for global fixed-priority scheduled systems, at the cost of a bit of pessimism. For symmetry let us also model every non-overlapping task as a main/copy task pair, $\tau_j$ and $\tau'_j$, with the copy task having $C'_j = 0$.

Then, a task $\tau_i$ is deemed schedulable under Case 1 if $R^\varnothing_i \leq D_i$ where

$$R^\varnothing_i = \begin{cases} C_i, & \text{if } |hp(i)| + |hp_{ov}(i)| < m \\ C_i + I^\varnothing_i(R^\varnothing_i), & \text{otherwise} \end{cases} \quad (5)$$

where $hp_{ov}(i) \subseteq hp(i)$ is the subset of higher-priority tasks that are overlapping.

To simplify the notation when computing $I^\varnothing_i$ let $\tau_{n+j}$ denote $\tau'_j$ and $hp^*(i) \stackrel{\text{def}}{=} \cup_{j \in hp(i)} \{\tau_j, \tau'_j\}$. (Since the following equations do not differentiate between main/copy tasks, this allows us to refer to the members of $hp^*(i)$ using a single non-ambiguous index.)

To compute $I^\varnothing_i$, we calculate $W^{\varnothing|NC}_j(t)$ and $W^{\varnothing|CI}_j(t)$, for each $j \in hp^*(i)$, similarly as before:

$$W^{\varnothing|NC}_j(t) = \left\lfloor \frac{t}{T_j} \right\rfloor \cdot C_j + [\![ t \bmod T_j ]\!]^{C_j} \quad (6)$$

---

[2]It is possible that, instead of a main job, it is the copy job of some overlapping task that is affected by a core failure. Yet, this case is always dominated by the case that the main job is the one wasted – simply because it is impossible for a copy job to have received more execution time than its corresponding main job, at the time of the failure. Hence, we only need to consider the main job being wasted.

$$W^{\varnothing|CI}_j(t) = \overbrace{\left\lfloor \frac{[\![ t - C_j ]\!]_0}{T_j} \right\rfloor \cdot C_j}^{\text{body}} + \overbrace{C_j}^{\text{carry-out}}$$
$$+ \underbrace{[\![ [\![ t - C_j ]\!]_0 \bmod T_j - (T_j - R^\varnothing_j) ]\!]^{C_j-1}_0}_{\text{carry-in}} \quad (7)$$

Specifically for tasks in $hp^*(i)$ that are copy tasks of the original task set, we need to define $R^\varnothing$, for use in the above equations. Since a copy job, under Case 1 (i.e. in Mode S), is terminated short of completion, upon the completion of its corresponding main job, then as response time of a copy task $\tau'_j$ (also denoted by $\tau_{n+j}$), we can consider the quantity $R^\varnothing_{n+j} \stackrel{\text{def}}{=} R^\varnothing_j - O_j$. As for non-overlapping tasks, the workload of their copies is, by definition for Case 1, equal to zero.

Then, reasoning similarly as before

$$I^\varnothing_i(R^\varnothing_i) = \left\lfloor \frac{1}{m} \Big( \sum_{j \in hp^{\varnothing|CI}(i)} [\![ W^{\varnothing|CI}_j(R^\varnothing_i) ]\!]^{R^\varnothing_i - C_i + 1} \right.$$
$$\left. + \sum_{j \in hp^{\varnothing|NC}(i)} [\![ W^{\varnothing|NC}_j(R^\varnothing_i) ]\!]^{R^\varnothing_i - C_i + 1} \Big) \right\rfloor (8)$$

where $hp^{\varnothing|CI}(i)$ is the subset, of cardinality $m - 1$, of $hp^*(i)$ for which

$$\sum_{j \in hp^{\varnothing|CI}(i)} \Big( [\![ W^{\varnothing|CI}_j(R^\varnothing_i) ]\!]^{R^\varnothing_i - C_i + 1}$$
$$- [\![ W^{\varnothing|NC}_j(R^\varnothing_i) ]\!]^{R^\varnothing_i - C_i + 1} \Big)$$

is maximised; and $hp^{\varnothing|NC}(i) \stackrel{\text{def}}{=} hp^*(i) \setminus hp^{\varnothing|CI}(i)$.

Similarly to the original analysis [13], $hp^{\varnothing|CI}(i)$ and $hp^{\varnothing|NC}(i)$ need to be computed anew at every iteration of the recurrence relation when solving the response time equation.

*2) Schedulability analysis for the case of a core failure affecting some higher priority task $\tau_k$ (Case 2):* A core failure by affecting some $\tau_k \in hp(i)$ creates additional interference on $\tau_i$ in the short term. In the worst case, when $\tau_k$ fails just before completion, its main and copy jobs combined could execute for up to $2C_k - \epsilon$ time units within its period of $T_k$ (with $\epsilon$ arbitrarily small). Additionally, in the case of a permanent core failure, the system would be left with one core less ($m' = m - 1$) with which to process the workload, making it even harder to meet deadlines. On the other hand, dropping any other copy jobs currently executing and no longer releasing any of them in the future, after a core failure, eases up on the workload, especially on the longer term. In practice, this means that Case 2 has to be analysed separately for each higher-priority task and also that schedulability cannot be inferred from Case 1 either.

Let us now try to bound the interference onto $\tau_i$ both from $\tau_k$ (the task affected by the failure) and also from every other higher-priority task.

**Higher-priority task activations whose deadlines fall before the core failure:** This includes all activations whose deadlines fall earlier than the time instant of a core failure. Given that these activations complete at a time when the system is still in Mode S, our earlier reasoning for Case 1 applies. Hence, the respective main and copy jobs of the same task $\tau_j$ can be modelled as originating from distinct tasks $\tau_j$ and $\tau_j'$, with the latter's execution time being

$$
C_j' = \begin{cases} 0 & \text{if } \tau_j \text{ is non-overlapping} \\ \min(C_j, R_j^{\varnothing} - O_j) & \text{if } \tau_j \text{ is overlapping} \end{cases} \tag{9}
$$

Note that the above derivation requires the value of $R_j^{\varnothing}$ (from Case 1) to have been already computed for each higher-priority task $\tau_j$. This has implications for the order in which the task WCRTs need to be computed for the different cases (1-3) under consideration. We revisit this subject in Sections IV-D and IV-E.

Again each copy task $\tau_j'$ is mapped to a task $\tau_{n+j}$, for ease of referencing in equations, and $hp^{\mathcal{K}}(i) \stackrel{\text{def}}{=} \cup_{j \in hp(i)} \{\tau_j, \tau_j'\}$.

**Higher-priority task activations whose deadlines fall after the core failure:** When considering the higher-priority task activations whose deadlines fall after the core failure, we need to distinguish between the task $\tau_k$ directly affected by the failure and every other higher-priority task $\tau_j \neq \tau_k$.

For some higher-priority task $\tau_j \neq \tau_k$, we note that all activations released after the time instant of failure $t_f$ exert workload of up to $C_j$ time units (i.e. from their main job only), since copy jobs are no longer released. However, activations released before $t_f$ have workload both from the main and the copy job, in the general case. Then, it is safe to use Equation 9 to upper-bound for the respective per-job workloads. To further simplify the analysis (at the cost of some additional pessimism) let us use Equation 9 to also bound the workloads of jobs released after $t_f$, for every higher-priority job not directly affected by the failure. This approach allows us to obtain a safe upper-bound on the workload, without the need to identify the time instant $t_f$, which leads to the worst-case scenario.

As for the task $\tau_k$ that was directly affected by the failure, the same reasoning as for $\tau_j \neq \tau_k$ applies, albeit with one crucial difference. The last interfering copy job by $\tau_k$, i.e. the one that completes filling in for the main job that was terminated early by the failure, may execute up to the entire $C_k$ – not $C_k'$, if smaller. This effect can be incorporated in the modelling of workload as follows:
– If $\tau_k'$ (the copy task of $\tau_k$) has carry-in workload, we set in the equations its carry-out workload equal to $C_k$, not $C_k'$. Note that this may have a "knock-on" effect on the length of the body and carry-in intervals.
– If it has no carry-in workload, exactly one of the jobs contributing to the workload within the interval in consideration is modelled as having execution requirement of $C_k$ and all other jobs are modelled with an execution requirement of $C_k'$. To simplify the derivation, potentially at the cost of some pessimism, we shift the execution requirement $C_k$ to the instance at the start of the "body".

Based on the above reasoning we get:

For every $\tau_j \in hp^*(i)$, $j \notin \{k, n+k\}$:

$$
W_j^{\mathcal{K}|NC}(t) = \left\lfloor \frac{t}{T_j} \right\rfloor \cdot C_j + [\![t \bmod T_j]\!]^{C_j} \tag{10}
$$

$$
W_j^{\mathcal{K}|CI}(t) = \overbrace{\left\lfloor \frac{[\![t - C_j]\!]_0}{T_j} \right\rfloor \cdot C_j}^{\text{body}} + \overbrace{C_j}^{\text{carry-out}} \\
+ \underbrace{[\![[\![t - C_j]\!]_0 \bmod T_j - (T_j - R_j^{\varnothing})]\!]_0^{C_j - 1}}_{\text{carry-in}} \tag{11}
$$

Analogously as in Case 1, for each $\tau_{n+j} \in hp^*(i)$, (i.e. the copy task of $\tau_j \in hp(i)$), we assume $R_{n+j}^{\varnothing} \stackrel{\text{def}}{=} R_j^{\varnothing} - O_j$, if $\tau_j$ is overlapping.

For the main task $\tau_k$:

$$
W_k^{\mathcal{K}|NC}(t) = \left\lfloor \frac{t}{T_k} \right\rfloor \cdot C_k + [\![t \bmod T_k]\!]^{C_k} \tag{12}
$$

$$
W_k^{\mathcal{K}|CI}(t) = \overbrace{\left\lfloor \frac{[\![t - C_k]\!]_0}{T_k} \right\rfloor \cdot C_k}^{\text{body}} + \overbrace{C_k}^{\text{carry-out}} \\
+ \underbrace{[\![[\![t - C_k]\!]_0 \bmod T_k - (T_k - R_k^{\varnothing})]\!]_0^{C_k - 1}}_{\text{carry-in}} \tag{13}
$$

For its copy, denoted as $\tau_k'$ or equivalently as $\tau_{n+k}$:

$$
W_{n+k}^{\mathcal{K}|NC}(t) = [\![t]\!]_0^{C_k} + \left\lfloor \frac{t - T_k}{T_k} \right\rfloor_0 \cdot C_k' \\
+ \left[\![ [\![t - T_k]\!]_0 \bmod T_k \right]\!]^{C_k'} \tag{14}
$$

$$
W_{n+k}^{\mathcal{K}|CI}(t) = \overbrace{\left\lfloor \frac{[\![t - \mathbf{C_k}]\!]_0}{T_k} \right\rfloor \cdot C_k'}^{\text{body}} + \overbrace{\mathbf{C_k}}^{\text{carry-out}} \\
+ \underbrace{[\![[\![t - \mathbf{C_k}]\!]_0 \bmod T_k - (T_k - R_{n+k}^{\varnothing})]\!]_0^{C_k' - 1}}_{\text{carry-in}} \tag{15}
$$

Note the use of $C_k$ and $C_k'$ in Equation 15. Specifically, for the jobs of $\tau_k'$ contributing to the carry-out part of the workload, the entire $C_k$ is assumed, which is necessary, due to the failure that affected $\tau_k$. Conversely, for the jobs constituting the body and the carry-in parts, only $C_k'$ is used.

Putting this all together, we have

$$
I_i^{\mathcal{K}}(R_i^{\mathcal{K}}) = \left\lfloor \frac{1}{m'} \Big( \sum_{j \in hp^{\mathcal{K}|CI}(i)} [\![W_j^{\mathcal{K}|CI}(R_i^{\mathcal{K}})]\!]^{R_i^{\mathcal{K}} - C_i + 1} \\
+ \sum_{j \in hp^{\mathcal{K}|NC}(i)} [\![W_j^{\mathcal{K}|NC}(R_i^{\mathcal{K}})]\!]^{R_i^{\mathcal{K}} - C_i + 1} \Big) \right\rfloor \tag{16}
$$

where $hp^{\mathcal{K}|CI}(i)$ is the subset, of cardinality $m - 1$, of $hp^*(i)$ for which

$$\sum_{j \in hp^{k|CI}(i)} \left( \quad [\![W_j^{k|CI}(R_i^k)]\!]^{R_i^k - C_i + 1} \right.$$
$$\left. - \quad [\![W_j^{k|NC}(R_i^k)]\!]^{R_i^k - C_i + 1} \right)$$

is maximised; and $hp^{k|NC}(i) \stackrel{\text{def}}{=} hp^*(i) \setminus hp^{k|CI}(i)$.

Accordingly

$$R_i^k = \begin{cases} C_i, & \text{if } |hp(i)| + |hp_{ov}(i)| < m' \\ C_i + I_i^k(R_i^k) & \text{otherwise} \end{cases} \quad (17)$$

Note the use of term $m'$ rather than $m$ in Equations 16 and 17. The symbol $m'$ refers to the number of available cores after the failure, namely

$$m' = \begin{cases} m, & \text{in case of a non-permanent core failure} \\ m - 1, & \text{in case of a permanent core failure} \end{cases}$$

We pessimistically assume that the number of usable cores is $m'$ throughout the activation of $\tau_i$ in the WCRT equations, to be on the safe side without having to identify the instant of failure $t_f$ leading to the worst-case scenario. Similarly, although the number of higher-priority tasks with carry-in can be at most $m' - 1$ in Mode D, consistent with the reasoning of Guan et al. [13], we pessimistically assume $m - 1$ such tasks throughout (i.e. in the definition of $hp^{k|CI}(i)$), in order to be on the safe side.

*3) Schedulability analysis for the case of a core failure affecting the task $\tau_i$ under analysis (Case 3):* This case concerns the schedulability of some task $\tau_i$, when its main job is terminated early, due to a core failure. Then its corresponding copy job must complete within the same absolute deadline. The condition for schedulability is:

$$R_{n+i}^l \leq \begin{cases} D_i - O_i, & \text{if } \tau_i \text{ is overlapping} \\ D_i - R_i^{\varnothing}, & \text{if } \tau_i \text{ is non-overlapping} \end{cases} \quad (18)$$

where $R_{n+i}^l$ is the worst-case response time of a job by the copy task $\tau_{n+i}$ (equivalently denoted as $\tau_i'$) assuming that its corresponding main job was directly affected by the core failure. The latest that this copy job can be released is $O_i$ time units after the release of its corresponding main job, if $\tau_i$ is overlapping, or respectively, $R_i^{\varnothing}$ time units, if non-overlapping. This reduces accordingly the effective relative deadline for the copy job, leading to the condition of Equation 18.

Bounding the workloads by higher-priority jobs is analogous to the previous cases. For every $\tau_j \in hp^*(i)$:

$$W_j^{l|NC}(t) = \left\lfloor \frac{t}{T_j} \right\rfloor \cdot C_j + [\![t \bmod T_j]\!]^{C_j} \quad (19)$$

$$W_j^{l|CI}(t) = \overbrace{\left\lceil \frac{[\![t - C_j]\!]_0}{T_j} \right\rceil \cdot C_j}^{\text{body}} + \overbrace{C_j}^{\text{carry-out}}$$
$$+ \underbrace{[\![[\![t - C_j]\!]_0 \bmod T_j - (T_j - R_j^{\varnothing})]\!]_0^{C_j - 1}}_{\text{carry-in}} \quad (20)$$

However, we also need to consider the potential interference from the main job of $\tau_i$.

*Lemma 2:* The interfering workload that some copy job $\tau_i'$ can suffer by its corresponding main job, assuming that no higher-priority task $\tau_j$ has been affected by a core failure, is at most $C_i'$.

*Proof:* If $\tau_i$ is non-overlapping, its main and copy jobs can never concurrently exist. Otherwise, if $\tau_i$ is overlapping, the main can execute for at most $\min(C_i, R_i^{\varnothing} - O_i)$ time units after the release of its copy before it completes (or is terminated due to failure). Given the definition of $C_i'$ (Equation 9), the claim then holds in both cases. ∎

Lemma 2 forms an upper-bound on interfering workload upon $\tau_i'$ from any job by $\tau_i$ including the one affected by the core failure. Additionally, we know that at most one job by $\tau_i$ (i.e. the one terminated early due to the failure) interferes with the copy job $\tau_i'$ under analysis. Therefore the interfering workload from $\tau_i$ is $C_i'$ and

$$I_{n+i}^l(R_{n+i}^l) = \left\lfloor \frac{1}{m'} \left( \sum_{j \in hp^{l|CI}(i)} [\![W_j^{l|CI}(R_{n+i}^l)]\!]^{R_{n+i}^l - C_i + 1} \right. \right.$$
$$\left. \left. + \sum_{j \in hp^{l|NC}(i)} [\![W_j^{l|NC}(R_{n+i}^l)]\!]^{R_{n+i}^l - C_i + 1 + \mathbf{C_i'}} \right) \right\rfloor \quad (21)$$

where $hp^{l|CI}(i)$ and $hp^{l|NC}(i)$ are computed analogously as before. Note that the copy job $(\tau_i')$, in this case, needs to execute for the entire $C_i$ – not just $C_i'$, because it is not terminated early. Thus:

$$R_{n+i}^l = \begin{cases} C_i, & \text{if } |hp(i)| + |hp_{ov}(i)| + \text{is\_ov}(i) < m' \\ C_i + I_{n+i}^l(R_{n+i}^l) & \text{otherwise} \end{cases} \quad (22)$$

where is\_ov$(i) = 1$ if and only if $\tau_i$ is overlapping; otherwise it is zero.

*D. Offset selection for copy tasks*

To test the schedulability of some $\tau_i$, we need as inputs the $R_j^{\varnothing}$ of each $\tau_j \in hp(i)$ and knowledge of whether that $\tau_j$ is overlapping (and if so, its copy job offset $O_j$). This means that tasks have to be analysed for schedulability in a top-down priority order, assuming a given priority assignment (which itself is being tested for feasibility). By inspecting the WCRT equations, we also note that, $R_i^{\varnothing}$ of each task $\tau_i$ serves as input for the calculation of $R_i^k$s, $\forall k \in hp(i)$ and also $R_{n+i}^l$; so it has to be computed before those.

```
1.    int test_task_schedulability(int i, int p)
2.    {if !(τ_i is not schedulable at priority p when no core fails) return FAILURE;
3.     if !(τ_i is not schedulable at priority p when a core fails and τ_k is killed, ∀τ_k ∈ hp(i)) return FAILURE;
4.     if !(τ_i is not schedulable at priority p when a core fails and τ_i itself is killed) return FAILURE;
5.     // The test of the above line also computes O_i, in the process, if successful.

6.     return SUCCESS;
7.    }
```

Fig. 4: Pseudocode for testing the core-failure-resilient schedulability of a task $\tau_i$ at a priority level $p$. This assumes that the set $hp(i)$ of higher-priority tasks, and their priorities and copy jobs offets are specified, and also and that every · $\tau_j \in hp(i)$ is core-failure-resiliently schedulable, according to the same test.

```
1.    int calculate_O(τ_i) //in tandem with R^j_{n+i}
2.    {O_i:=R^∅_i; //initialisation
3.     calculate R^j_{n+i};
4.     while (O_i+R^j_{n+i} >D_i)
5.     {O_i:=D_i-R^j_{n+i}; //adjusting offset
6.      if (O_i<0) return FAILURE;
7.      calculate R^j_{n+i}; //using new O_i
8.     }
9.     return SUCCESS;
10.   }
```

Fig. 5: Calculation of $O_i$ and $R^j_{n+i}$ in tandem.

It is only at the stage of testing the schedulability of $\tau'_i$ (Case 3) that we need to consider whether or not $\tau_i$ is overlapping – and if so, what its copy offset $O_i$ is. For the purposes of schedulability testing, a non-overlapping $\tau_i$ can be equivalently modelled as having a fixed offset of $O_i = R^∅_i$. Given that the selection of $O_i$ entails a tradeoff, as earlier discussed, its value is thus best decided in tandem with the schedulability testing at Case 3 for the task in consideration. A task should be non-overlapping, if possible, to avoid interference from redundant execution onto lower priority tasks. But if a non-overlapping arrangement is not schedulable, then it is desirable to set $O_i$ to the highest value that makes $\tau'$ schedulable, in order to minimise the overlap-related interference onto lower-priority tasks. We optimally identify this value using the algorithm of Figure 5, as proven by Theorem 1.

*Theorem 1:* The algorithm of Figure 5 optimally selects the copy offset $O_i$.

*Proof:* If the task is schedulable without overlap, then the loop is never entered and SUCCESS is declared. But if it is unschedulable without overlap, then let $O^v$ denote the value that $O_i$ is updated to within the $v^{th}$ loop iteration and $R^v$ denote the $R^j_{n+i}$ calculated using $O^v$ during the same iteration. $O^{v-1}$ and $R^{v-1}$ then refer to the respective previous values.

Then, in line 5, $O^{v+1} = O^v - ((O^v + R^v) - D_i)$, i.e., the offset is decreased by the amount of time that the deadline is exceeded. And since, by inspection $R^j_{n+i}$ cannot decrease when $O_i$ decreases, this means that all offsets in the range $[O^{v+1}, O^v)$ would have been infeasible. Hence, each iteration disqualifies an infeasible subrange of $[0, R^∅_i]$, from right to left. And if a feasible offset is identified, then this will be the greatest such offset. Moreover, because task parameters are integers, after a finite number of iterations, either a feasible offset is found or the entire range for the offset $[0, R^∅_i]$ is found infeasible (line 6). ∎

Note that, although all candidate offsets greater than the one derived (in case of success) by the algorithm of Figure 5 are provably infeasible (according to Theorem 1), not all smaller offsets will be feasible in the general case, because for some offset ranges, $R^j_{n+i}$ may increase more than the amount by which the offset is decreased. Therefore, the optimal offset could not have been identified, e.g., by simpler approaches such as iterative binary (dichotomic) search. Applying binary search would either result in false negatives regarding the schedulability of the copy task or in suboptimal offsets (i.e., unnecessarily big overlap).

*E. Priority assignment*

For sporadic global fixed-priority scheduling, no exact schedulability tests are known. Therefore, it makes sense to speak of an optimal priority assignment scheme only in the context of a given, sufficient schedulability test [10]. Then a priority assignment scheme is optimal in the context of a test $T$ if it always finds a priority assignment for which the task is proven schedulable using the test $T$, if such an assignment exists at all. An obvious but often intractable $(O(n!))$ optimal priority assignment scheme, for any test $T$ is exhaustive enumeration.

Davis and Burns [10] demonstrated that Audsley's OPA [2], a bottom-up scheme with pseudo-polynomial $O(n^2)$ complexity, is optimal in the above sense for global scheduling if the schedulability test fulfils some conditions. One of those requirements is that the schedulability of a task, using the test in consideration, must depend only on the set of higher-priority tasks – but not their relative priorities. Our schedulability analysis does not meet this condition because, as discussed earlier, to test the schedulability of $\tau_i$ at a given priority (Figure 4) it requires $R^∅_j, \forall j \in hp^*(i)$ – i.e. it depends on the relative priorities of higher-priority tasks.[3] Hence the most practical *optimal* priority assignment scheme for our test is top-down, branch-and-bound enumeration, which is still intractable in the general case $(O(n!))$.

Fortunately, our experiments showed that suboptimal tractable heuristics exist that in practice perform very close to optimal. Namely, the DkC scheme [1], which assigns priorities (high to low) in order of increasing $D_i - k \cdot C_i$, where $k$ is a tunable constant. In particular, a value of $k = 1.1$, which roughly corresponds to a Slack Monotonic priority assignment, seemed to work best and significantly outperformed Deadline Monotonic (DM), in the context of our scheduling approach. Trying a few different DkC priority assignments, corresponding to different values of $k$, until success, barely underperformed the optimal branch-and-bound exhaustive enumeration. More detailed explanations are given in Section V and our technical report [16].

---

[3]Modifying the test to use $D_j$ instead of $R^∅_j$ as input would make OPA optimal in its context but would be counter-productive due to the resulting pessimism in the WCRT derivations.

## V. EVALUATION

We evaluate our approach by testing the schedulability of thousands of task sets, generated using UUnifast-Discard [10] (extension of [7])[4]. We plot the scheduling success ratios as a function of the system utilisation ($U_s \overset{\text{def}}{=} \frac{1}{m} \sum_{i=1}^{n} \frac{C_i}{T_i}$) for:

**OursH/OursS:** Our approach, assuming permanent/transient core failures, respectively.

**Guan:** The state-of-the-art test [13] for global fixed priorities without provision for fault resilience, purely for reference.

**Dupl-Part-FP:** A fault-resilient partitioned arrangement, with full task duplication, that assigns the two copies of the same task to different cores, via Best-Fit bin-packing[5], and uses Deadline Monotonic (DM) priority assignment (optimal, on each core).

**Dupl-Part-EDF:** The same, but using partitioned EDF.

Different configurations of the above approaches:

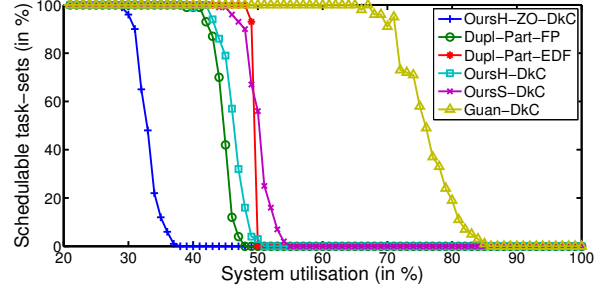**ZO:** Forcing full overlap for all tasks, i.e., $O_i = 0$, $\forall \tau_i'$.

**DkC:** Using DkC priority assignment, and trying all $k \in \{0, 0.1, \ldots, 1.9, 2\}$ until success.

Figure 6a contains the plots for systems with $m = 8$ cores and sets of $n = 16$ tasks. The schedulability improvement from optimally picking copy task offsets is considerable. OursH-DkC outperforms the duplicated partitioned fixed-priority arrangement and in case of transient core failures (OursS-DkC) the lead increases, as expected. In the experiment of Figure 6b ($n=40$), we note that the Ours-* curves improve, whereas for the other curves for fixed priorities, this is not the case or even a slight performance deteriotation is noted. The trend is amplified for $n=80$ (Figure 6c), which shows that our approach tends to be more efficient, for higher $n/m$. A higher $n/m$ ratio means smaller average task utilisations, which one might expect would improve schedulability for all approaches, since it would mean smaller bin-packing-related fragmentation on average. However, no such effect occurs for the approaches with full duplication, because the average processor utilisation is twice the nominal system utilisation, so there is little room for improvement anyway. Additionally, the subset of tasks that never suffer any interference (i.e., because they are always guarranteed a processor due to their high priority) becomes a smaller fraction of the overall task set, for higher $n/m$; this increases the potential pessimism in the schedulability test, which explains our observations. In the case of OursH-DkC and OursS-DkC however, the performance improves for higher $n/m$, because the smaller per-task average utilisation allows more tasks to be accommodated without overlap. The scheduling performace gains from the low overlap can be assessed by comparing with OursH-ZO (the naive approach with full overlap), whose performace does not improve with higher $n/m$. The results for $m = 16$ cores were similar, so we do not include them.
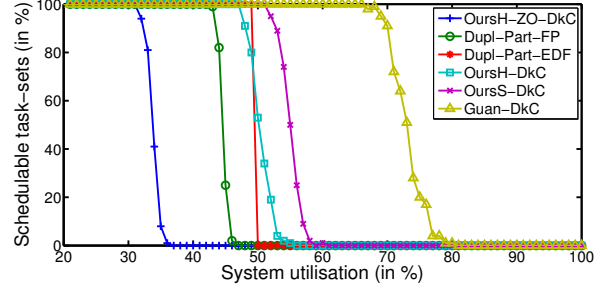
Dupl-Part-EDF outperforms our approach for task sets with lower $n/m$, which we attribute to the better scheduling potential associated with EDF, but this reverses for higher $n/m$. As noted earlier, this reversal occurs because the amount of redundant execution under our approach tends to decrease

[4]Each percentage point represents 200 implicit-deadline (D=T) task sets. Task periods were chosen uniformly over [30000, 100000] μsec.
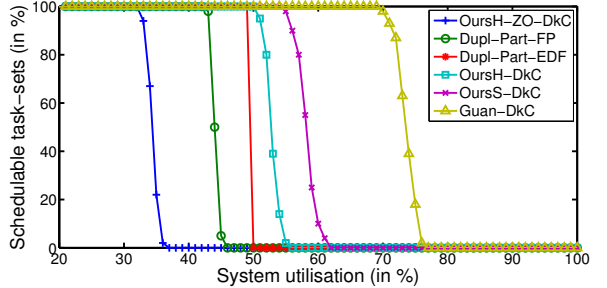
[5]In our experiments it outperforms First-Fit.

(a) $m = 8$, $n = 16$



(b) $m = 8$, $n = 40$



(c) $m = 8$, $n = 80$

Fig. 6: Experimental evaluation (part 1 of 2)

when the average task utilisation is smaller. Also, there is no room for better performance by Dupl-Part-EDF anyway, because neither this nor any other approach with full task duplication is capable of scheduling *any* task set with utilisation above half the system capacity, hence the sharp cutoff at 50%.

Next, we inspect the scheduling arrangements resulting from the application of our approach. We consider $m = 8$ cores, implicit-deadline tasks and permanent core failures. Figure 7a plots the average "effective" utilisation ($U^* \overset{\text{def}}{=} (1/m) \sum_{i=1}^{n} (C_i + C_i')/T_i$) in Mode S for *schedulable* tasks sets, according to the nominal system utilisation $U_s$. Higher $U^*/U_s$ indicates more overlap. We note that hardly any overlap is needed to schedule low-utilisation task sets. This means that there is capacity for additional background soft-real time tasks without resilience guarantees. But even for high utilisations, the additional utilisation taken up by copies ($U^* - U_s$), does not exceed 40% of $U_s$ on average.

Figure 7b, assuming $n = 40$ tasks plots the average degree of overlap $C_j'/C_j$ according to the task priority. Three different system utilisations were considered. We note that middle-priority tasks barely require any overlap to be schedulable.

For higher-utilised systems there may be some overlap in the higher or lower spectrum. High-priority tasks with overlap are mostly tasks that *must* have overlap (i.e., tasks with $u_j > 0.5$), but which also can tolerate very little interference. Obviously, such tasks are rarer in lower-utilisation task sets (line $U = 20\%$). In any case, DkC does a good job at assigning high priorities to those. At the other end, it is natural that overlap tends to be higher at lower priorities, because these tasks suffer more interference and their copies need a commensurate "head start" (small $O_j$) to complete on time. Moreover, overlap at higher priorities tends to penalise schedulability because the additional interference is exerted onto more lower-priority tasks, in a cascade effect. Therefore, an efficient priority assignment would intuitively "shift" the overlap to the lower priorities. By inspection, DkC does that.
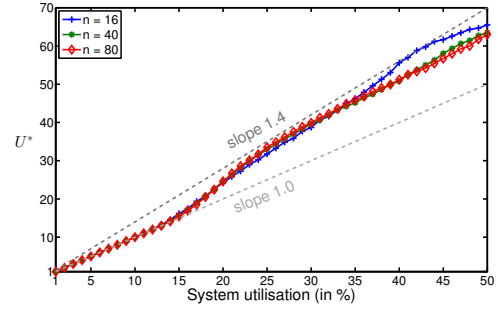
Figure 7c, examines the degree of overlap according to the task utilisation. By inspection, lower-utilisation tasks tend to require little overlap. Past 40%, the utilisation correlates positively with overlap. This behaviour matches our observations for the previous two graphs, especially with respect to the priority assignment output by DkC, and also explains why an approximate Slack Monotonic usually performs well.
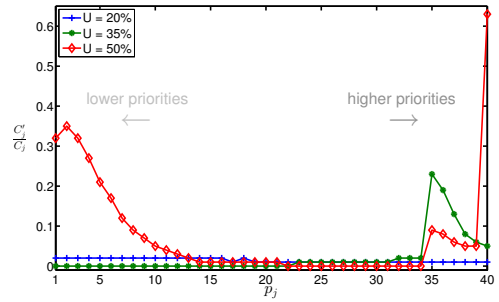
## VI. CONCLUSION

We introduced a new form of schedulability guarantees (for surviving a transient or a permanent core failure), a global fixed-priority-based scheduling arrangement for achieving them and novel analysis for their derivation. The preliminary evaluation indicates the efficiency of our approach. Many directions for future work exist, such as: (i) to consider a global EDF policy, (ii) to augment the model with arbitrary deadlines and shared resources, (including surviving a failure *during* the resource access), (iii) to reduce the analysis pessimism, and (iv) to extend the approach, so as to make it also applicable to scenarios with deferred (delayed) fault detection.
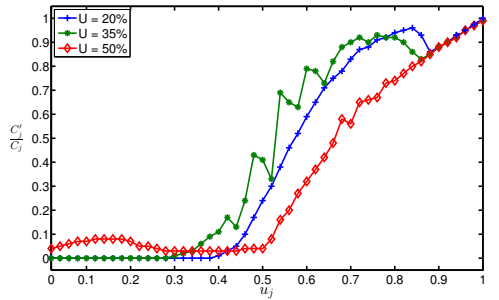
## REFERENCES

[1] B. Andersson and J. Jonsson. Fixed-priority preemptive multiprocessor scheduling: to partition or not to partition. In *Proc. RTCSA*, 2000.

[2] N. Audsley. On priority assignment in fixed priority scheduling. *Inf. Proc. Letters*, 79(1):39–44, 2001.

[3] N. Audsley, A. Burns, M. Richardson, K. Tindell, and A. J. Wellings. Applying new scheduling theory to static priority pre-emptive scheduling. *Software Engineering Journal*, 8(5):284–292, 1993.

[4] S. Baruah. Techniques for multiprocessor global schedulability analysis. In *Proc. RTSS*, 2007.

[5] S. Baruah and Z. Guo. Mixed-criticality scheduling upon varying-speed processors. In *Proc. RTSS*, 2013.

[6] M. Bertogna and M. Cirinei. Response-time analysis for globally scheduled symmetric multiprocessor platforms. In *Proc. RTSS*, 2007.

[7] E. Bini and G. Buttazzo. Measuring the performance of schedulability tests. *Real-Time Systems*, 30(1-2):129–154, 2009.

[8] A. Burns and R. Davis. Mixed criticality systems: A review. TR. Computer Science, U. of York, UK, 2013.

[9] M. Cirinei, E. Bini, G. Lipari, and A. Ferrari. A flexible scheme for scheduling fault-tolerant real-time tasks on multiprocessors. In *Proc. IPDPS*, 2007.

[10] R. I. Davis and A. Burns. Priority assignment for global fixed priority pre-emptive scheduling in multiprocessor real-time systems. In *Proc. RTSS*, 2009.

[11] Federal Aviation Authority. CAST-32: Multi-core processors. https://www.faa.gov/, 2014.

(a) Effective utilisation $U^* \overset{\text{def}}{=} \sum\limits_{i=1}^{n} \frac{C_i + C_i'}{T_i}$ vs nominal utilisation.



(b) Degree of overlap $\frac{C_j'}{C_j}$ according to the task priority.



(c) Degree of overlap $\frac{C_j'}{C_j}$ according to task utilisation $u_j \overset{\text{def}}{=} \frac{C_j}{T_j}$.

Fig. 7: Experimental evaluation (part 2 of 2)

[12] S. Ghosh, R. Melhem, and D. Mosse. Fault-tolerant scheduling on a hard real-time multiprocessor system. In *Proc. IPPS*, 1994.

[13] N. Guan, M. Stigge, W. Yi, and G. Yu. New response time bounds for fixed priority multiprocessor scheduling. In *Proc. 30th RTSS*, 2009.

[14] J. Lee and K. G. Shin. Schedulability analysis for a mode transition in real-time multi-core systems. In *Proc. RTSS*, 2013.

[15] L. Lundberg. Multiprocessor scheduling of age constraint processes. In *Proc. 5th RTCSA*, 1998.

[16] B. Nikolić, K. Bletsas, and S. M. Petters. Hard real-time multiprocessor scheduling resilient to core failures. Technical report. Available at: http://www.cister.isep.ipp.pt/people/Borislav+Nikolic/publications/.

[17] R. M. Pathan. Fault-tolerant real-time scheduling using chip multiprocessors. In *Proc. Suppl. vol. EDCC*, 2008.

[18] R. M. Pathan. Fault-tolerant and real-time scheduling for mixed-criticality systems. *Real-Time Systems*, 50(4):509–547, 2014.

[19] Y. Sun, G. Lipari, N. Guan, and W. Yi. Improving the response time analysis of global fixed-priority multiprocessor scheduling. In *Proc. RTCSA*, 2014.

[20] S. Vestal. Preemptive scheduling of multi-criticality systems with varying degrees of execution time assurance. In *Proc. RTSS*, 2007.