# Task Placement and Selection of Data Consistency Mechanisms for Real-Time Applications on Multicores

Zaid Al-bayati[†], Youcheng Sun[‡], Haibo Zeng[¶], Marco Di Natale[‡], Qi Zhu[§], and Brett Meyer[†]

McGill University[†], Scuola Superiore Sant'Anna[‡], Virginia Tech[¶], Univ. of California, Riverside[§]

E-Mails: {zaid.al-bayati@mail.,brett.meyer@}mcgill.ca, {y.sun,marco}@sssup.it, hbzeng@vt.edu, qzhu@ee.ucr.edu

*Abstract*—**Multicores are today used in several automotive, controls and avionics systems supporting real-time functionality. When real-time tasks allocated on different cores cooperate through the use of shared communication resources, they need to be protected by mechanisms that guarantee access in a mutual exclusive way with bounded worst-case blocking time. Lock-based mechanisms such as MPCP and MSRP have been developed to fulfill this demand, and research papers are today tackling the problem of finding the optimal task placement in multicores while trying to meet the deadlines against blocking times. In this work, we provide a method that improves on existing task placement algorithms for systems that use MSRP to protect shared resources. Furthermore, we leverage an additional opportunity provided by wait-free methods as an alternative data consistency mechanism in those cases where the shared resource is communication or state memory. The selective use of wait-free communication can further extend the range of schedulable systems and consequently the design space, at the cost of memory.**

## I. INTRODUCTION

Multicore architectures have become commonplace in general computing and multimedia applications, and are rapidly advancing in typical embedded computing systems, including automotive and controls. Partitioning computing tasks over multiple (on-chip) cores presents several advantages with respect to power consumption, reliability, and scalability, but it often requires significant changes to the design flow to leverage the availability of parallel processing.

Among the possible scheduling options, *partitioned fixed-priority scheduling*, where tasks are statically assigned to cores and each core is scheduled by a local fixed-priority scheduler, is widely used in embedded multicore real-time systems today. Such scheduling policy is supported by the AUTOSAR standard for automotive systems [1], as well as most commercial RTOSes, including VxWorks, QNX, LynxOS, and all POSIX-compliant ones. The designer's job for a partitioned system with real-time constraints is to define a static mapping of tasks to cores such that all tasks are guaranteed to meet their deadlines. When tasks share (communication) resources, the problem requires the consideration of the possible blocking times, which are a function of the task allocation and the selection of the data consistency mechanisms. The problem of finding such a feasible solution is demonstrated to be NP-hard (even in the special case of no shared resources, the problem is an instance of bin-packing [14]).

In real-time multicore architectures, communication resources can be shared using lock-based or wait-free methods.

**Lock-based methods** include locks with suspension, as in MPCP (Multiprocessor Priority Ceiling Protocol [24]), and spin locks, as in MSRP (Multiprocessor Stack Resource Policy [17]). They both provide a bounded worst-case blocking time, with MSRP being simpler to implement and providing better performance for short critical sections (see e.g. [7]). When the shared resource is a **communication buffer** (memory used for communicating data), another possibility is to use **wait-free methods**. The writer and readers are protected against concurrent access by replicating the communication buffer and leveraging information on the time instant and order (priority and scheduling) of the buffer access [12], [20].

Of course, wait-free methods requires the replication of the shared resource, and are not applicable if the shared resource is a hardware device. Also, wait-free methods do not directly substitute for lock-free in other ways, even when the shared resource is a simple memory buffer. We will discuss how the two methods compare in the next sections.

When using MSRP, a full MILP (Mixed Integer Linear Programming) encoding of the allocation problem (in principle capable of providing an optimal solution) is provided in [27]. As with all MILP solutions, the solver is faced with an NP-hard problem and possibly scalability issues. Hence, the Greedy Slacker (GS) heuristic is presented in [27] to obtain good-quality solutions in a much shorter time. GS assigns tasks to cores sequentially according to a simple greedy policy that tries to maximize the least slack (i.e., the difference between the deadline and the worst case response time).

*State of the Art*

Allocation problems are very common in multicore systems and in most cases they are proven to be special instances of the general bin-packing [21] problem. In the context of real-time systems and not considering shared resources, several heuristics have been proposed, e.g., [6], [14], [16]. For the case of independent tasks, Baruah presented a polynomial-time approximation scheme [4], and Chattopadhyay and Baruah showed how to leverage lookup tables to enable fast, yet arbitrarily accurate partitioning [10]. Finally, Baruah and Bini presented an exact MILP based approach for partitioning [5]. However, in all these cases, tasks are assumed to be independent and shared resources are not considered. Similarly, solutions for the partitioned scheduling of independent sporadic tasks are presented in [15], [4], [13].

MPCP [24] and MSRP [17] are protocols for sharing resources with predictable blocking times in multicore platforms. The response time bound in [17] for MSRP has been recently

imrpoved using a MILP formulation, as presented in [28]. MPCP and MSRP have been compared in a number of research works (with respect to the worst-case timing guarantees) with general consensus that MSRP performs best for short (global) critical sections and MPCP for large ones. A more detailed discussion on the design options and the characteristics of lock-based resource protection algorithms can be found in [8]. Wait-free mechanisms are an alternative option for preserving data consistency in multicore communication [12]. They can be extended to guarantee semantics preservation of synchronous models [29] and sized according to the time properties of communicating tasks [25], [26].

With respect to the real-time task placement (partitioning) problem, when considering the possible access to global resources, Lakshmanan et al. [22] presented a partitioning heuristic tailored to MPCP. This heuristic organizes tasks sharing resources into groups in order to assign them to the same processor. In subsequent work, Nemati et al. presented BPA [23], another partitioning heuristic for MPCP following a similar approach. It tries to identify communication clusters in such a way that globally shared resources are minimized and the blocking time reduced.

The first work to provide solutions for the partitioning problem when using spin locks (MSRP) is [27]. Two solutions are proposed in the paper: one is the MILP optimization formulation that can provide the optimal solution, but its runtime is exponentially increasing with the problem size; the other is the GS heuristic. [18] provides implementations of MPCP, MSRP, and wait-free methods on two open source RTOSes and the selection of protection mechanisms that satisfies the schedulability constraints, but the task placement is assumed to be given.

*Our Contributions*
In this work, we first explore the possibility of improving on the Greedy Slacker algorithm. The proposed algorithm, defined as Communication Affinity and Slack with Retries (CASR), attempts to improve on the GS allocation strategy in two ways: it considers task (communication) affinity in the allocation and includes a recover-and-retry mechanism in those cases where the slack-based policy fails to find a feasible allocation.

Next, we consider an additional design option when the shared resource is a communication buffer, by selectively replacing MSRP-protected resources with those protected by a wait-free mechanism. Wait-free methods have virtually no blocking time (in reality often negligible) and may enhance the schedulability of the system at the cost of additional memory. We combine these two methods for managing global resources to find a system configuration that is schedulable with the minimum amount of required memory. We present two approaches that try to compute the optimal task placement and selection between MSRP and wait-free methods for global communication resources. The first (GS-WF) simply extends the GS by using wait-free methods when GS fails to allocate the next task. The algorithm starts by using MSRP for all global resources, and uses wait-free methods only when needed. The second approach (MPA, Memory-aware Partitioning Algorithm) is a heuristic that leverages wait-free methods from the start (while still trying to minimize the overall memory cost). Experiments show that the latter heuristic provides better results in terms of both schedulability and memory.

The paper is organized as follows. We define our system

model and summarize the basic principles and analysis results for the methods considered in this work in Sections II and III. In Section IV, we present the CASR algorithm. In Section V we describe the two approaches for selecting the data consistency mechanisms. We give an illustrative example in Section VI. In Section VII, we evaluate and compare the results of the four algorithms (GS, CASR, their extension with wait-free methods, and MPA) in terms of schedulability and memory cost. Finally, Section VIII concludes the paper.

## II. System Model

The system under consideration consists of $m$ cores $\mathcal{P} = \{p_1, p_2, \ldots, p_m\}$. $n$ tasks $\mathcal{T} = \{\tau_1, \tau_2, \ldots \tau_n\}$ are statically allocated on them and scheduled by static priority. Each task $\tau_i$ *is activated by a periodic or sporadic event stream with period or minimum interarrival* $T_i$. The execution of task $\tau_i$ is defined as a set of alternating critical sections and sections in which the task executes without using a (global or local) shared resource, defined as *normal execution segments*. The worst case execution time (WCET) is defined by a tuple $\{C_{i,1}, C'_{i,1}, C_{i,2}, C'_{i,2}, ..., C'_{i,s(i)-1}, C_{i,s(i)}\}$, where $s(i)$ is the number of normal execution segments, and $s(i) - 1$ is the number of critical sections. $C_{i,j}$ ($C'_{i,j}$) is the WCET of the $j$-th normal execution segment (critical section) of $\tau_i$. $\pi_i$ denotes the nominal priority of $\tau_i$ (*the higher the number, the lower the priority*, thus $\pi_i < \pi_j$ means $\tau_i$ has a higher priority than $\tau_j$), and $P_i$ the core where it executes. The global shared resource associated to the $j$-th critical section of $\tau_i$ is $\mathcal{S}_{i,j}$. The WCET $C_i$ of $\tau_i$ is

$$C_i = \sum_{1 \leq j \leq s(i)} C_{i,j} + \sum_{1 \leq j < s(i)-1} C'_{i,j} \qquad (1)$$

The worst case response time $R_i$ of $\tau_i$ can be calculated as the least fixed-point solution of the formula (2), where $B_i^l$ is the local blocking time, and $B_i^r$ the remote blocking time, that is, the worst-case length of the time interval in which the task has to wait because of critical sections executed by lower priority tasks on local and global resources, respectively. $C_i^*$ is the worst-case execution time for task $\tau_i$ after considering the possible additional time spent while spinning on a global lock (in the case of spin lock based protocols).

$$R_i = C_i^* + B_i^l + B_i^r + \sum_{\pi_h < \pi_i \wedge P_h = P_i} \left\lceil \frac{R_i}{T_h} \right\rceil C_h^* \qquad (2)$$

## III. MSRP and Wait-free Communication: Blocking Time vs. Memory

We propose to combine MSRP and wait-free methods to improve system schedulability while keeping the memory overhead low. We first introduce them and highlight the tradeoff between the associated blocking time and memory cost.

### A. Wait-free Communication Buffers

The objective of wait-free methods is to avoid blocking by ensuring that each time a writer needs to update the communication data, it is reserved with a new buffer. Readers are free to use other dedicated buffers. Figure 1 shows the typical stages performed by the writer and the readers in a wait-free protocol implementation.
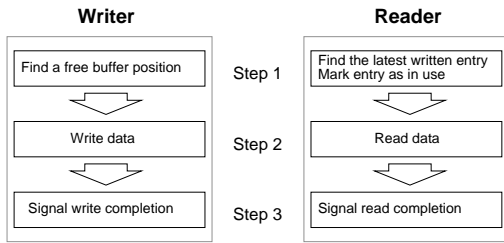
Fig. 1: Writer and readers stages in wait-free methods.

These stages have been implemented in [11]. The algorithm makes use of three global sets of data. An array of buffers is sized so that there is always an available buffer for the writer to write new data. An array keeps in the $i$-th position the buffer index in use by the $i$-th reader. Finally, a variable keeps track of the latest buffer entry that has been updated by the writer. Each reader looks for the latest entry updated by the writer, stores its index in a local variable, and then reads its contents.

Consistency in the assignment of the buffer indexes can be guaranteed by any type of hardware support for atomic operations, including Compare-And-Swap (CAS) (as in the original paper by Chen and Burns [11]), but also Test and Set, or other operations (the required support for atomicity is discussed at length in [19]) and also hardware semaphores, as available in several architectures today. The use of hardware support for atomicity is not limited to wait-free methods. Any implementation of MSRP needs similar instructions for the low-level protection of the lock data structures (including the task FIFO).

The implementation of a wait-free mechanism is possible in constant time (independent from the number $n_r$ of the readers, as opposed to the $O(n_r)$ solution in [11]) using the Temporal Concurrency Control Protocol (TCCP) discussed in [26], without the need of *additional blocking times* and limited overhead (as discussed next). However, it requires an array of buffers, thus *additional memory*, if compared with lock-based methods. In TCCP, the number of required buffers is

$$\max_{j \in \mathbb{R}} \left\lceil \frac{l_j}{T_w} \right\rceil \tag{3}$$

where $\mathbb{R}$ is the set of readers, $T_w$ is the period of the writer, and $l_j$ is the lifetime of the reader $j$ (the sum of its worst-case response time $R_j$ and the largest time distance $O_{w,j}$ between the any activation instant of the reader and the previous activation of the writer).

### B. Wait-free vs. Lock based, limitation and performance issues

Wait-free methods are not a simple replacement for lock-based methods. They can only be used for memory resources (communication or state buffers) In those cases in which the access to the shared resource cannot be simply isolated as a set of reads, clearly separated from writes into the memory buffer, such as a counter increment, wait-free methods still apply, upon condition that the task that performs the operation reads the content from a buffer item and writes its updated content onto a different buffer. A simple `x = x+1` wrapped inside a lock is a simpler solution, and requires a single lock-overhead compared with the sum of the overheads for a read and a write in the wait-free case. However, the counter example

may be misleading. The operations that affect more the worst-case performance are those that require large access times, where a read-then-write pattern is less frequent.

As for the performance, even though the wait-free implementation requires multiple copies of the buffer, it does not require additional copy operations. Simply, instead of overwriting the same memory section, the writer writes its updates on a different memory buffer every time (in cyclic fashion). Also, readers get a reference to the buffer item they are allowed to use (the freshest one), which is guaranteed not to be overwritten by the writer, and therefore can use the data *in place, without the need of additional copies.*

In [18], a measurement-based comparison of wait-free (TCCP) and lock-based (MSRP) implementations on the Erika open source OSEK OS, running on a Freescale FADO processor (a heterogeneous dual-core) has been performed. The overhead of TCCP is no larger than 43 CPU cycles for TCCP, or 0.36 microsecond if the CPU frequency is at 120MHz, including the time to execute the hardware semaphore instruction used to achieve atomicity in the implementation. In all the measures, TCCP had an overhead one third of the corresponding MSRP overhead or less (in MSRP to set up the data structures for acquiring or releasing the lock). Of course this does not related to the blocking time, which is the time spent *after the lock acquisition and before its release.*

Hence, with a limited penalty for the wait-free implementation, we assume that overheads even out or are negligible when compared to task executions and resource access times, and let $C_i^* = C_i$ for wait-free. The task response time can be computed using the general formula (2) by also setting $B_i^l = 0$ and $B_i^r = 0$, as there is no local or global blocking.

### C. Multiprocessor Stack Resource Policy

MSRP (Multiprocessor Stack Resource Policy) [17] is a multiprocessor synchronization protocol, derived by extension from the single-core Stack Resource Policy (SRP) [3]. In SRP, each resource has a ceiling equal to the highest priority of any tasks that may access it. At runtime, a task executing a critical section immediately changes its priority into the ceiling of the corresponding resource. In SRP tasks never block while executing when requesting a shared resource, but can possibly wait for lower priority tasks to terminate their critical section while being in the ready queue. For single-core, SRP allows to bound the local blocking time $B_i^l$ as the longest critical section executed by a lower priority task on a resource with a ceiling higher than the task priority.

$$B_i^l = \max_{k:\pi_k > \pi_i \wedge P_k = P_i} \left\{ \max_{1 \leq m < s(k)} C'_{k,m} \right\} \tag{4}$$

In the multicore extension MSRP, global resources are assigned a ceiling that is higher than that of any local resource. In MSRP, a task that fails to lock a global resource (in use by another task) *spins* or the resource lock until it is freed, keeping its processor busy (in the MPCP protocol, for comparison, the task is suspended and yields the CPU). To minimize the spin lock time (wasted CPU time), tasks cannot be preempted when executing a global critical section, in an attempt to free the resource as soon as possible. MSRP uses a First-Come-First-Serve queue (as opposed to a priority-based queue in MPCP) to manage the tasks waiting on a lock for a given busy resource.

In this paper, we adopt the succifient analysis in [17]. The spin time $L_{i,j}$ that a task $\tau_i$ may spend for accessing a global resource $\mathcal{S}_{i,j}$ can be bounded by [17]

$$L_{i,j} = \sum_{E \neq E_i} \left\{ \max_{\tau_k : P_k = P_i, 1 \leq m < s(k)} C'_{k,m} \right\} \tag{5}$$

This is the time increment to the $j$-th critical section of $\tau_i$. Thus, the total worst case execution time $C_i^*$ is

$$C_i^* = C_i + \sum_{1 \leq j < s(i)} L_{i,j} \tag{6}$$

MSRP maintains the same basic property of SRP, that is, once a task starts execution it cannot be blocked. The local blocking time $B_i^l$ is the same as in SRP (Equation (4)) and the worst-case remote blocking time $B_i^r$ is [17]

$$B_i^r = \max_{k : \pi_k > \pi_i \wedge P_k = P_i} \left\{ \max_{1 \leq m < s(k)} (C'_{k,m} + L_{k,m}) \right\} \tag{7}$$

Recently, a more accurate analysis of the blocking time in MSRP based on a MILP formulation has been proposed [28]. However, for its use in an optimization problem in which many solutions need to be evaluated quickly, it requires significantly more time making it unscalable to large systems as shown in Section VII.

and most likely the better accuracy does not impact the selection of the solution by the heuristics.

## IV. CASR: Task Placement with Affinity

The Greedy Slacker (GS) algorithm [27] presents a method to obtain a resource-aware assignment of tasks to cores. It is based on a simple rule: tasks are ordered by decreasing task density $(C_i/D_i)$ and assigned sequentially to cores. At each step, the candidate task is assigned to the core that guarantees the maximum value for the minimum slack of the already allocated tasks. To determine the minimum slack assuming a task is assigned to a specific core, GS uses a modified version of Audsley's optimal priority assignment scheme [2]. The main problem with GS is that slacks cannot be computed unless all tasks are allocated. Therefore, at each stage, the task slacks can only be estimated by using the subset of tasks already allocated. In addition, GS does not include recursion or retries. If at any point, a task fails to be assigned to any core, no alternative solutions are tried and the algorithm fails.

The Communication Affinity and Slack with Retries (CASR), listed in Algorithm 1, is an improvement on the Greedy Slacker (GS) [27] heuristic with the following features:
– CASR tries to assign tasks sharing the same resource on the same core by considering their communication affinity.
– When a task is unschedulable on any core, CASR has a recovery mechanism that de-allocates tasks and retries.

A task $\tau_i$ is **affine** to another task $\tau_j$ if the two share some resource. A core $p$ is **affine** to $\tau_i$ if at least one task affine to $\tau_i$ has been assigned to $p$. Task affinity is not transitive. For example, if $\tau_i$ uses $\{r_1, r_2\}$, $\tau_j$ uses $\{r_2, r_3\}$ and $\tau_k$ uses $\{r_3, r_4\}$, then $\tau_i$ and $\tau_k$ are both affine to $\tau_j$. However, $\tau_i$ and $\tau_k$ are not affine to each other.

CASR has a main loop in which tasks are considered for allocation in order of their density. For each task $\tau_i$ in the set to be allocated $\mathcal{NA}$, CASR tries to allocate $\tau_i$ on its affine cores

first. Intuitively, this is an attempt at trading global blocking for local blocking. However, without further constraints, a core could be overloaded by a large cluster of affine tasks and the partitioning procedure results in failure. To prevent such a scenario, an affine core $p$ can only be added to the set of candidate cores $\mathcal{Q}$ considered for the allocation of the task if its total utilization $U(p)$ is lower than a predefined utilization bound $(U_b)$ (lines 5-9). If there is no available affine core, e.g. they are all overloaded, CASR simply considers all cores as candidates to host $\tau_i$ (line 11).

After finding the set of candidate cores $\mathcal{Q}$, the following steps (lines 12-20) are the same as in GS. The tryAssign method is iteratively called on each of these cores. A task is assigned to the core where tryAssign returns the maximum least normalized slack ($s$). Like [8], tryAssign uses the Audsley's algorithm [2] to find a priority order of the tasks.

The second salient feature of CASR with respect to GS is the recovery procedure that is invoked when no core is available for $\tau_i$ (in the algorithm, when the set $C$ is empty). Every time the allocation of a task $\tau_i$ fails, CASR de-allocates all its affine tasks (previously allocated to some core) and puts them back in the list $\mathcal{NA}$ (lines 28-29). This recovery procedure, however, is not used unconditionally. The first time a task fails to find a suitable core, it is put in the black-list $\mathcal{BL}$ (line 26); the second time it fails, the task is put in the post-black-list $\mathcal{PBL}$ (line 23). When the post-black-list $\mathcal{PBL}$ is used first, continuing to the assign a task to its affine cores will most likely lead to another failure. Hence, when a task enters the list $\mathcal{PBL}$, the CASR algorithm stops considering the affinity between cores and tasks, by setting the aff_chk flag to false (line 24). Only if the allocation of a task inside the $\mathcal{PBL}$ fails, the CASR task partitioning procedure returns failure (line 21).

$U_b$ is a tunable parameter for the algorithm. Unfortunately, the algorithm is sensitive to the value of $U_b$ and there is no single value that performs best in all our experiments. A value that performed very well on average and provides for a good baseline is $U_b = U_\mathcal{T}/m$, where $U_\mathcal{T}$ is the total task utilization, with the obvious intuitive meaning of trying to achieve load balancing. In our experiments (Section VII) this value performs consistently better than GS. However, given that the runtime of the algorithm allows for multiple executions in a reasonable time, it is possible to execute the algorithm with a set of different values for $U_b$ (e.g., from 0 to 1 in steps of 0.25 or 0.1) and select the best outcome among them. This allows to further and sensibly improve upon the GS solution.

## V. Using Wait-free Methods to Increase Schedulability

As summarized in Section III, MSRP provides data consistency for shared buffers at the price of a blocking time and a negative impact on the schedulability of the system. The blocking time increases with the number and the length of the critical sections on global shared resources. Wait-free methods provide an alternative way to ensure consistent data access to shared resources. They have practically no blocking time, but this comes at the price of additional memory cost. To leverage the complementary characteristics of the two approaches, we propose to use a combination of MSRP and wait-free methods. We propose two algorithms to handle the selection of the partitioning and protection mechanism in multicore systems.

**Algorithm 1:** Communication Affinity and Slack with Retries

```
 1: 𝒩𝒜 ← 𝒯; aff_chk=true; ℬℒ ← {}; 𝒫ℬℒ ← {};
 2: while 𝒩𝒜 ≠ ∅ do
 3:     τ_i ← HighestDensity(𝒩𝒜); 𝒬 ← ∅
 4:     if aff_chk then
 5:         for all p ∈ 𝒫 do
 6:             if Affine(p, τ_i) and U(p) ≤ U_b then
 7:                 insert p in 𝒬
 8:             end if
 9:         end for
10:     end if
11:     if 𝒬 = ∅ then 𝒬 = 𝒫
12:     C ← ∅
13:     for all p ∈ 𝒬 do
14:         if tryAssign(τ_i, p) then C ← (s, p)
15:     end for
16:     if C ≠ ∅ then
17:         p = FindMaxS(C)
18:         Allocate(τ_i, p); 𝒩𝒜 = 𝒩𝒜 \ {τ_i}
19:         continue
20:     end if
21:     if τ_i ∈ 𝒫ℬℒ then return failure
22:     if τ_i ∈ ℬℒ then
23:         𝒫ℬℒ = 𝒫ℬℒ ∪ τ_i
24:         aff_chk = False
25:     else
26:         ℬℒ = ℬℒ ∪ τ_i
27:     end if
28:     𝒟 = AffineSet(τ_i);
29:     for all τ_k ∈ 𝒟 do DeAllocate(τ_k); 𝒩𝒜 = 𝒩𝒜 ∪ {τ_k}
30: end while
```

The first is a wait-free extension of the GS and CASR algorithms. The second is an entirely novel memory-aware partitioning heuristic.

### A. Extending Greedy Slacker with Wait-free Methods

Our first algorithm (GS-WF) extends the greedy slacker policies by providing an additional recovery option. Once a task fails to be assigned to any core, either in the original GS or after the latest attempt from the $\mathcal{PBL}$ in CASR, GS-WF uses wait-free mechanisms for all the global resources accessed by the task and tries to assign the task to each of the cores. If the use of wait-free methods makes the task schedulable on more than one core, GS-WF selects the one with the maximum value of the smallest normalized slack. If the task can not be assigned to any core even after using wait-free methods, the algorithm fails.

### B. Memory-aware Partitioning Algorithm

Extending the greedy slacker algorithm with wait-free methods enhances schedulability. However, our extension only enhances the assignments that the GS or CASR defined in case of their failure. Given that the allocation decisions are inherently based on minimizing slack and do not consider the memory cost, the end result can be schedulable but quite inefficient in terms of memory requirements. Hence, we propose the Memory-aware Partitioning Algorithm (MPA). The algorithm consists of two phases:
1) Finding an initial feasible solution (possibly of good quality, that is, low memory use);
2) Improving the solution by exploring other possible solutions using a local search.

In the first phase, the algorithm tries to obtain an initial schedulable solution. In order to maximize the probability of finding such a solution, all global resources are initially protected by wait-free methods. Local resources are managed using the SRP policy [3]. The second phase reduces the memory cost resulting from the use of wait-free methods by selectively changing the data consistency mechanism to MSRP and using local search. We first describe the concept of *assignment urgency*.

*1) Assignment Urgency:* In the proposed algorithm, tasks are allocated to cores in a sequential order. The order has a significant impact on the schedulability and the overall cost of the task allocation. We propose the concept of **Assignment Urgency** ($AU$), which is an estimate of the penalty (in schedulability or memory) if a task is not the next to be assigned. For a task $\tau_i$, $AU_i$ is defined as follows:

$$AU_i = \begin{cases} M_{max} & \tau_i \text{ schedulable on 1 core} \\ \left| \min_{1 \leq j \leq m} MC(\tau_i, p_j) - \min_{\substack{1 \leq k \leq m \\ k \neq j}} MC(\tau_i, p_k) \right| & \tau_i \text{ schedulable on } > 1 \text{ cores} \end{cases}$$
(8)

where $MC(\tau_i, p_j)$ is the memory cost of assigning $\tau_i$ to $p_j$ using wait-free for all its global resources, and $M_{max}$ is a value higher than the worst case memory cost of assigning any task to any core in the system.

$$M_{max} = \max_{\forall \tau_i \in \mathcal{T}} \max_{\forall p_j \in \mathcal{P}} MC(\tau_i, p_j) + 1$$
(9)

Defining $M_{max}$ in this way ensures that tasks schedulable on only one core will always have higher $AU$ values than tasks with more assignment options. If a task can be scheduled on more than one core, then there is usually enough allocation freedom to defer its assignment. However, this might come with a cost in memory since the task may have to be scheduled on a different core that could not guarantee the lowest memory cost. The possible memory penalty is quantified by the (absolute) difference in memory between the core where it has best memory cost and the one with the second best.

The function `TaskSort()` in Algorithm 2 computes the assignment urgencies and sorts the unassigned tasks in the system by decreasing AU values. $\mathcal{LT}$ contains the list of tasks that have not been assigned yet. The feasible cores for a given unassigned task $\tau_i$ are stored in a list of candidate solutions $c\mathcal{P}$. Each entry in $c\mathcal{P}$ is a pair of core and cost *(p, c)*, where $c$ is the memory cost of assigning $\tau_i$ to the core $p$, as computed by the function MemoryCost(). If there are no feasible solutions for a given task, the algorithm reports failure (line 10). Otherwise, the task assignment urgency is calculated according to Equation (8) (lines 12-17). The function also computes the best candidate core for a given task $\tau_i$ as $BP_i$ (line 18). It then sorts the list $\mathcal{LT}$ by decreasing assignment urgency (line 20). In case of a tie in the AU values, the task with the highest density is assigned first.

The main algorithm for task allocation and resource protection selection works in two phases, as shown in Algorithm 3. In the first phase, we try to obtain an initial assignment of tasks to cores with a low memory cost. In the second phase, the solution is improved by local search.

*2) Phase 1:* The first phase is a greedy algorithm that assumes the use of wait-free methods for all global resources and tries to assign each task to the core on which it has the least memory cost. The algorithm uses the function `TaskSort()`

**Algorithm 2:** Calculating assignment urgencies

```
 1: Function TaskSort(𝓛𝓣)
 2: for all τᵢ in 𝓛𝓣 do
 3:     c𝓟 = {}
 4:     for all pⱼ in 𝓟 do
 5:         if isSchedulable(τᵢ,pⱼ) then
 6:             c𝓟 ← (pⱼ, MemoryCost(τᵢ,pⱼ))
 7:         end if
 8:     end for
 9:     if c𝓟 = ∅ then
10:         return failure
11:     end if
12:     if size(c𝓟) = 1 then
13:         AUᵢ = Mₘₐₓ
14:     else
15:         sortByCost(c𝓟)
16:         AUᵢ = c𝓟[1].c-c𝓟[0].c
17:     end if
18:     BPᵢ = c𝓟[0].p
19: end for
20: sortByAU(𝓛𝓣)
21: return success
```

to sort the set of unassigned tasks in the list $\mathcal{LT}$ by their AUs, and compute the best core for the allocation of each task (where it causes the least memory cost). The task with the highest assignment urgency is then assigned to its best candidate core (lines 13-14). At any time, if a task has no feasible core, the algorithm resets the task assignments and tries a resource-oblivious Best-Fit Decreasing (BFD) policy (lines 5-11). If BFD also fails, the algorithm fails. At the end of the first phase, all tasks should be assigned to a core in a task allocation scheme $\mathcal{TA}$.

**Algorithm 3:** Memory-aware Placement Algorithm

```
 1: Function AllocationAndSynthesis(𝒯)
 2: Phase 1:
 3: 𝓛𝓣 ← 𝒯
 4: while 𝓛𝓣 ≠ ∅ do
 5:     if TaskSort(𝓛𝓣) = failure then
 6:         Reset partitioning
 7:         if tryBFD() = true then
 8:             Goto Phase 2
 9:         else
10:             return failure
11:         end if
12:     else
13:         τₖ = ExtractFirst(𝓛𝓣)
14:         𝒯𝒜 ← Allocate(τₖ, BPₖ)
15:     end if
16: end while
17:
18: Phase 2:
19: optimizeResources(𝒯𝒜)
20: CurOpt = 𝒯𝒜; 𝒩𝒜 ← 𝒯𝒜
21: while 𝒩𝒜 ≠ ∅ do
22:     Th = MemoryCost(GetLast(𝒩𝒜))
23:     curSys = ExtractFirst(𝒩𝒜)
24:     𝓛𝒩 = generateNeighbors(curSys)
25:     for all 𝒜𝒮 in 𝓛𝒩 do
26:         optimizeResources(𝒜𝒮)
27:         if IsSchedulable(𝒜𝒮) and MemoryCost(𝒜𝒮) < Th and not visited(𝒜𝒮)
            then
28:             𝒩𝒜 ← 𝒜𝒮
29:             if size(𝒩𝒜) > n then RemoveLast(𝒩𝒜)
30:             if cost(𝒜𝒮) < cost(curOpt) then curOpt=𝒜𝒮
31:             Th = MemoryCost(GetLast(𝒩𝒜))
32:         end if
33:     end for
34:     if NoChange(curOpt, #iter) OR cost(curOpt) ≤ tgtCost then
35:         return CurOpt
36:     end if
37: end while
```

*3) Phase 2:* In the second phase, the solution obtained in the first phase is improved with respect to its memory cost. This is done in an iterative manner by exploring selected neighbors (with small memory cost) of candidate solutions starting from the initial solution obtained at the end of phase 1. Phase 2 uses two sub-functions:
– optimizeResources(), which optimizes data consistency mechanisms for shared resources;
– generateNeighbors(), which generates neighboring solutions for a given task assignment.

optimizeResources() is the function responsible for assigning a data consistency mechanism to all globally shared resources. The function optimizeResources() changes the protection mechanism of global resources from wait-free to MSRP (for as many resources as possible) while retaining schedulability. It is called at the start of phase 2 (line 19, Algorithm 3) and whenever a new candidate solution is found (line 26, Algorithm 3), to optimize its memory cost.

The determination of the optimal resource protection mechanism for each globally shared resource would require an exhaustive analysis of all the possible configurations with complexity $2^r$, where $r$ is the number of global resources in the system, which could easily be impractical. Therefore, a heuristic procedure (detailed in Algorithm 4) is used. The procedure changes the protection mechanism for all the global resources to wait-free and places them in a list $\mathcal{GR}$ (lines 2-5). The procedure then sorts the resources in $\mathcal{GR}$ by decreasing memory cost of their wait-free implementation (line 6). Then, the protection mechanism of the first resource in $\mathcal{GR}$ (the resource with highest cost) is changed to MSRP and the system schedulability is checked. If the system becomes unschedulable, the protection mechanism is reversed back to wait-free. The procedure iterates through the resources in $\mathcal{GR}$, changing the protection mechanism from wait-free to MSRP whenever possible (lines 7-12).

**Algorithm 4:** Determining resource protection mechanisms

```
 1: Function optimizeResources(𝒯𝒜)
 2: 𝒢ℛ = FindGlobalResources(𝒯𝒜)
 3: for all rᵢ in 𝒢ℛ do
 4:     setProtocol(rᵢ, WAITFREE)
 5: end for
 6: sortByMemoryCost(𝒢ℛ)
 7: for all rᵢ in 𝒢ℛ do
 8:     setProtocol(rᵢ, MSRP)
 9:     if isSchedulable(𝒯𝒜) = false then
10:         setProtocol(rᵢ, WAITFREE)
11:     end if
12: end for
```

The other function invoked in the main loop of phase 2 is generateNeighbors(), shown in Algorithm 5. generateNeighbors() generates neighbors and places them in the list of neighbors $\mathcal{LN}$. It computes the set of neighboring solutions to be evaluated in the context of the local search for improving the current solution. A *neighbor* of a given task allocation solution can be obtained by a re-assignment of a task $\tau_i$ allocated on core $P_a$ to a different core $P_b \neq P_a$ that can accommodate it, either directly (1-move neighbor) or by removing a task $\tau_j$ with equal or higher utilization from $P_b$ and assigning $\tau_j$ to another core $P_c \neq P_b$ (2-move neighbor).

The **main loop** of phase 2 (lines 21-37 of Algorithm 3) is similar to a branch-and-bound algorithm. It performs a best-

**Algorithm 5:** Generating neighboring solutions

```
 1: Function generateNeighbors(curSys)
 2:   LN.clear()
 3:   GR = FindGlobalResources(curSys)
 4:   for all r_i in GR do
 5:     setProtocol(r_i, WAITFREE)
 6:   end for
 7:   for all tasks τ_i in curSys do
 8:     P_init = τ_i.core
 9:     for all cores P_a ≠ P_init do
10:       assign(τ_i, P_a)
11:       LN ← curSys
12:     end for
13:     Allocate(τ_i, P_init)
14:     for all cores P_a ≠ P_init do
15:       for all tasks τ_j on P_a do
16:         if τ_j.util ≥ τ_i.util then
17:           for all cores P_b ≠ P_a do
18:             Allocate(τ_j, P_b)
19:             Allocate(τ_i, P_a)
20:             LN ← curSys
21:             Allocate(τ_j, P_a)
22:             Allocate(τ_i, P_init)
23:           end for
24:         end if
25:       end for
26:     end for
27:   end for
28:   return LN
```

| Task | Period (ms) | WCET (ms) | Readers | Comm. size (Bytes) |
|------|-------------|-----------|---------|--------------------|
| $\tau_0$ | 10 | 1 | 1,3 | 256 |
| $\tau_1$ | 100 | 8 | 0,5 | 128 |
| $\tau_2$ | 400 | 117 | 4 | 48 |
| $\tau_3$ | 40 | 6 | 1,6 | 128 |
| $\tau_4$ | 20 | 7 | 2 | 48 |
| $\tau_5$ | 1000 | 394 | 6 | 256 |
| $\tau_6$ | 20 | 7 | 1,5 | 128 |

TABLE I: Task parameters for the example system

first search among candidate solutions. Candidate solutions for exploration are placed in a list $\mathcal{NA}$ sorted by increasing memory cost of the solution. The first solution (one with lowest cost) in $\mathcal{NA}$ is then further explored by branching (generating its neighbors).

The difficulty in exploration arises from the fact that an estimate of the quality of the solutions that may be found under a given branch (and therefore a possible pruning decision) is very difficult. The algorithm allows the exploration of solutions (neighbors) with both lower and higher costs than the current optimum to avoid getting stuck at a local optimum. However, to avoid infinite searches, the exploration is bounded by a condition on the number of iterations without improvement (line 34, Algorithm 3). It is also bounded by the size of $\mathcal{NA}$ which is at most equal to the number $n$ of tasks in the system (line 29, Algorithm 3). Essentially, the best $n$ unexplored candidates at any stage are kept in $\mathcal{NA}$. In our experiments, larger sizes for $\mathcal{NA}$ such as $2n$ or $4n$ do not improve the quality of the obtained solution but result in substantially longer execution times. To avoid loops, recently visited neighbors are discarded.

The solution space exploration is depicted in lines 21-37 of Algorithm 3. The first solution in $\mathcal{NA}$ (solution with minimum cost) is removed from the list (line 23) and considered as the new base ($curSys$) for further exploration. Lines 24-33 show the generation and exploration of neighbors. All feasible neighbors of the current base $curSys$ are generated (line 24). However, not all of them are further explored. A threshold value $Th$ is used as an acceptance criteria for new neighbors generated from $curSys$, which is set to be the memory cost of the last solution in $\mathcal{NA}$ (the unexplored solution with the $n$-th highest cost, as in line 31). Only solutions with lower costs than $Th$ are accepted for further exploration and added to $\mathcal{NA}$ (lines 27-28). The cost of any new solution is considered only after resource optimization (line 26), such as when comparing with $Th$ (line 27), or when comparing with the current optimum $curOpt$ (line 30).

If there are no more promising solutions to explore ($\mathcal{NA}$ becomes empty), or the solution is not improving for a certain number of iterations, or a solution with the desired quality ($tgtCost$, typically depending on the available RAM memory) is found, the algorithm terminates (lines 34-35).

## VI. AN ILLUSTRATIVE EXAMPLE

We provide a simple example to illustrate the operation of the GS, CASR, GS extended with wait-free (GS-WF), and MPA algorithms. The system consists of seven tasks to be partitioned on two cores, as shown in Table I. All critical sections have a WCET of 1ms, except that task $\tau_0$ has a duration of the critical sections as 0.15ms. Task deadlines are equal to their periods. The task-to-core assignment during the execution of the four algorithms is shown in Figure 2.

GS (Figure 2a) arranges tasks by density and assigns them to maximize the minimum slack. Task $\tau_5$ has the highest density and is assigned first. GS then attempts at assigning task $\tau_4$ to both cores. On $p_1$, the value of the smallest normalized slack after assigning task $\tau_4$ to $p_1$ is 0.389. Alternatively, assigning $\tau_4$ to $p_2$ achieves a minimum normalized slack of 0.65. Hence, GS chooses to assign $\tau_4$ to $p_2$. Following the same procedure, $\tau_6$ is assigned to $p_1$ then $\tau_2$, $\tau_3$, and $\tau_0$ are assigned to to $p_2$. Greedy slacker then fails to assign task $\tau_1$ to any core. At this point, the algorithm fails.

For CASR (Figure 2b), we assume $U_b = 0.858$ (the total task utilization divided by 2). CASR first assigns $\tau_5$ and $\tau_4$ to $p_1$ and $p_2$ respectively, similar to GS. Then, CASR allocates $\tau_6$ to its affine core $p_1$ and $\tau_2$ to its affine core $p_2$. Next, CASR assigns $\tau_3$ to $p_1$ because of the communication affinity, whereas GS would place it on $p_2$ because of the slack-based rule. At this point, the total task utilization of $p_1$ exceeds $U_b$. Thus, when allocating the next task $\tau_0$, $p_1$ loses its privilege as an affine core. CASR would try both $p_1$ and $p_2$, and allocates $\tau_0$ to $p_2$ to maximize the least normalized slack. Finally, $\tau_1$ is assigned to $p_2$, and the CASR succeeds in partitioning the task set.

GS-WF (Figure 2c) performs the first six task assignment (5, 4, 6, 2, 3, 0) in a similar way to GS. The algorithm then attempts and fails to assign $\tau_1$ using MSRP. At this point, a second attempt is made to assign $\tau_1$ using wait-free for all its global resources. This attempts succeeds with a memory penalty of 1664 bytes. The memory cost is calculated according to Equation (3), leveraging information about task parameters to reduce the buffer sizes.

The MPA algorithm (Figure 2d) starts by calculating the initial assignment urgencies according to Equation (8). Tasks are then sorted by their assignment urgency values. Since no task has been assigned, all assignment urgencies will be equal to zero. The task density is used as a tie-breaker and $\tau_5$ is selected as the first task to be assigned. Once $\tau_5$ is
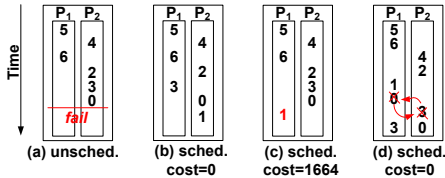
Fig. 2: Task assignment for the example system: (a) Greedy slacker, (b) CASR, (c) GS-WF, (d) MPA

assigned, all remaining tasks are schedulable on the two cores, therefore assignment urgencies depend solely on the memory costs resulting from sharing resources with $\tau_5$. Assigning any task to core $p_1$ has no memory cost since all shared resources defined until this stage will be local and managed using SRP. However, the costs of assigning tasks to $p_2$ will differ among tasks. The highest memory cost comes from $\tau_6$ (cost is 1152) and hence it has the highest assignment urgency of 1152 (by Equation (8)). $\tau_6$ will then be assigned to the core where it causes the least memory cost ($p_1$).

Next, the check on schedulability, performed by `TaskSort()`, reveals that tasks $\tau_4$ and $\tau_2$ are now unschedulable on $p_1$. Since these tasks can only be to assigned to one core, their assignment urgencies becomes $M_{max}$ (according to Equation (8)). All other tasks are schedulable on both cores and hence have lower AU values. Task $\tau_4$ is now at the top of the list $\mathcal{LT}$ (using its density to break the tie with task $\tau_2$) and is assigned immediately to core $p_2$ followed by task $\tau_2$. Phase 1 of the algorithm continues in a similar way, until the last task $\tau_3$ is assigned. All global resources use wait-free methods as the data consistency mechanism, and the cost of this initial solution is 1280 bytes.

Phase 2 of MPA starts by running the function `optimizeResources()`, which reduces the memory cost to 768 bytes by changing the protection mechanism for the resource written by $\tau_3$ to MSRP. This solution is the current optimal solution ($CurOpt$ in line 20 of Algorithm 3). Phase 2 then generates and examines neighbors to this solution. One of the feasible neighbors can be obtained by swapping tasks 0 and 3. This proves to be an effective move as `optimizeResources()` is able to reduce the memory cost to 0. Since this is an optimal solution in terms of memory, the algorithm terminates.

## VII. Experimental Results

In this section, the proposed algorithms are evaluated in terms of schedulability and memory cost (if applicable). We start by providing an overall evaluation in Section VII-A comparing GS with CASR, GS-WF and MPA in terms of schedulability. We also compare the memory costs of GS-WF and MPA. In Section VII-B, we study the effect of increasing the communication between tasks. Finally, in Section VII-C, we compare MPA with an exhaustive search to validate its ability to minimize memory cost.

We adopt a task generation scheme similar to the one used in [27]. We consider systems with 4 or 8 cores. The periods of the tasks are generated according to a log-uniform distribution from one of two different ranges [10, 100] ms and [3-33] ms. The average utilization of the task is selected from the set {0.05, 0.1, 0.2, 0.3}. The worst case execution time is

then derived from the values of the period and utilization. The critical section lengths are randomly generated in either [0.001, 0.1] ms or [0.001, 0.015] ms.

Tasks in the system share 1-40 resources. For each experiment, a resource sharing factor from the set {0.1, 0.25, 0.5, 0.75} is selected. The resource sharing factor represents the portion of tasks in the system sharing a given resource. A resource sharing factor of 0.1 means that each resource is shared by 10% of the tasks in the system. The tasks that share a certain resource are randomly selected and are independently generated from other resources. For each parameter configuration, 100 systems are randomly generated. The size of communication data was chosen among a set of possible values (in bytes) with probability values: 1 (with probability $p$ =10%), 4 ($p$ =20%), 24 ($p$ =20%), 48 ($p$ =10%),128 ($p$ =20%), 256 ($p$ =10%), and 512 ($p$ =10%). In all experiments, for MPA, the number of iterations is set to be $10n$ where $n$ is the number of tasks in the system.

### A. General Evaluation

For the first evaluation, we compare the algorithms with the greedy slacker algorithm (similar settings as in [27]). A variable number of tasks $n$ in the range [40, 76] is generated to be scheduled on 4 and 8 cores. The average utilization of each task is 0.1. Task periods, resources and critical sections are selected as defined earlier.

We compare the percentage of schedulable solutions obtained by GS with those obtained from: 1) a single run of CASR with the utilization bound $U_b = U_{\mathcal{T}}/m$ (denoted as CASR (single Ub) in Figure 3); and 2) selecting the best solution from multiple (five) runs of CASR with a set of $U_b$ values $\{0, 25\%, 50\%, 75\%, 100\%\}$ (CASR (multiple Ub) in the figure). The results are shown in Figures 3 and 4, where a clear schedulability improvement from CASR is demonstrated.

Table II, show results from the case with 8 cores, 4 shared resources and 25% of the tasks using each resource. The table illustrates the effectiveness of the recovery strategy. For each $U_b$ value, the table shows the results of the CASR execution for all the generated systems, as the percentage of systems that are found schedulable after the first and the second recovery stage (over all the schedulable ones). In each table entry, the first number is the percentage of systems successfully partitioned (schedulable) only after using the post-black-list $\mathcal{PBL}$ and the second number is the percentage of configurations found only after using the $\mathcal{BL}$ list and with empty $\mathcal{PBL}$. For systems with more than 68 tasks, CASR fails to find any feasible solution and the corresponding rows are omitted in Table II. The recovery strategy is more effective when the utilization bound value is large or the system utilization is high.

Figure 3 also shows the improvement obtained using wait-free methods. The GS-WF algorithm, as expected, contributes to a significant increase in the schedulability of GS. MPA performs even better, scheduling more systems at high utilizations . Figure 4 compares the additional memory cost needed by wait-free methods for both GS-WF and MPA. The figure shows that for systems with a small number of tasks, GS-WF performs slightly better. However, as the number of tasks (and hence utilization) increases, MPA outperforms GS-WF. For systems with 60 tasks (or an average core utilization of about 0.75), these algorithms perform approximately equally; at higher utilization, MPA tends to have a lower memory cost.

| #tasks | $U_b = 0$ | $U_b = 25\%$ | $U_b = 50\%$ | $U_b = 75\%$ | $U_b = 100\%$ |
|---|---|---|---|---|---|
| 40 | (0, 0) | (0, 0) | (0, 0) | (3, 2) | (27, 3) |
| 43 | (0, 0) | (0, 0) | (0, 0) | (3, 2) | (24, 5) |
| 46 | (0, 0) | (0, 0) | (0, 1) | (7, 10) | (30, 4) |
| 50 | (0, 0) | (0, 0) | (0, 0) | (9, 5) | (35, 6) |
| 53 | (0, 1) | (0, 1) | (0, 1) | (7, 12) | (36, 7) |
| 56 | (0, 4) | (0, 5) | (1, 5) | (17, 8) | (40, 7) |
| 60 | (1, 11) | (3, 13) | (0, 19) | (19, 15) | (37, 11) |
| 63 | (10, 31) | (8, 30) | (4, 30) | (38, 21) | (53, 17) |
| 66 | (0, 0) | (0, 0) | (25, 25) | (30, 20) | (75, 25) |

TABLE II: Percentage of systems recovered by $\mathcal{PBL}$ and $\mathcal{BL}$ in CASR among the schedulable ones (8 cores, 4 resources)
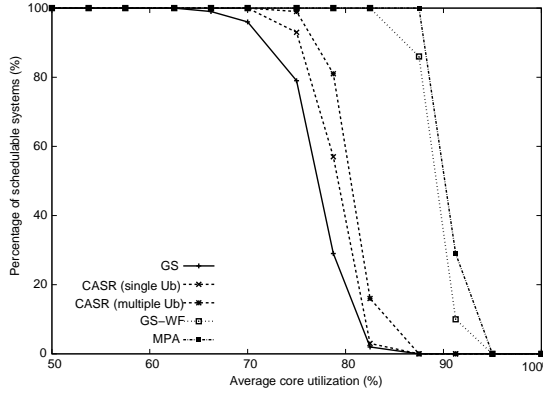


Fig. 3: Comparison of schedulability (8 cores, 4 resources)

## B. Effect of Other Parameters

In many applications, tasks can share a larger number of resources than in the communication scheme used in the first experiment. For the basic setting, we use a system with 4 cores and 20 resources. The number of tasks is varied in [20-40]. All other parameter settings are kept similar to those in Section VII-A. Figures 5 and 6 show the results of this experiment. The general trend in schedulability remains the same with CASR performing better than GS and MPA significantly improving upon GS-WF in both schedulability and memory cost.

To evaluate the applicability of the algorithms to systems with a small or large number of resources, we performs the same experiment. However, this time we fix the number of tasks at 28 and vary the number of resources in the range 1-40. The schedulability results are shown in Figure 7. Both GS-
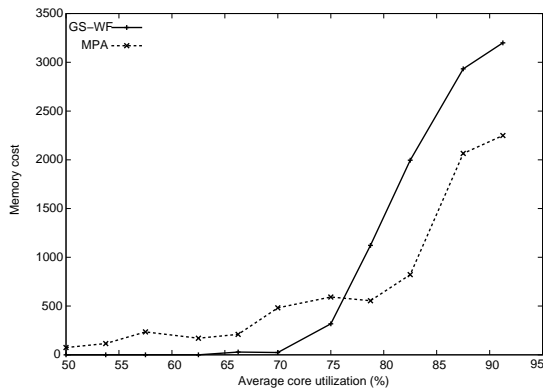


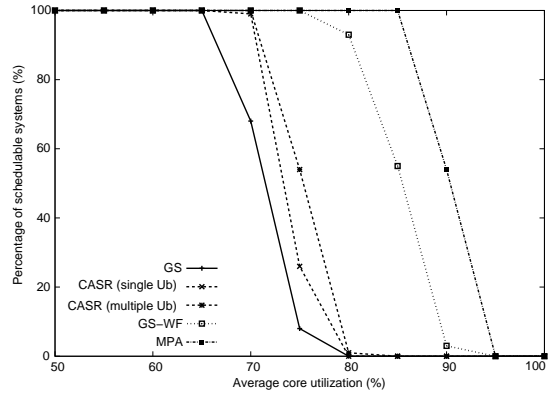Fig. 4: Comparison of memory cost (8 cores, 4 resources)



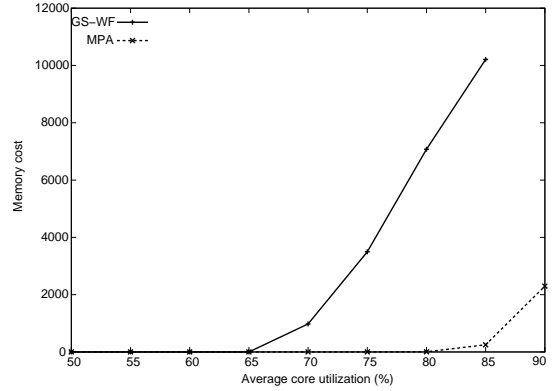Fig. 5: Comparison of schedulability (4 cores, 20 resources)



Fig. 6: Comparison of memory cost (4 cores, 20 resources)

WF and MPA always schedule all tasks and are thus omitted from the figures. Some systems can have some heavily shared resources. (i.e. high resource sharing factor). To evaluate the effect of changing the resource sharing factor, another experiment was performed using the same parameter settings while keeping the number of resources at 20 and changing the resource sharing factor (rsf). Table III present the result of this experiment.

These experiments show that CASR schedules more tasks than GS. As communication increases, CASR outperforms GS while requiring no additional memory resources. When the use of wait-free resources is permitted, the MPA algorithm performs significantly better than both GS and GS-WF as tasks communicate more often. MPA alsp succeeds in keeping the memory usage to minimum requiring just 841 bytes of memory at rsf=0.75 compared to 11934 for GS-WF.

In all experiments, schedulability analysis was performed using the approach in [17]. An improved analysis based on ILP formulations was presented in [28]. This approach can be

| | rsf=0.1 | rsf=0.25 | rsf=0.5 | rsf=0.75 |
|---|---|---|---|---|
| GS | 100% | 68% | 0% | 0% |
| CASR (single Ub) | 100% | 94% | 0% | 0% |
| CASR (multiple Ub) | 100% | 100% | 0% | 0% |
| GS-WF | 100%/0 | 100%/1136 | 98%/9628 | 67%/11493 |
| MPA | 100%/0 | 100%/0 | 100%/104 | 100%/841 |

TABLE III: Schedulability and average memory cost (GS-WF and MPA only, in bytes) for different resource uses
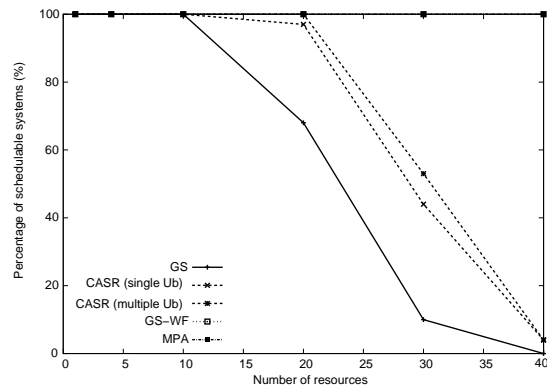
Fig. 7: Comparison of schedulability with a variable number of resources

combined with the algorithms presented, however, since the algorithms check schedulability, using the new analysis in our algorithms can lead to longer run times. We compared the use of both schedulability approaches for systems consisting of 14 tasks with an average utilization of 0.25 per task (overall utilization =87.5%) while keeping other parameters similar to their defaults in the previous experiments. Using MPA, both approaches scheduled 73% of the generated systems. The ILP-based approach reduced the average memory requirement from 11.2 bytes to 0 by judging more solutions to be schedulable thus allowing MPA more alternatives. However the ILP-based approach was on average 311 times slower. Increasing the task count to 20, the ILP-based approach was about 14000 times slower. Therefore, we used the approach in [17] when generating the experimental results in this section, however, using the improved approach in [28] is a viable option.

### C. MPA vs. Exhaustive Search

Finally, we use exhaustive search to validate the quality of the proposed heuristic MPA. Exhaustive search can be time consuming due to the fact that partitioning is a bin packing problem which is NP-hard. Therefore, we have restricted our comparison with exhaustive partitioning to small systems with 12-16 tasks and 3 cores. To compare the effectiveness of the MPA algorithm in terms of memory cost, we computed the results of exhaustive search for partitioning 100 systems at a utilization of 80%. The average memory cost was 19% worse for MPA. In terms of schedulability, MPA successfully scheduled all 100 systems. However, MPA runs significantly faster than the exhaustive search: the average runtime for MPA is about 6.3 seconds, while the exhaustive search takes about 2988 seconds on average.

### VIII. CONCLUSION

In this paper, we provide a method improving on existing task placement algorithms for systems that use MSRP to protect shared resources. Furthermore, we leverage an additional opportunity provided by wait-free methods as an alternative mechanism to protect data consistency for shared buffers. The selective use of wait-free methods significantly extends the range of schedulable systems at the cost of memory, as shown by experiments on random task sets.

REFERENCES

[1] The AUTOSAR consortium. *The AUTOSAR Standard, specification version 4.1.* [Online] http://www.autosar.org.

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

[3] T. Baker. "A stack-based resource allocation policy for realtime processes." In *Proc. 11th IEEE Real-Time Systems Symposium*, 1990.

[4] S. Baruah. "The partitioned EDF scheduling of sporadic task systems." In *Proc. IEEE Real-Time Systems Symposium*, 2011.

[5] S. Baruah and E. Bini. "Partitioned scheduling of sporadic task systems: An ILP based approach." In *Proc. Conf. on Design and Architectures for Signal and Image Processing*, 2008.

[6] S. Baruah and N. Fisher. "The partitioned multiprocessor scheduling of sporadic task systems." In *Proc. 26th IEEE Real-Time Systems Symposium*, 2005.

[7] B. Brandenburg, J. Calandrino, A. Block, H. Leontyev, and J. Anderson. "Real-time synchronization on multiprocessors: To block or not to block, to suspend or spin?" In *Proc. 14th IEEE Real-Time and Embedded Technology and Applications Symposium*, 2008.

[8] B. Brandenburg. "A Fully Preemptive Multiprocessor Semaphore Protocol for Latency-Sensitive Real-Time Applications In *Proc. 25th Euromicro Conference on Real-Time Systems*, 2013.

[9] E. Bini and G. Buttazzo. "Measuring the Performance of Schedulability Tests." *Real-Time Systems*, 30(1-2):129–154, May 2005.

[10] B. Chattopadhyay and S. Baruah. "A lookup-table driven approach to partitioned scheduling." In *Proc. 17th IEEE Real-Time and Embedded Technology and Applications Symposium*, 2011.

[11] J. Chen and A. Burns. "A fully asynchronous reader/write mechanism for multiprocessor real-time systems." *Technical Report YCS 288, Department of Computer Science, University of York*, May 1997.

[12] J. Chen and A. Burns. "Loop-free asynchronous data sharing in multiprocessor real-time systems based on timing properties." In *Proc. 6th Int. Conference on Real-Time Computing Systems and Applications*, 1999.

[13] R. Davis and A. Burns. "A Survey of Hard Real-Time Scheduling for Multiprocessor Systems." *ACM Computing Surveys*, 43(4), Oct. 2011.

[14] S. Dhall and C. Liu. "On a real-time scheduling problem." *Operations Research*, 26(1):127–140, 1978.

[15] N. Fisher. "The multiprocessor real-time scheduling of general task systems." *Ph.D. dissertation, The University of North Carolina at Chapel Hill*, 2007.

[16] N. Fisher and S. Baruah. "The Partitioned Scheduling of Sporadic Tasks According to Static Priorities." In *Proc. 18th Euromicro Conference on Real-Time Systems*, 2006.

[17] P. Gai, G. Lipari, and M. D. Natale. "Minimizing memory utilization of real-time task sets in single and multi-processor systems-on-a-chip." In *Proc. 22nd IEEE Real-Time Systems Symposium*, 2001.

[18] G. Han, H. Zeng, M. Di Natale, X. Liu, and W. Dou. "Experimental Evaluation and Selection of Data Communication Mechanisms in Multicore Platforms." *IEEE Transactions on Industrial Informatics*, vol.PP, no.99, pp.1, 2013.

[19] M. Herlihy. A methodology for implementing highly concurrent structures. In *Proc. of ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, 1990.

[20] H. Huang, P. Pillai, and K. G. Shin. "Improving wait-free algorithms for interprocess communication in embedded real-time systems." In *Proc. USENIX Annual Technical Conference*, 2002.

[21] D. Johnson. "Near-optimal bin packing algorithms." *Ph.D. dissertation, Massachusetts Institute of Technology*, 1973.

[22] K. Lakshmanan, D. de Niz, and R. Rajkumar. "Coordinated task scheduling, allocation and synchronization on multiprocessors." In *Proc. 30th IEEE Real-Time Systems Symposium*, 2009.

[23] F. Nemati, T. Nolte, and M. Behnam. "Partitioning real-time systems on multiprocessors with shared resources." In *Proc. Conference on Principles of distributed systems*, 2010.

[24] R. Rajkumar. "Real-time synchronization protocols for shared memory multiprocessors." In *Proc. 10th International Conference on Distributed Computing Systems*, 1990.

[25] C. Sofronis, S. Tripakis, and P. Caspi. "A memory-optimal buffering protocol for preservation of synchronous semantics under preemptive scheduling." In *Proc. International conference on Embedded software*, 2006.

[26] G. Wang, M. Di Natale, and A. Sangiovanni-Vincentelli. "Improving the size of communication buffers in synchronous models with time constraints." *IEEE Transactions on Industrial Informatics*, 5(3):229–240, Aug. 2009.

[27] A. Wieder and B. Brandenburg. "Efficient Partitioning of Sporadic Real-Time Tasks with Shared Resources and Spin Locks." In *Proc. 8th IEEE International Symposium on Industrial Embedded Systems*, 2013.

[28] A. Wieder and B. Brandenburg. "On Spin Locks in AUTOSAR: Blocking Analysis of FIFO, Unordered, and Priority-Ordered Spin Locks." In *Proc. 34th IEEE Real-Time Systems Symposium*, 2013.

[29] H. Zeng and M. Di Natale. "Mechanisms for Guaranteeing Data Consistency and Flow Preservation in AUTOSAR Software on Multi-Core Platforms." In *Proc. 6th IEEE International Symposium on Industrial Embedded Systems*, 2011.