

# Optimal Fixed Priority Scheduling in Multi-Stage Multi-Resource Distributed Real-Time Systems

Niraj Kumar, Chuanchao Gao, and Arvind Easwaran

School of Computer Science and Engineering, Nanyang Technological University Singapore

Email: niraj.kumar@ntu.edu.sg, gaoc0008@e.ntu.edu.sg, arvinde@ntu.edu.sg

**Abstract**—This work studies fixed priority (FP) scheduling of real-time jobs with end-to-end deadlines in a distributed system. Specifically, given a multi-stage pipeline with multiple heterogeneous resources of the same type at each stage, the problem is to assign priorities to a set of real-time jobs with different release times to access a resource at each stage of the pipeline subject to the end-to-end deadline constraints. Note, in such a system, jobs may compete with different sets of jobs at different stages of the pipeline depending on the job-to-resource mapping. To this end, following are the two major contributions of this work. We show that an OPA-compatible schedulability test based on the *delay composition algebra* can be constructed, which we then use with an *optimal priority assignment* algorithm to compute a *priority ordering*. Further, we establish the versatility of *pairwise priority assignment* in such a multi-stage multi-resource system, compared to a total priority ordering. In particular, we show that a pairwise priority assignment may be feasible even if a priority ordering does not exist. We propose an integer linear programming formulation and a scalable heuristic to compute a pairwise priority assignment. We also show through simulation experiments that the proposed approaches can be used for the holistic scheduling of real-time jobs in edge computing systems.

**Index Terms**—Real-time Distributed Systems, Delay Composition Algebra, Optimal Priority Assignment, Edge Computing

## I. INTRODUCTION

*Fixed priority* (FP) real-time scheduling has attracted significant attention in academia and broader acceptance in the industry due to its predictability and low runtime overhead [1], [2]. FP uniprocessor scheduling is NP-hard in general. However, due to the *optimal priority assignment* (OPA) proposed by Audsley [3], the hardness is attributed to the schedulability test. More specifically, OPA is an optimal<sup>1</sup> FP scheduling algorithm to compute a *priority ordering*<sup>2</sup> of preemptive sporadic tasks with arbitrary release times on a uniprocessor. An optimal priority assignment with OPA only requires  $O(n^2)$ , for  $n$  tasks, invocations to an OPA-compatible<sup>3</sup> schedulability test.

Schedulability tests for distributed real-time systems are inherently more complex. Thus, a common approach to analyze the schedulability in distributed systems is to decompose the end-to-end deadline to a per-stage deadline and then independently analyze each stage as a uniprocessor system [4], [5]. Alternatively, *delay composition algebra* (DCA) is proposed

to transform a multi-stage system into an equivalent single-stage system and then analyze schedulability by applying a uniprocessor schedulability test [6], [7]. As an intermediate step, a delay composition rule is also presented to compute an upper bound of end-to-end delay for each job in a multi-stage pipeline. In this work, we use the delay composition rule of DCA only; henceforth, the term DCA refers to the same.

This work aims to compute an FP schedule for real-time jobs with release offsets and end-to-end deadlines in a distributed system. In particular, we focus on a multi-stage multi-resource (MSMR) system in which each stage corresponds to a resource type with multiple heterogeneous resources of the same type. Each real-time job accesses one of the resources at every stage in the pipeline and must exit the pipeline within the specified deadline. A simplified version of this problem, particularly multi-stage single-resource scheduling, reduces to flow shop scheduling of jobs with a deadline, which is known to be an NP-hard problem in general [8]. Furthermore, the problem addressed in this work is more complex as we consider a generic model with heterogeneity among resources of the same type and jobs with arbitrary release times. Moreover, jobs are allowed to compete with distinct sets of jobs at different stages of the pipeline depending on the job-to-resource mapping. The major contributions of this work are summarized as follows.

- To our knowledge, this work is the first to use OPA and DCA for optimal FP scheduling of real-time jobs with end-to-end deadlines in an MSMR system. Specifically, we show that an OPA-compatible schedulability test, based on DCA, can be constructed to compute an optimal priority ordering<sup>4</sup> of real-time jobs in an MSMR system. The proposed algorithm can be applied for both preemptive and non-preemptive FP scheduling in an MSMR system for a given job-to-resource mapping.
- We establish that *pairwise priority assignment*<sup>5</sup> and priority ordering are not necessarily equivalent in MSMR systems. In particular, we show that a feasible pairwise priority assignment may exist for a set of real-time jobs that does not admit a priority ordering. Further, we propose an ILP formulation and a *deadline-monotonic & repair*

**Funding:** This work was supported by the MoE Tier-2 grant MOE-T2EP20221-0016.

<sup>1</sup>An FP scheduling algorithm  $\mathcal{P}$  is said to be *optimal* with respect to a schedulability test  $\mathcal{S}$  provided any schedulable taskset by some FP scheduling algorithm (under  $\mathcal{S}$ ) is also schedulable by  $\mathcal{P}$ .

<sup>2</sup>A permutation or a total order of tasks that indicates priority.

<sup>3</sup>Discussed, along with delay composition algebra, in Section III.

<sup>4</sup>It is worth mentioning that DCA is not originally intended to compute a priority assignment. Instead, given a priority assignment, DCA bounds the end-to-end delay of each job in a multi-stage system to obtain an equivalent single-stage system for schedulability analysis.

<sup>5</sup>Priorities are assigned to each pair of jobs that share at least one resource in the pipeline as opposed to priority ordering in which a total order is established among all jobs.

based heuristic to compute an optimal and a sub-optimal pairwise priority assignment, respectively.

- Through simulation experiments, we show that the proposed approaches can be applied for holistic scheduling of real-time jobs with end-to-end deadlines in edge computing systems. Notably, these approaches alleviate the challenges associated with computing and imposing a virtual deadline, as end-to-end deadlines are used directly for the priority assignment.

## II. SYSTEM MODEL AND PROBLEM STATEMENT

Consider a multi-stage multi-resource (MSMR) system with  $N$  stages in which each stage corresponds to a resource of a specific type (for instance, edge servers). Moreover, multiple resources of heterogeneous capabilities are available at each stage. Let  $\mathcal{J}$  be the given set of  $n$  real-time jobs. Each job  $J_i$  is specified with the following parameters (i)  $A_i$ : arrival time, (ii)  $P_{i,j}$ : processing time<sup>6</sup> of  $J_i$  in stage  $S_j$ , (iii)  $D_i$ : end-to-end deadline, and (iv)  $R_{i,j}$ : the resource to which  $J_i$  is mapped in stage  $S_j$ . Let  $\mathcal{M}_{i,j}$  be the set of all jobs mapped with  $J_i$  in stage  $S_j$  (and thus, compete for  $R_{i,j}$ ). Moreover, the set of all jobs that compete with  $J_i$  is  $\mathcal{M}_i = \cup_j \mathcal{M}_{i,j}$ .

Let  $\Delta_i$  denote the *end-to-end delay* of  $J_i$ , which is the time it exits the pipeline since its arrival. Given an MSMR system and a set  $\mathcal{J}$  of  $n$  real-time jobs  $J_i \langle A_i, \{P_{i,j}\}, D_i, \{R_{i,j}\} \rangle$ , then the problem of computing a

- P1. **Priority Ordering** is to devise an OPA-compatible schedulability test  $\mathcal{S}^{DCA}$  based on DCA and then use it with OPA to assign a unique *global* (i.e., valid across all stages) priority  $\rho_i \in [1, n]$  to each job  $J_i$  such that the deadline constraint is satisfied  $\Delta_i \leq D_i$ .
- P2. **Pairwise Priority Assignment** is to assign a *global* fixed priority  $J_i > J_k$  (indicating  $J_i$  has a higher priority than  $J_k$ ) or  $J_k > J_i$  for each job pair  $\langle J_i, J_k \rangle | J_k \in \mathcal{M}_i$  such that  $\Delta_i \leq D_i$ .

If a feasible priority ordering exists, solving P1 is sufficient (since OPA is optimal); however, otherwise, solving P2 is required (as a pairwise priority assignment may still be feasible, further details in Section V).

We define a *segment* corresponding to a job pair  $\langle J_i, J_k \rangle$  as one of the longest sequences of consecutive stages for which  $J_i$  and  $J_k$  are mapped to the same resources. A pipeline may have multiple segments corresponding to each job pair. Let  $m_{i,k}$  denote the number of segments corresponding to the job pair  $\langle J_i, J_k \rangle$ . For instance,  $m_{i,b} = 0$  in Figure 1(a),  $m_{i,b} = 1$  in Figures 1(b)-(d), and  $m_{i,b} = 2$  in Figure 1(e).

Let  $\mathcal{H}_i$  (and  $\mathcal{L}_i$ ) denote the set of higher (and lower) priority jobs of  $J_i$ . Also, define  $\mathcal{Q}_i = \mathcal{H}_i \cup J_i$ . We assume that any job  $J_k$  which cannot interfere  $J_i$ , i.e.,  $[A_k, A_k + D_k]$  and  $[A_i, A_i + D_i]$  do not overlap, is already excluded from  $\mathcal{H}_i$  (and  $\mathcal{L}_i$ ). We assume a lower value (of  $\rho_i$ ) indicates a higher priority. Both preemptive and non-preemptive FP scheduling are considered in this work. Firstly, essential background is discussed in Section III. Subsequently, problems P1 and P2 are addressed in

<sup>6</sup>We use processing time in a generic sense, thus, for instance, for a networking resource it means transmission time.

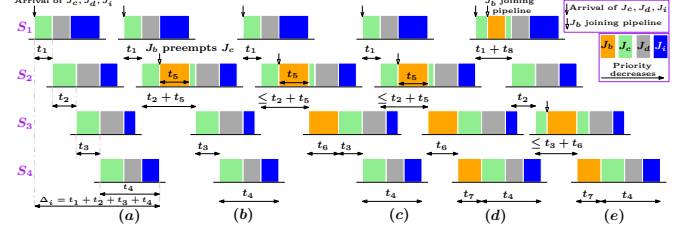


Fig. 1: End-to-end delay of  $J_i$  in a four-stage pipeline. (a)  $\Delta_i = t_1 + t_2 + t_3 + t_4$ . (b) One job-additive term due to  $J_b$  as it delays  $J_i$  at one stage only,  $\Delta_i = t_1 + t_2 + t_5 + t_3 + t_4$ . (c) & (d) Two job-additive terms if  $J_b$  delays  $J_i$  for two or more stages. (e)  $m_{i,b} = 2$  and up to two job-additive terms of  $J_b$  corresponding to each segment.

Sections IV and V, respectively. The commonly used notations are listed in the Table I.

## III. BACKGROUND

### A. Delay Composition Algebra (DCA)

The end-to-end delay bound of a job (henceforth, of  $J_i$ ) in a multi-stage single-resource pipeline<sup>7</sup> can be computed using DCA as explained in the following (reproduced from [6]).

Let  $t_{i,x}$  denote the  $x^{th}$  maximum stage-processing time of  $J_i$ , for instance,  $t_{i,1} = \max_j \{P_{i,j}\}$ . Then, in preemptive scheduling, as illustrated in Figure 1(a),  $\Delta_i$  can be bounded by (i) *stage-additive component*: corresponding to each stage, except the last, to incorporate the overlapped execution of jobs of  $\mathcal{Q}_i$ , and (ii) *job-additive component*: corresponding to each higher priority job, as it may delay  $J_i$  while leaving the pipeline. Let  $\mathcal{H}_i^a \subseteq \mathcal{H}_i$  be the set of higher priority jobs that join the pipeline after  $J_i$  i.e. a job  $J_k \in \mathcal{H}_i^a$  provided  $J_k \in \mathcal{H}_i$  and  $A_i < A_k$ . Then,  $J_k \in \mathcal{H}_i^a$  will also contribute a job-additive component while joining the pipeline. Referring to Figure 1(c) for an illustration,  $J_b$  contributes two terms, one while joining the execution pipeline of  $J_i$  at  $S_2$  and another while leaving at  $S_3$ . Thus, combining these interferences, the end-to-end delay of  $J_i$  in preemptive scheduling is (Eq. 1, [6])

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} t_{k,1} + \sum_{J_k \in \mathcal{H}_i^a} t_{k,2} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{P_{k,j}\} \quad (1)$$

In non-preemptive scheduling, one lower priority job may delay  $J_i$  in each stage; however, jobs  $J_k \in \mathcal{H}_i^a$  cannot preempt  $J_i$ . Thus, in this case, the delay of  $J_i$  is (Eq. 2, [6])

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} t_{k,1} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{P_{k,j}\} + \sum_{j=1}^N \max_{J_k \in \mathcal{L}_i} \{P_{k,j}\} \quad (2)$$

The extended version of DCA [7] can be used to analyze the delay in an MSMR system with multiple segments for each job pair. We define  $\tilde{p}_{k,j} = P_{k,j}$  if  $J_i$  and  $J_k$  share the  $j^{th}$  stage, otherwise  $\tilde{p}_{k,j} = 0$ . Correspondingly, we also define  $\tilde{t}_{k,x}$ , for instance,  $\tilde{t}_{k,1} = \max_j \{\tilde{p}_{k,j}\}$ . Analogous to Eq. 1, a higher priority job  $J_k$  contributes two job-additive terms

<sup>7</sup>All jobs compete with each other for every resource in every stage.

TABLE I: Notations

Symbol	Meaning
$N$	No. of stages in the pipeline
$S_j$	$j^{th}$ stage of the pipeline
$J_i$	$i^{th}$ job in the set of jobs $\mathcal{J}$
$D_i/\Delta_i$	End-to-end deadline/delay of $J_i$
$P_{i,j}$	Processing time of $J_i$ in stage $S_j$
$\tilde{p}_{k,j}$	$P_{k,j}$ if $J_i$ and $J_k$ share $S_j$ , else 0
$\mathcal{H}_i/\mathcal{L}_i$	Set of higher/lower priority jobs of $J_i$
$\mathcal{Q}_i$	Set of higher priority jobs of $J_i$ and itself ( $\mathcal{H}_i \cup J_i$ )
$t_{i,x}$	$x^{th}$ maximum stage-processing time of $J_i$
$\tilde{t}_{k,x}$	$t_{k,x}$ but computed only for the stages shared with $J_i$
$m_{i,k}$	No. of segments corresponding to job pair $\langle J_i, J_k \rangle$
$w_{i,k}$	No. of job-additive components of $J_k$ that delays $J_i$
$\mathcal{M}_{i,j}/\mathcal{M}_i$	Set of jobs competing $J_i$ in $S_j$ /in pipeline

corresponding to each of the  $m_{i,k}$  segments in a preemptive scheduling, and thus, the delay is (Eq. 1, [7])

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} 2m_{i,k}\tilde{t}_{k,1} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{\tilde{p}_{k,j}\} \quad (3)$$

Likewise, Eq. 2 is rewritten to (Eq. 5, [7])

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} m_{i,k}\tilde{t}_{k,1} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{\tilde{p}_{k,j}\} + \sum_{j=1}^N \max_{J_k \in \mathcal{L}_i} \{\tilde{p}_{k,j}\} \quad (4)$$

#### B. Optimal Priority Assignment (OPA)

OPA computes a priority ordering of  $n$  tasks by assigning priorities from  $n$  (lowest) to 1 (highest) in that order. The current priority is assigned to one of the feasible tasks ( $\Delta_i \leq D_i$ ), which is yet to be assigned a priority. If no task is feasible for a priority, the taskset is declared unschedulable. The process is repeated until all tasks are assigned a priority or no feasible task is found for a priority.

A schedulability test is used with OPA to determine the feasibility of a task for a priority. However, the schedulability test must be OPA-compatible [2], [9]. Specifically, a schedulability test is OPA-compatible, provided it satisfies the following three necessary and sufficient conditions. The **first (second) condition** requires that the schedulability of a task may depend on the set of higher (lower) priority tasks but not on their relative priority. Whereas the **third condition** states that if the priorities of two tasks with adjacent priorities are swapped, then the task being assigned the lower (higher) priority remains unschedulable (schedulable) if it was unschedulable (schedulable) earlier. In the following, we use OPA to assign priority to the real-time jobs.

### IV. OPTIMAL FIXED PRIORITY SCHEDULING IN DISTRIBUTED REAL-TIME SYSTEMS

#### A. $\mathcal{S}^{DCA}$ – A schedulability test based on DCA

We define  $\mathcal{S}^{DCA}(J_i, \mathcal{H}_i, \mathcal{L}_i)$  to be a schedulability test based on DCA to determine the feasibility of job  $J_i$  given  $\mathcal{H}_i$  and  $\mathcal{L}_i$ . However, firstly, we refine the DCA analysis to ensure OPA-compatibility and to reduce the pessimism.

**Observation IV.1.**  $\mathcal{S}^{DCA}$  based on Eq. 1 is OPA-compatible.

*Proof.* The first two conditions of OPA-compatibility are satisfied trivially as only the set of higher priority jobs is required (and not the relative priority among them) to compute  $\Delta_i$ . Assume  $J_c$  and  $J_d$  are two jobs with priorities  $\rho_i - 1$  and  $\rho_i + 1$ , respectively. Then, the following three cases must be investigated to verify the third condition. Case 1. Priorities of  $J_i$  and  $J_c$  are swapped.  $J_i$  gets a higher priority. As a result,  $\mathcal{Q}_i \leftarrow \mathcal{Q}_i \setminus J_c$ , and  $\mathcal{H}_i^a \leftarrow \mathcal{H}_i^a \setminus J_c$  if  $J_c \in \mathcal{H}_i^a$ , and thus,  $\Delta_i$  would decrease. Case 2. Priorities of  $J_i$  and  $J_d$  are swapped.  $J_i$  gets a lower priority. Thus,  $\mathcal{Q}_i \leftarrow \mathcal{Q}_i \cup J_d$  and  $\mathcal{H}_i^a \leftarrow \mathcal{H}_i^a \cup J_d$  if  $A_i < A_d$ , and hence,  $\Delta_i$  would increase. Case 3. Priorities of two jobs with adjacent priorities are swapped such that both are either higher or lower priority jobs than  $J_i$ . Both  $\mathcal{Q}_i$  and  $\mathcal{H}_i^a$  remains unchanged in this case, and so does  $\Delta_i$ . Thus, the third condition is also satisfied.  $\square$

**Observation IV.2.**  $\mathcal{S}^{DCA}$  based on Eq. 2 is OPA-incompatible.

For this observation, consider the following example, for which the third condition of OPA-compatibility is violated.

**Example 1.** Consider a distributed system of three stages and four jobs:  $J_1 \langle 5, 7, 15 \rangle^8$ ,  $J_2 \langle 7, 9, 17 \rangle$ ,  $J_3 \langle 6, 8, 30 \rangle$ , and  $J_4 \langle 2, 4, 3 \rangle$ . Assuming  $J_1, J_2, J_3, J_4$  to be the priority ordering from highest to lowest,  $\Delta_2 = 92$ , however if the priorities of  $J_2$  and  $J_3$  are swapped then  $\Delta_2$  reduces to 87.

Observations IV.1 and IV.2 hold for Eqs. 3 and 4, respectively. Since  $\mathcal{L}_i \subseteq \mathcal{J} \setminus J_i$ , we get the following, more pessimistic but OPA-compatible, end-to-end delay for non-preemptive scheduling (by using the set  $\mathcal{J} \setminus J_i$  rather than  $\mathcal{L}_i$  in the  $3^{rd}$  term of Eq. 4)

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} m_{i,k}\tilde{t}_{k,1} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{\tilde{p}_{k,j}\} + \sum_{j=1}^N \max_{J_k \in \mathcal{J} \setminus J_i} \{\tilde{p}_{k,j}\} \quad (5)$$

We also note pessimism in Eq. 3 as one job-additive term is sufficient for a segment of one stage as the job arrives and leaves at the same stage (also illustrated in Figures 1(b) and 1(e)). Let  $J_i$  and  $J_k$  share  $u_{i,k}$  segments of exactly one stage each and  $v_{i,k}$  segments of two or more stages. Then, the maximum number of stage-processing terms of  $J_k \in \mathcal{H}_i$  that contributes to  $\Delta_i$  is  $w_{i,k} = u_{i,k} + 2v_{i,k}$ . Assuming  $w_{i,i} = 1$ , Eq. 3 can be rewritten as

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} \sum_{x=1}^{w_{i,k}} \tilde{t}_{k,x} + \sum_{j=1}^{N-1} \max_{J_k \in \mathcal{Q}_i} \{\tilde{p}_{k,j}\} \quad (6)$$

The delay bound reported by (refined) DCA is used in  $\mathcal{S}^{DCA}$  as the value of  $\Delta_i$  to determine feasibility. That is to say,  $\mathcal{S}^{DCA}$  determines  $\Delta_i$  using Eq. 5 (or Eq. 6) for non-preemptive (or for preemptive) scheduling, and  $J_i$  is deemed feasible under this priority provided  $\Delta_i \leq D_i$ .

#### B. Optimal Priority Assignment based on $\mathcal{S}^{DCA}$ (OPDCA)

OPDCA, Algorithm 1, assigns optimal priority (based on OPA) using  $\mathcal{S}^{DCA}$ . If no task, which is yet to be assigned

<sup>8</sup>Values inside  $\langle \dots \rangle$  indicate respective stage-processing times.

---

**Algorithm 1: OPDCA**


---

```

1 Set  $\mathcal{U} \leftarrow \mathcal{J}$ , and  $\ell \leftarrow \phi$ 
2 for each  $prior = n$  to 1 do
3    $assigned \leftarrow FALSE$ 
4   for each  $J_i \in \mathcal{U}$  do
5      $\mathcal{H}_i \leftarrow \mathcal{U} \setminus J_i$ ,  $\mathcal{L}_i \leftarrow \ell$ 
6     if  $S^{DCA}(J_i, \mathcal{H}_i, \mathcal{L}_i)$  is true then
7        $\rho_i \leftarrow prior$ ;  $\mathcal{U} \leftarrow \mathcal{U} \setminus J_i$ ;  $\ell \leftarrow \ell \cup J_i$ 
8        $assigned \leftarrow TRUE$ ; Break
9   if  $assigned$  is  $FALSE$  then
10    Declare  $\mathcal{J}$  infeasible

```

---

a priority, is feasible for the current priority level  $prior$ , in Steps 4-8,  $\mathcal{J}$  is declared infeasible in Step 10, otherwise  $prior$  is assigned to a feasible task in each iteration of the outer *for* loop (Steps 2-10). In the worst case, for each iteration of the outer *for* loop, the inner *for* loop is executed for all jobs that are yet to be assigned a priority. Thus,  $S^{DCA}$  is called  $O(n^2)$  times and the complexity of each call is  $O(nN)$ . As a result, the overall complexity of Algorithm 1 is  $O(n^3N)$ .

**Observation IV.3.** *OPDCA is an optimal FP scheduling algorithm to compute a priority ordering with respect to  $S^{DCA}$ .*

The proof of this observation follows from the optimality proof of OPA [9], [10]. If there exists any FP scheduling algorithm  $\mathcal{P}$  that can compute a priority ordering with respect to  $S^{DCA}$ , OPDCA can also compute one such priority ordering since in any feasible priority ordering, at least one job is feasible for each priority. Moreover, OPA is also optimal for non-preemptive scheduling [2]; the same holds for OPDCA since  $S^{DCA}$  (based on Eq. 5) is OPA-compatible.

## V. PAIRWISE PRIORITY ASSIGNMENT

A total priority ordering is suitable for single-stage or multi-stage single-resource systems. However, such a total priority ordering in MSMR systems mandates relative priority assignment even to jobs that may not share a single resource in the pipeline. For instance, the priority of  $J_1$  with respect to  $J_4$  (and vice-versa) is inconsequential for the given job-to-resource mapping shown in Figure 2(a).

**Observation V.1.** *A pairwise priority assignment can exist for a set of jobs that does not admit a total priority ordering.*

Consider the distributed system used in Example 1 with deadlines  $\{D_i\} = \{60, 55, 55, 50\}$ . Furthermore, assume preemptive scheduling and the same arrival time for all jobs. Then, there does not exist a total priority ordering for the job-to-resource mapping shown in Figure 2(a); however, a pairwise priority assignment, such as shown in Figure 2(b), exists.

The pairwise priority assignment is more flexible; for instance, in Figure 2(b),  $J_3$  is a lower priority (compared to  $J_4$ ) job at  $S_1$ , but a higher priority (compared to  $J_2$ ) job at  $S_2$ . Such flexibility is not viable with priority ordering as, for instance, one job is the lowest priority job at all stages. Conclusively,

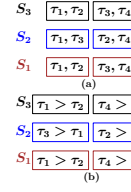


Fig. 2: (a). Job-to-resource mapping. (b) Pairwise priority assignment.

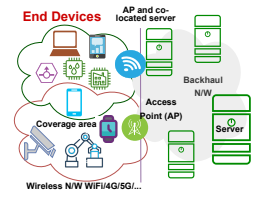


Fig. 3: Edge Computing System

with Observation V.1, it is established that the optimality of OPDCA ceases to exist with respect to problem P2. On the other hand, if a feasible total priority order exists, a pairwise priority assignment is also ensured.

As  $S^{DCA}$  only requires the set of higher (and lower) priority jobs, it can also be used for computing a pairwise priority assignment. Furthermore, Eq. 4 can be used to compute a pairwise priority assignment as OPA is not being used (hence, OPA-compatibility is not essential) to assign priority.

### A. ILP Formulation to Compute a Pairwise Priority Assignment

In the following, we present an ILP formulation, hereafter referred to as OPT, to compute pairwise priority assignment for preemptive scheduling (the same can also be extended for non-preemptive scheduling). For every job  $J_k \in \mathcal{M}_i$ , let  $X_{i,k}$  be a binary variable such that  $X_{i,k} = 1$  if  $J_i > J_k$ , otherwise  $X_{i,k} = 0$ . Thus,

$$X_{i,k} + X_{k,i} = 1 \quad (7)$$

Moreover, the end-to-end delay of  $J_i$  under preemptive scheduling (Eq. 6) is

$$\Delta_i = t_{i,1} + \sum_{J_k \in \mathcal{M}_i} X_{k,i} \sum_{x=1}^{w_{k,i}} \tilde{t}_{k,x} + \sum_{j=1}^{N-1} \theta_{i,j} \quad (8)$$

where  $\theta_{i,j} = \max_{J_k \in \mathcal{Q}_i} \{\tilde{p}_{k,j}\}$  and computed as described in the following. Let  $P = \max_{i,j} \{P_{i,j}\}$  and  $\mathcal{Z}_{i,j} = \mathcal{M}_{i,j} \cup J_i$ . Define  $|\mathcal{Z}_{i,j}|$  ( $|\cdot|$  denotes the cardinality of a set) auxiliary binary variables  $\{b_y | y \in [1, |\mathcal{Z}_{i,j}|]\}$ . Let  $z_{i,j,y} \in \mathcal{Z}_{i,j}$  denote the  $y^{th}$  job in  $\mathcal{Z}_{i,j}$ .

$$\theta_{i,j} \geq X_{k,i} \tilde{p}_{k,j}, \forall J_k = z_{i,j,y} | z_{i,j,y} \in \mathcal{Z}_{i,j} \quad (9a)$$

$$\theta_{i,j} \leq X_{k,i} \tilde{p}_{k,j} + (1 - b_y)P, \forall J_k = z_{i,j,y} | z_{i,j,y} \in \mathcal{Z}_{i,j} \quad (9b)$$

$$\sum_y b_y = 1 \quad (9c)$$

Then, the problem of pairwise priority assignment is to determine  $X_{i,k}$  subject to constraints (i) specified in Eqs. 7 and 9, and (ii) timing constraints ( $\Delta_i \leq D_i | \forall J_i \in \mathcal{J}$ ) are satisfied.

### B. Deadline-Monotonic & Repair (DMR) based Heuristic

DMR computes a feasible pairwise priority assignment. DMR begins with a Deadline-Monotonic (DM) based pairwise priority assignment<sup>9</sup> followed by a *repair* phase in case of

<sup>9</sup>DM is not optimal even in a multi-stage single-resource system for preemptive jobs with the same arrival time. The same can be verified by fixing  $D_1 = 60$  in Example 1.  $J_1$  gets the lowest priority, and thus,  $\Delta_1 = 82$ .

---

**Algorithm 2: DMR**

---

```
1 for each job pair  $\langle J_i, J_k \rangle | i \in [1, n-1], k \in [i+1, n],$   
    $J_k \in \mathcal{M}_i$  do  
2    $\lfloor$  Set  $J_i > J_k$  if  $D_i \leq D_k, J_k > J_i$  otherwise  
3 Compute  $\Delta_i, \forall J_i \in \mathcal{J}$   
4 for each job  $J_i | \Delta_i > D_i$  do  
5   Compute  $\mathcal{F}_i \subseteq \mathcal{M}_i$  such that  $J_k \in \mathcal{F}_i$  provided  
    $J_k > J_i$  &  $\Delta_k < D_k$   
6   Sort jobs of  $\mathcal{F}_i$  in the decreasing order of  $D_k - \Delta_k$   
7   for each job  $J_k \in \mathcal{F}_i$  do  
8     Set  $J_i > J_k$  if feasible for  $J_k$   
9     Recompute  $\Delta_i$ . If  $\Delta_i \leq D_i$  go to Step 4  
10   $\lfloor$  If  $\Delta_i > D_i$ , declare  $\mathcal{J}$  infeasible
```

---

timing constraint violation. In the repair phase, for a job pair  $\langle J_i, J_k \rangle | J_k \in \mathcal{M}_i, \Delta_i > D_i, \Delta_k < D_k, J_k > J_i$ , the priority is reversed to  $J_i > J_k$  provided it is feasible for  $J_k$ . Pseudo-code for DMR is listed in Algorithm 2.

## VI. EXPERIMENTAL EVALUATION

The performance of the proposed approaches<sup>10</sup> is evaluated for holistic scheduling real-time jobs in edge computing systems, a typical example of MSMR systems.

### A. Simulation Setup

*System Architecture:* We consider an edge computing system as shown in Figure 3. Deadline-constrained jobs generated by *End Devices* (EDs) are offloaded to an *access point* (AP) through heterogeneous wireless networks using protocols such as 5G-URLLC (Ultra-Reliable Low Latency Communication). APs, in turn, forward the workloads to one of the edge/cloud servers through a wired backhaul network. Subsequently, the computed result is returned to the corresponding ED via an AP. Considering a wired high-speed and contention-free backhaul, we assume that the deadline has already been adjusted to incorporate the backhaul delay. Thus, the number of stages is  $N = 3$  (wireless network for offloading, compute resource at server, and wireless network for downloading). For a given set of jobs and corresponding job-to-resource mapping<sup>11</sup>, the problem is to assign fixed priority to jobs for processing at servers and APs so that timing constraints are satisfied. Furthermore, considering a generic model, preemption is allowed at a server but prohibited during offloading (and downloading) at APs.

In the edge/cloud architecture, the schedule is computed at periodic intervals, and all previously arrived and unscheduled jobs are simultaneously considered. Thus, a lower-priority job begins offloading only after the higher-priority jobs; hence,  $\mathcal{H}_i^a = \emptyset$ . However, a job can be delayed by a lower-priority job

during downloading. Thus, Eq. 6 can be extended by including the delay due to non-preemptive downloading as follows.

$$\Delta_i \leq \sum_{J_k \in \mathcal{Q}_i} \sum_{x=1}^{w_{i,k}} \tilde{t}_{k,x} + \max_{J_k \in \mathcal{Q}_i} \tilde{p}_{k,1} + \max_{J_k \in \mathcal{Q}_i} \tilde{p}_{k,2} + \max_{J_k \in \mathcal{L}_i} \tilde{p}_{k,3} \quad (10)$$

Being a 3-stage system with non-preemption only at the last stage, Eq. 10 is OPA-compatible even if the 4<sup>th</sup> term (delay component due to non-preemption) is computed over  $\mathcal{L}_i$ .

*System Parameters:* The system parameters are selected in line with existing works [11], [12]. The number of APs (for offloading and downloading), servers, and jobs are fixed to 25, 20, and 100, respectively. Moreover, the uplink and downlink speeds of APs, computation speed of servers, and job parameters are set such that offloading, processing, and downloading times (in milliseconds) of a job are in the ranges [2, 200], [50, 500], and [2, 100], respectively.

We define heaviness of a job  $J_i$  at stage  $S_j$  as  $h_{i,j} = \frac{P_{i,j}}{D_i}$ . Let  $\chi_{y,j}$  be the sum of the heaviness of all jobs mapped to the  $y^{\text{th}}$  resource at  $S_j$ . Then, we define the heaviness of the set of jobs as  $H = \max_{y,j} \{\chi_{y,j}\}$ . Note that  $H$ , with the number of stages, resembles total utilization. Then, the performance of the proposed approaches is reported for varying: (i) *heaviness threshold* ( $\beta$ ):  $J_i$  is referred to as *heavy* at  $S_j$  provided  $h_{i,j} \geq \beta$ , (ii) *per-stage heaviness requirement* ( $[h_1, h_2, h_3]$ ):  $h_j$  denotes the ratio of jobs that are heavy at  $S_j$ , and (iii) *heaviness bound* ( $\gamma$ ) such that  $H \leq \gamma$ . The maximum heaviness of a job at any stage is always bounded by  $2\beta$ . Furthermore, the default values are  $\beta = 0.15$ ,  $h_1 = h_2 = 0.05$ ,  $h_3 = 0.01$ , and  $\gamma = 0.7$ .

To investigate the usefulness of the repair phase in DMR, we use a baseline, referred to as DM, which is like DMR but without any repair phase. A decomposition-based approach, referred to as DCMP, is used as another baseline, which involves decomposing the end-to-end deadline of each job to a virtual deadline corresponding to each stage. Let  $\Upsilon_{i,j}$  denotes the sum of heaviness of all jobs mapped to  $R_{i,j}$  (the resource to which  $J_i$  is mapped in stage  $S_j$ ); then, the virtual deadline of  $J_i$  at  $S_j$  is  $D_i \times \frac{\Upsilon_{i,j}}{\sum_j \Upsilon_{i,j}}$ . As a schedulability test that can be applied even to the decomposed jobs in this setting does not exist, we simulate the processing of decomposed jobs by assigning priorities in the inverse order of the deadline.

### B. Results and Discussion

We define *acceptance ratio* (AR) of an approach as the percentage of accepted test cases. Moreover,  $\text{AR}_{DM}$ , for instance, denotes the AR of DM. As shown in Figure 4, for a lightly loaded system such as  $\beta = 0.05$  (i.e.,  $J_i$  with  $P_{i,j} \geq 0.05 \times D_i$  is heavy at  $S_j$ ),  $\text{AR}_{DM}$  is substantially high ( $> 97\%$ ). However, as  $\beta$  increases, the processing time of all jobs, heavy as well as those which are not heavy, increases; the same can be observed with the rapidly deteriorating  $\text{AR}_{DM}$ . However,  $\text{AR}_{OPDCA}$  (feasibility of computing a priority ordering) remains relatively stable except for heavily loaded systems. For instance, for  $\beta = 0.15$ , a priority ordering can be computed (by OPDCA) for 94% of the test cases, whereas 71% and 74% test cases are schedulable by DM and DCMP, respectively.  $\text{AR}_{OPDCA}$

<sup>10</sup>Gurobi has been used to solve OPT. All source codes are available at <https://github.com/CPS-research-group/CPS-NTU-Public/tree/DATE2024>

<sup>11</sup>Due to the highly complex nature of holistic resource allocation and scheduling problems, the two problems are commonly addressed separately. Many works, such as [11], [12], focus on the former. Surveys including [13] may be referred for related works.



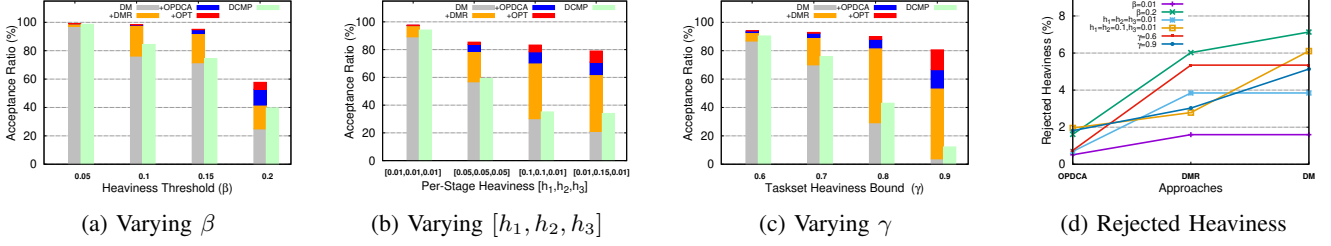


Fig. 4: Acceptance ratios (ARs) for varying (a)  $\beta$ , (b)  $[h_1, h_2, h_3]$ , and (c)  $\gamma$ . The base of the stacked histogram represents the AR of DM. The increment in the AR of DMR, OPDCA, and OPT compared to the AR of DM, DMR, and OPDCA, respectively, are stacked (upward) in this sequence. (d) Performance as admission controller.

is only marginally higher than  $AR_{DM}$  or  $AR_{DCMP}$  for lightly loaded system such as  $\beta = 0.05$  (or  $\gamma = 0.6$  or  $h_j = 0.01$ ). However, a substantial gap is observed with increasing load.

Likewise, for a lightly loaded system, AR for pairwise priority assignment using DMR is comparable to that of the DM. However,  $AR_{DM}$  decreases rapidly with increasing load, while a relatively stable  $AR_{DMR}$  indicates the usefulness of the repair phase. We observe a higher AR for computing a priority ordering (by OPDCA, an optimal algorithm) compared to computing a pairwise priority assignment (by DMR, a heuristic). However,  $AR_{OPT}$  (to compute an optimal pairwise priority assignment) is higher than  $AR_{OPDCA}$ , supporting Observation V.1. For a moderate heaviness, such as  $\gamma = 0.8$ , a pairwise priority assignment can be computed for 90% of the test cases by OPT compared to the priority ordering for 88% by OPDCA. However, the gap in the performance of OPDCA and OPT becomes more evident with higher load. Moreover, the simulation results also reveal the rapidly degrading performance of DCMP with increasing load.

We also observe the performance of the proposed approaches by running them as admission controllers to accept as many jobs as possible. Thus, in Step 10 of Algorithm 1, rather than declaring the entire set of jobs infeasible, a job (with the largest  $\Delta_i - D_i$ ) is discarded, and priority assignment is reattempted for the jobs which are yet to be assigned a priority. Likewise, in Step 10 of Algorithm 2, the job with the largest  $\Delta_i - D_i$  is discarded. Figure 4d illustrates the *rejected heaviness*, defined as the percentage of the heaviness of rejected jobs with respect to the heaviness of all jobs. In general, even for the parameters corresponding to high system load, OPDCA performs reasonably well.

Conclusively, we observe higher acceptance ratios for OPT and OPDCA compared to other approaches. However, the performance gap increases with the system load, indicating the suitability of pairwise priority assignment. Further, this gap is likely to grow with the number of stages, resources, and jobs.

## VII. CONCLUSION AND FUTURE WORKS

In this work, the delay composition rule proposed in [6], [7] is refined and shown to be OPA-compatible for optimal FP scheduling of real-time jobs with end-to-end delay in a multi-stage multi-resource (MSMR) distributed system. Furthermore, the suitability of pairwise priority assignment is also established in MSMR systems. Moreover, an ILP formulation

and a deadline-monotonic repair-based heuristic have also been presented to compute a pairwise priority assignment. Finally, simulation results are discussed for the holistic scheduling of real-time jobs for compute and network resources in an edge-computing system using the proposed approaches.

This work presents a starting point for future research to apply and investigate the use of DCA for scheduling, in addition to the originally intended schedulability analysis, in distributed systems. Future research directions also include further exploration of pairwise priority assignment strategies, including developing algorithms with performance guarantees.

## REFERENCES

- [1] R. J. Bril, S. Altmeyer, M. M. van den Heuvel, R. I. Davis, and M. Behnam, "Fixed priority scheduling with pre-emption thresholds and cache-related pre-emption delays: integrated analysis and evaluation," *Real-Time Systems*, vol. 53, pp. 403–466, 2017.
- [2] R. I. Davis, L. Cucu-Grosjean, M. Bertogna, and A. Burns, "A review of priority assignment in real-time systems," *Journal of systems architecture*, vol. 65, pp. 64–82, 2016.
- [3] N. C. Audsley, "Optimal priority assignment and feasibility of static priority tasks with arbitrary start times," 1991.
- [4] J. C. Palencia and M. G. Harbour, "Offset-based response time analysis of distributed systems scheduled under edf," in *15th Euromicro Conference on Real-Time Systems, 2003. Proceedings.* IEEE, 2003, pp. 3–12.
- [5] S. Hong, T. Chantem, and X. S. Hu, "Meeting end-to-end deadlines through distributed local deadline assignments," in *2011 IEEE 32nd Real-Time Systems Symposium.* IEEE, 2011, pp. 183–192.
- [6] P. Jayachandran and T. Abdelzaher, "Delay composition in preemptive and non-preemptive real-time pipelines," *Real-Time Systems*, vol. 40, no. 3, pp. 290–320, 2008.
- [7] —, "Transforming distributed acyclic systems into equivalent uniprocessors under preemptive and non-preemptive scheduling," in *2008 Euromicro Conference on Real-Time Systems.* IEEE, 2008, pp. 233–242.
- [8] R. Bettati and J. W.-S. Liu, "End-to-end scheduling to meet deadlines in distributed systems," in *ICDCS.* Citeseer, 1992, pp. 452–459.
- [9] R. I. Davis and A. Burns, "Priority assignment for global fixed priority pre-emptive scheduling in multiprocessor real-time systems," in *2009 30th IEEE Real-Time Systems Symposium.* IEEE, 2009, pp. 398–409.
- [10] N. C. Audsley, "On priority assignment in fixed priority scheduling," *Information Processing Letters*, vol. 79, no. 1, pp. 39–44, 2001.
- [11] C. Gao, A. Shaan, and A. Easwaran, "Deadline-constrained multi-resource task mapping and allocation for edge-cloud systems," in *GLOBECOM 2022-2022 IEEE Global Communications Conference.* IEEE, 2022, pp. 5037–5043.
- [12] W. Chu, P. Yu, Z. Yu, J. C. Lui, and Y. Lin, "Online optimal service selection, resource allocation and task offloading for multi-access edge computing: A utility-based approach," *IEEE Transactions on Mobile Computing*, 2022.
- [13] S. Ramanathan, N. Shivaraman, S. Suryasekaran, A. Easwaran, E. Borde, and S. Steinhorst, "A survey on time-sensitive resource allocation in the cloud continuum," *it-Information Technology*, vol. 62, no. 5-6, pp. 241–255, 2020.