

# On Schedulability and Time Composability of Data Aggregation Networks

Fatemeh Saremi\*, Praveen Jayachandran<sup>†</sup>, Forrest Iandola\*, Md Yusuf Sarwar Uddin\*, Tarek Abdelzaher\*, and Aylin Yener<sup>‡</sup>

\*Department of Computer Science, University of Illinois, Urbana, IL 61801

<sup>†</sup>IBM Research, India

<sup>‡</sup>Department of Electrical Engineering, Pennsylvania State University, University Park, PA 16802

Email: saremi1@illinois.edu, prjayach@in.ibm.com, {iandola1, mduddin2, zaher}@illinois.edu, yener@ee.psu.edu

**Abstract**—This paper develops a framework to analyze the latency and delay composition of workflows in a real-time networked aggregation system. These workflows are characterized by different sensor inputs that are processed along parallel branches that eventually merge or fuse to compute the aggregation result. The results for each flow must be produced within certain end-to-end deadlines or else the information would become stale, inaccurate and useless. We extend results developed by the authors recently and consider an end-to-end view of the aggregation system that allows us to derive a much tighter analysis of the end-to-end delay compared to traditional analysis techniques. We then provide a reduction of the aggregation network system to an equivalent hypothetical uniprocessor for the purposes of schedulability analysis. Extensive simulations show that latency bound obtained from the analysis framework is significantly more accurate than that of traditional analysis techniques.

## I. INTRODUCTION

Multisensor data aggregation is a discipline concerned with collecting and processing data from multiple sources in order to produce information that is accurate and relevant. An overarching challenge in such systems, in which the data collected is oftentimes real-time in nature, is to guarantee that the end-to-end aggregation process can be completed within certain time constraints. For instance, in defense surveillance systems, detection, tracking and identification of the threat have to complete within stringent deadline constraints in order to enable timely action to counter the threat.

The workflows we consider in this paper are characterized by different sensory inputs processed along parallel branches that merge together to produce the aggregation results. When parallel branches merge together, the aggregation results can be computed only when *all* incoming branches have completed execution. A workflow could have multiple such successive merges. The aggregation results of individual jobs of workflows need to be computed within certain end-to-end deadlines. This system model presents a deviation from existing literature and presents a new “MERGE” primitive that is characteristic of aggregation systems. Composition-based analysis of such aggregation workflows is a challenging problem, and we will explain why seemingly-intuitive approaches do not yield correct results.

In order to guarantee that time constraints will be met while resources not being severely underutilized, an accurate

analysis technique with little pessimism is required to estimate latency resulting from a particular allocation of resources to sensor data processing streams. In this paper, we investigate timing composability and schedulability of workflows in such multisensor data aggregation systems. We develop an analysis framework which extends recent work by the authors in developing a delay composition theory [1]. The theory was designed for jobs that traverse a sequence of resource stages in a distributed system and is not applicable to aggregation systems characterized by the “MERGE” primitive described above.

Our framework develops theoretical rules bounding latency in data aggregation systems under non-preemptive and preemptive scheduling. We consider an end-to-end view of the system and this allows us to develop intuition into the main system and task parameters that contribute towards increasing the end-to-end delay, and to express our delay bound in terms of these parameters. This results in a much tighter analysis of the latency and schedulability of jobs compared to existing analysis techniques, especially for large systems.

We evaluate our theoretical framework through extensive simulations based on two related metrics. The first is the capacity or the maximum workload due to jobs of workflows that can be guaranteed to meet their end-to-end deadlines. This is important in ensuring that the system operates at its capacity limit without wasting resources serving jobs that do not meet their end-to-end deadlines. The second metric measures the ratio of the average end-to-end delay of jobs of a workflow to the analytically computed worst-case end-to-end delay bound. This gives us a measure of the tightness of the analysis.

The rest of this paper is organized as follows. Section II reviews related work and section III describes the system model and the problem addressed in this paper. In section IV, we explain why seemingly intuitive approaches to calculate the worst-case latency are not accurate for aggregation systems. We then proceed to formally derive the delay bound under non-preemptive and preemptive scheduling in section V. Section VI presents reduction rules that can be employed to transform the entire data aggregation system into a single hypothetical stage for the purposes of schedulability analysis. Then the performance of the analysis framework is extensively evaluated using simulations in section VII. Section VIII concludes the paper.

## II. RELATED WORK

Existing data aggregation schedulability literature often focuses on a specialized system model and a specific scheduling algorithm. For example, Li and Cao proposed a non-preemptive scheduling algorithm for data aggregation systems called Coordinated Workload Scheduling (CWS) [2]. They also analyze schedulability, but the schedulability results in that study only apply when a particular scheduling algorithm is used. CWS constrains the system model to contain just three workflows: sensing, communication, and computation. In the same vein, in an earlier paper we analyzed schedulability of real-time flows following a fluid task model and under non-preemptive EDF scheduling [3].

There also exist several studies that are applicable to more general data aggregation system models. Several algorithms have been developed to perform optimal offline scheduling of workflow sets in distributed systems [4] [5]. These algorithms work by constructing the complete schedule for all the jobs executing on all the stages. The schedule is then used to determine the schedulability of the workflow set. However, these optimal algorithms require NP-hard computation.

To reduce the computational complexity of schedulability testing, Kao et al. [6] and Zhang et al. [7] presented techniques to divide the end-to-end deadline into per-stage deadlines. Uniprocessor schedulability tests are then used to determine whether each stage is schedulable. If all the stages are schedulable, the system is deemed to be schedulable. These techniques allow for a generalized system model and are not constrained to a particular scheduling algorithm. However, they do not accurately account for the inherent parallelism in the execution of different stages of a data aggregation system. As a result, they tend to exhibit extreme pessimism when estimating schedulability boundaries for pipelined and merging workflows, especially for large systems.

Holistic analysis [8] [9] and real-time calculus [10] [11] comprise a “middle ground” between pessimistic bounds and NP-hard optimal calculations. Unlike some of the related work we have discussed, holistic analysis and network calculus allow for a general system model. A study by Koubaa and Song found that holistic analysis is less pessimistic than network calculus for most system configurations [12]. For this reason, we use holistic analysis as a benchmark for our evaluation. Yet, as the holistic approach analyzes one stage at a time, the pessimism of the analysis grows with system scale. In contrast, in this paper we develop a delay bound for aggregation systems by considering an end-to-end view of the delay. This allows us to develop a bound that does not become increasingly pessimistic with system scale. Further, by reducing the problem of schedulability analysis of the distributed system to that on an equivalent uniprocessor, we significantly reduce the complexity of analysis, making it extremely suitable for large systems.

## III. AGGREGATION MODEL AND PROBLEM STATEMENT

In this paper, we consider an abstract model of a data aggregation system, comprising of  $m$  data workflows (e.g.,

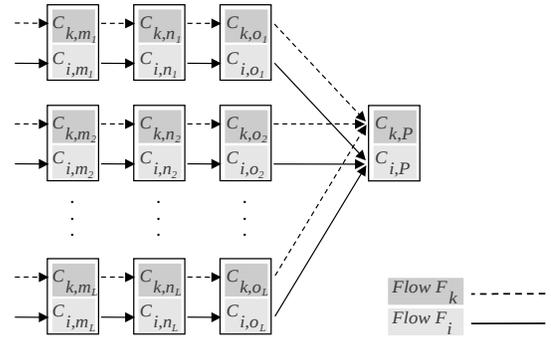


Figure 1: An example of merging workflows.

audio and video data from camera and microphone sensors, speed and proximity sensor information in a next-generation automobile), each requiring several stages of processing (such as monitoring, threat analysis, actuation and display). Each processing stage is handled by a single resource that is scheduled in priority order. For instance, a resource could refer to processors that schedule and serve arriving tasks, or a network link on which enqueued packets are transmitted. Since each stage represents one resource, in the rest of this paper, we use the terms *stage* and *resource* interchangeably. Each workflow could potentially have multiple branches that ultimately merge forming an aggregation tree. Collectively, the stages of processing of all workflows taken together are organized in a graph called the *workflow graph*. It determines precedence constraints among the different stages. For simplicity, we only consider workflow graphs that are trees (no two task paths can split and re-merge with one another).

Consider a flow  $F_i$  that has  $L$  parallel pipelines, described by disjoint paths  $p_1, \dots, p_L$ , followed by a common parent stage,  $P$ , as shown in Figure 1. Jobs from different workflows are assumed to be aperiodic and may have different arrival offsets. Let  $Off_i$  denote the offset from time zero at which a job  $J_i$  of workflow  $F_i$  arrives at all the branches of its workflow (if different branches arrive at different offsets,  $Off_i$  can be set to the maximum of these offsets). A job does not become eligible to execute on the merge-stage (the common parent,  $P$ ) until all pipelines have finished processing it. A workflow could potentially have several such merge-segments.

In a *real-time* data aggregation system, each workflow,  $F_i$ , has an end-to-end latency constraint,  $D_i$ , denoting the maximum allowable latency between the arrival of a new job of workflow  $F_i$  into the system (to all its branches), and the completion of its processing on the last stage. We call this constraint the *end-to-end deadline* of workflow  $F_i$ .

Each job (invocation) of every workflow is assigned a priority and we derive results for both non-preemptive and preemptive scheduling. We also assume that the relative priority of each job remains the same across all the stages on which it executes. Let us assign priorities such that the lowest number corresponds with the highest priority. Thus, we let “ $i \leq k$ ” be shorthand for “ $J_i$  is at least as high-priority as  $J_k$ .” The main problem we address in this paper is to determine a worst-case bound on the end-to-end delay of jobs for an arbitrary workflow given the computation times of other workflows that

exist concurrently with it in the system.

#### IV. INTUITION

In this section, we develop some intuition that will assist us in deriving the delay bound of workflows in an aggregation network. We explore a simple and seemingly-intuitive approach to estimate the worst-case delay using an example and show why the approach is not accurate. This exercise points us to the right intuition into developing the worst-case delay bound in the next section.

Consider the following system. Let  $j_{merge}$  denote a merge stage in the workflow graph of job  $J_k$  which has  $l$  incoming branches from stages  $j^{m_1}, j^{m_2}, \dots$ , and  $j^{m_l}$ . By semantics of merge,  $J_k$  can execute on  $j_{merge}$  only when all its subjobs arrive from the  $l$  upstream stages. So,  $J_k$  can execute on stage  $j_{merge}$  after the last of the  $l$  sub-jobs arrives. One way to compute the worst-case end-to-end delay of  $J_k$ , is to consider all the possible paths from aggregation input to computing the aggregation result, and picking the path that incurs the largest end-to-end delay bound obtained through the delay composition theory. In the following, we explain why this intuition does not work.

Such a system with two merging branches is illustrated in Figure 2. Let us assume that scheduling is non-preemptive in nature. Job  $J_l$  is the job whose end-to-end delay we wish to bound. Job  $J_h$  is a higher priority job, and jobs  $J_{l_1}, J_{l_2}, J_{l_3}$  and  $J_{l_4}$  are lower priority jobs. Let us further suppose that  $J_l$  arrives at stages  $S_1$  and  $S_3$  concurrently at time  $\epsilon > 0$  ( $\epsilon$  can be arbitrarily small). Job  $J_h$  arrives at stages  $S_1$  and  $S_3$  at time 30 units. Job  $J_{l_1}$  executes only on stages  $S_1$  and  $S_3$  and arrives at time zero. Job  $J_{l_2}$  executes only on stage  $S_2$  and arrives at time  $20 - \epsilon$ . Job  $J_{l_3}$  executes only on stage  $S_4$  and arrives at time  $38 - \epsilon$ . Job  $J_{l_4}$  arrives at stage  $S_5$  at time  $58 - \epsilon$ . The computation times and execution trace of all the jobs are shown in Figure 2.

Now, let us use the delay composition theory developed in [1] along each path in the workflow graph of  $J_l$ , and

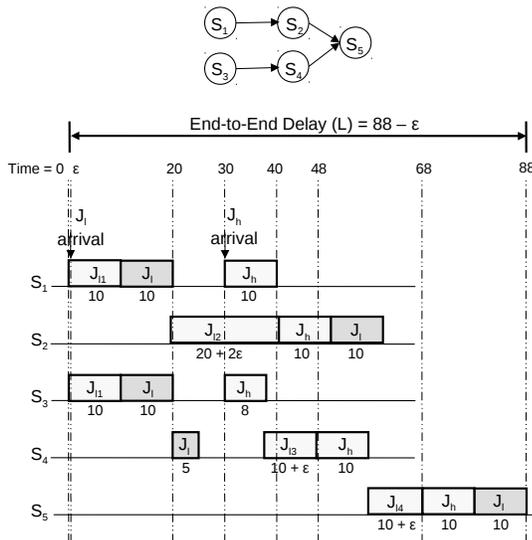


Figure 2: Workflow graph and execution trace for seemingly candidate aggregation composition approaches.

compute the maximum to obtain an estimate of the worst-case end-to-end delay bound. For the benefit of the reader, we state the delay composition result under non-preemptive scheduling from [1] here:

**Non-preemptive Delay Composition Theorem [1].** Assuming a non-preemptive scheduling policy with the same priorities across all stages for each job, the end-to-end delay of a job  $J_k$  of  $N$  stages can be composed from the execution parameters of other jobs that delay it as follows:

$$Delay(J_k) \leq \sum_{J_i \in H} C_{i,max} + \sum_{j \leq N} (\max_{J_i \in H} C_{i,j} + \max_{J_i \in L_j} C_{i,j}) \quad (1)$$

where  $H$  is the set of higher priority jobs and  $L_j$  is the set of lower priority jobs that execute on stage  $j$ .

The worst-case delay bound along the pipeline segment  $\langle S_1, S_2, S_5 \rangle$  can be computed using delay composition theory (ignoring the parallel branch along  $\langle S_3, S_4 \rangle$ ), as the sum of three terms. The first term is the sum of maximum computation times for each higher priority job across all the stages on which it executes. The second term, is the stage additive component, which is one maximum stage execution time over all higher priority jobs for each stage on the path. The third term is the blocking term, which is the sum of the maximum execution time of any lower priority job at each stage.

$$\begin{aligned} Delay_{\langle S_1, S_2, S_5 \rangle}(J_l) &\leq 10 + (3 \times 10) \\ &\quad + (10 + 20 + 2\epsilon + 10 + \epsilon) \\ &= 80 + 3\epsilon \end{aligned} \quad (2)$$

Similarly, the worst-case delay along the pipeline segment  $\langle S_3, S_4, S_5 \rangle$  can be computed as follows:

$$\begin{aligned} Delay_{\langle S_3, S_4, S_5 \rangle}(J_l) &= 10 + (3 \times 10) \\ &\quad + (10 + 10 + \epsilon + 10 + \epsilon) \\ &= 70 + 2\epsilon \end{aligned} \quad (3)$$

The maximum delay across both paths for  $J_l$  is therefore  $80 + 3\epsilon$ . However, as the execution trace suggests, the delay of  $J_l$  can be as large as 88 time units.

$$Delay_{trace}(J_l) = 88 - \epsilon$$

So, why didn't our intuition work in this case? Delay composition theory implicitly assumes that once a higher priority job preempts or overtakes a lower priority job, it will always execute ahead of the lower priority job on all future stages. However, when there are "merge" stages, this assumption breaks down. As illustrated in the example, while  $J_h$  delays  $J_l$  on one branch ( $\langle S_1, S_2 \rangle$ ), it executes after  $J_l$  in a different busy period on a parallel branch ( $\langle S_3, S_4 \rangle$ ) and takes longer to arrive at the merge stage. We call this a *revisit* event and a more formal definition is provided in section V. This is why current delay composition theory is not applicable to this case. We need to account for this revisit event.

A closer look at the above example reveals that the revisit event arises due to the fact that the higher priority job  $J_h$  completes execution sooner on one branch, while the lower priority job  $J_l$  completes sooner on another. This in turn requires that  $J_l$  arrives before  $J_h$  to the system. A closer scrutiny revealed that our previous delay composition result

works fine as long as  $J_l$  arrives no earlier than  $J_h$ . So, we should be able to calculate the delay starting from the arrival time of  $J_h$  using our previous result, and add the difference in the arrival times between  $J_l$  and  $J_h$ . This is how we derive our delay bound in section V. The end-to-end delay bound is expressed as the sum of two terms. The first is the maximum offset between the arrivals of  $J_l$  and any higher priority job. The second term is the maximum delay along any end-to-end path that constitutes the aggregation graph of  $J_l$ .

A natural question is, how do we determine the maximum offset of any higher priority job that can potentially interfere with the job under consideration? A simple, but pessimistic, solution is as follows. To begin with, assume that all higher priority jobs in the system can delay the job under consideration  $J_1$ . Use the analysis to determine the worst-case delay bound. If the arrival time of any higher priority job is greater than the worst-case delay bound, then that higher priority job cannot possibly delay  $J_1$ . Such higher priority jobs can be removed and the analysis can be repeated (resulting in a lower worst-case bound). At each iteration, we can discard some higher priority jobs, and the process is repeated until no higher priority jobs can be discarded from the interfering set of jobs. Clearly, this approach will only overestimate the set of interfering higher priority jobs. The analysis stands to gain if the worst-case arrival offset information for higher priority jobs is provided as system input together with the execution time characteristics of jobs.

## V. AN END-TO-END DELAY BOUND FOR AGGREGATION WORKFLOWS

In this section, we derive a worst-case bound on the end-to-end delay of an aggregation workflow in terms of the stage computation times of other workflows executing concurrently with it. We derive the result for non-preemptive scheduling in section V-A, and in the interest of brevity, only state the result for preemptive scheduling in section V-B. The analysis extends previous work by the authors in developing a theory called delay composition theory to analyze the worst-case end-to-end delay of jobs in various workflow topologies [1], [13]. The theory was developed for jobs which traverse a sequence of resource stages and so is not applicable to aggregation systems characterized by the ‘‘MERGE’’ primitive. This paper considers system graphs that contain aggregation nodes, with the property that execution on these nodes require the execution on *all* incoming branches to complete before it can be undertaken. As motivated in the previous section, this important distinction makes such systems difficult to analyze.

### A. Non-Preemptive Scheduling

As priority is assigned for each job, let us order all the jobs in decreasing priority order. Let us suppose that we wish to bound the worst-case end-to-end delay of a job  $J_k$ . Let  $S$  denote the set of all jobs executing concurrently with  $J_k$ . Let  $Off_i$  denote the offset from time zero at which a job  $J_i$  arrives at all the branches of its workflow (if different branches arrive at different offsets,  $Off_i$  can be set to the maximum of these

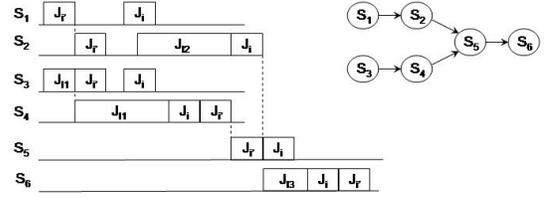


Figure 3: Example illustrating a revisit event

offsets). Let  $Delay_i$  denote the worst-case end-to-end delay of job  $J_i$  (we wish to determine  $Delay_k$ ). Further, let  $Delay_{i,j}$  denote the worst-case delay of a job  $J_i$  up to completing its execution on stage  $j$ . Let  $Paths_i$  denote the set of all paths from any source of a branch of job  $J_i$  to its sink (the union of all paths in  $Paths_i$  is the workflow tree of  $J_i$ ). We now state the main result of this paper:

**Non-Preemptive Delay Composition Theorem for Aggregation Workflows:** *Under a non-preemptive scheduling policy that assigns the same priority across all stages for each job, the worst-case end-to-end delay of a job of workflow  $F_k$  in an aggregation tree is bounded as,*

$$Delay_k \leq \max_{i \leq k} (Off_i) - Off_k + \sum_{i \leq k} C_{i,max} + \max_{p \in Paths_k} \sum_{j \in p} (\max_{i \in S} C_{i,j} + \max_{i \leq k} C_{i,j}) \quad (4)$$

where  $Off_i$  is the offset of job  $J_i$  from time zero.

**Proof:** At each merge-stage in its workflow tree, as a job needs to wait for all incoming branches to complete execution before it can execute on the merge-stage, the following sequence of events is possible (which is otherwise not possible in the absence of merge-stages). A higher priority job  $J_i$  delays a lower priority job  $J_{i'}$  and arrives ahead of it on one branch of merge-stage  $j_1$ . However, along another branch,  $J_{i'}$  completes execution and arrives ahead of  $J_i$  to stage  $j_1$ . Due to this reversal in the arrival order, it is possible for  $J_i$  to again delay  $J_{i'}$  at a downstream stage  $j_2$  (can be the same as stage  $j_1$ ). Let us call the instance where  $J_i$  again executes ahead of  $J_{i'}$  after the reversal, as a *revisit event* (if  $J_i$  always executes after  $J_{i'}$  on all remaining stages, the revisit event occurs at the last execution stage of  $J_i$ ). This is illustrated in Figure 3 ( $J_{l1}$ ,  $J_{l2}$ , and  $J_{l3}$  are of lower priority than job  $J_{i'}$ ).

We shall now prove the theorem using induction on the number of such revisit events between jobs  $J_i$  and  $J_{i'}$ , such that  $i < i' \leq k$ . The basis step is when there are no such events. This means that at every merge-stage each higher priority job that was ahead of job  $J_k$  remains ahead of job  $J_k$  throughout the system. Hence,  $Delay_k$  can be bounded by the maximum delay across any end-to-end pipeline path from a source to the sink. Applying the delay composition result (Equation 1) to each path  $p \in Paths_k$  and picking the maximum, we can obtain a bound on  $Delay_k$  as follows:

$$Delay_k \leq \sum_{i \leq k} C_{i,max} + \max_{p \in Paths_k} \sum_{j \in p} (\max_{i \in S} C_{i,j} + \max_{i \leq k} C_{i,j})$$

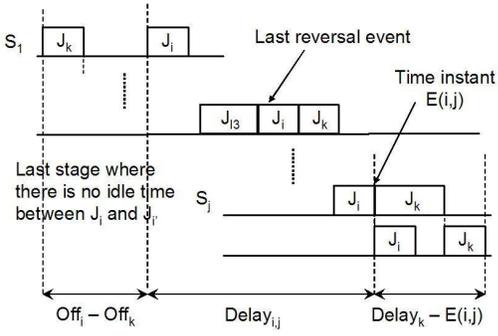


Figure 4: Example illustrating the proof

This proves the theorem for the basis step when there are no revisit events.

Let us assume that the result is true up to  $n-1$  revisit events,  $n \geq 1$ . We shall now prove that the result is true when there are  $n$  such events. Among all the  $n$  revisit events, consider the one that is latest in time. Let  $J_i$  be the higher priority job involved in this last revisit event. If  $J_i$  exits the system at some stage prior to the last stage on which  $J_k$  executes, let's add executions of zero computation time for  $J_i$  on the remaining stages. This operation cannot decrease the delay of  $J_k$ . Let stage  $j$  be the last stage where there is no idle time between the executions of  $J_i$  and  $J_k$ . Let  $E(i, j)$  denote the instant of time at which job  $J_i$  completes execution on stage  $j$ . We can now calculate the end-to-end delay of  $J_k$  as the sum of three terms: the offset of  $J_i$  relative to  $J_k$ , the delay of  $J_i$  up to its completion on stage  $j$ , and the delay of  $J_k$  from this instant on stage  $j$  to its sink. This is illustrated in Figure 4.

$$Delay_k \leq (Off_i - Off_k) + Delay_{i,j} + (Delay_k - E(i, j)) \quad (5)$$

The delay of job  $J_i$  up to stage  $j$ ,  $Delay_{i,j}$ , can be obtained from induction assumption as it has only  $n-1$  revisit events. Let  $H_1$  be the set of higher priority jobs including  $J_i$  that contribute to  $Delay_{i,j}$ .

$$Delay_{i,j} \leq \max_{i' \leq i} (Off_{i'} - Off_i) + \sum_{i' \in H_1} C_{i',max} + \max_{p \in Paths_k} \sum_{j' > j \in p} (\max_{i' \in S} C_{i',j'} + \max_{i' \leq i} C_{i',j'}) \quad (6)$$

In the above equation,  $j' \leq j$  denotes stages before stage  $j$  and  $j' > j$  denotes stages after stage  $j$  in path  $p$ . The delay of job  $J_k$  starting from instant  $E(i, j)$  on stage  $j$  does not encounter any revisit events and the delay can be obtained from the basis step. Let  $H_2$  be the set of higher priority jobs that contribute to the delay of  $J_k$  from stage  $j$  to the sink.

$$Delay_k - E(i, j) \leq \sum_{i' \in H_2} C_{i',max} + \max_{p \in Paths_k} \sum_{j' > j \in p} (\max_{i' \in S} C_{i',j'} + \max_{i' \leq k} C_{i',j'}) \quad (7)$$

Expanding Equation 5 using Equations 6 and 7 we get,

$$Delay_k \leq \max_{i' \leq k} (Off_{i'} - Off_k) + \sum_{i' \in H_1} C_{i',max} + \sum_{i' \in H_2} C_{i',max} + \max_{p \in Paths_k} \sum_{j' \in p} (\max_{i' \in S} C_{i',j'} + \max_{i' \leq k} C_{i',j'}) \quad (8)$$

Clearly,  $H_1 \cup H_2$  is a subset of the set of higher priority jobs of  $J_k$ . What remains to be shown is that  $H_1 \cap H_2 = \phi$ . Any higher priority job  $J_h$  in  $H_1$  has executed ahead of  $J_i$  on some stage  $j' \leq j$ . If  $J_h \in H_2$ , it must have arrived at stage  $j$  after  $J_i$  (there is an idle time between the executions of  $J_i$  and  $J_k$  after stage  $j$ ). This would be a revisit event that is after that of  $J_i$ 's, as  $J_h$  has executed ahead of  $J_i$  on a stage  $j'$  prior to  $j$ , leading to a contradiction. This proves the theorem.  $\square$

## B. Preemptive Scheduling

In the interest of brevity, we only state the result under preemptive scheduling. The proof is similar to the non-preemptive case.

**Preemptive Delay Composition Theorem for Aggregation Workflows:** Assuming a preemptive scheduling policy that assigns the same priority across all stages for each job, the worst-case end-to-end delay of a job of workflow  $F_k$  in an aggregation tree is bounded as,

$$Delay_k \leq \max_{i \leq k} (Off_i) - Off_k + \sum_{i \leq k} C_{i,max} + \max_{p \in Paths_k} \sum_{j \in p} \max_{i \leq k} C_{i,j} \quad (9)$$

where  $Off_i$  is the offset of job  $J_i$  from time zero.

## VI. SCHEDULABILITY ANALYSIS

In this section, we present a reduction of the aggregation workflow graph into an equivalent uniprocessor for the purposes of schedulability analysis. This reduction maintains the property that the worst-case delay of a job in the distributed aggregation graph is no more than the delay of the corresponding task in the equivalent uniprocessor. Thus, if the uniprocessor task completes execution within its deadline, then the corresponding distributed workflow will also complete execution within its end-to-end deadline. The hypothetical uniprocessor can be analyzed using any well-known schedulability analysis technique (e.g., response time analysis [14]).

Under non-preemptive scheduling, the schedulability of a job  $J_k$  of an aggregation workflow can be determined by analyzing the schedulability of a hypothetical uniprocessor constructed as follows:

- Each higher priority job  $J_i$  in the original distributed system is replaced with a job  $J_i^*$  on the hypothetical single stage, with a computation time equal to  $C_{i,max}$ , its maximum computation time across all stages in the distributed system and a deadline  $D_i$ , same as that of  $J_i$ .
- Job  $J_k$  is replaced with job  $J_k^*$  of lowest priority in the hypothetical single stage and computation time of  $\max_{i \leq k} (Off_i) - Off_k + C_{k,max} + \max_{p \in Paths_k} \sum_{j \in p} (\max_{i \in S} C_{i,j} + \max_{i \leq k} C_{i,j})$  and deadline equal to  $D_k$ .

Under preemptive scheduling, the hypothetical uniprocessor is constructed as follows:

- Each higher priority job  $J_i$  in the original distributed system is replaced with a job  $J_i^*$  on the hypothetical

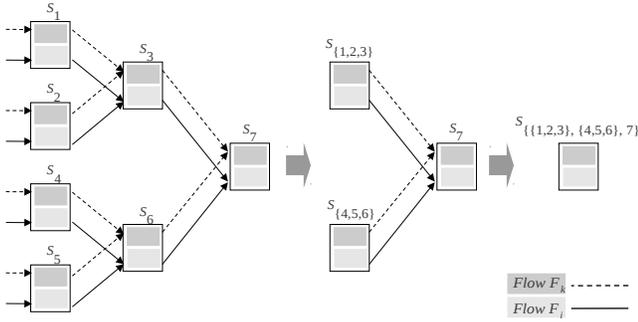


Figure 5: The reduction process on an aggregation tree in the delay composition framework

single stage, with a computation time equal to  $2C_{i,max}$  and deadline equal to  $D_i$ .

- Job  $J_k$  is replaced with job  $J_k^*$  of lowest priority in the hypothetical single stage and computation time of  $\max_{i \leq k}(Off_i) - Off_k + C_{k,max} + \max_{p \in Paths_k} \sum_{j \in p} \max_{i \leq k} C_{i,j}$  and deadline equal to  $D_k$ .

The reduction process in the delay composition framework is schematically depicted on an aggregation tree in Figure 5.

## VII. EVALUATION

In this section, we evaluate our approach to demonstrate the accuracy of our analysis in quantifying the end-to-end delays of workflows in an aggregation network. Specifically, we elaborate on how different system and load parameters can affect the effectiveness and accuracy of our approach. Towards this goal, two questions are explored: first, how efficiently our analysis utilizes the system resources, and second, how accurate it is in estimating the worst-case end-to-end delays.

We consider aggregation trees described as full binary trees for our simulations. So, a tree of height  $H$  has  $2^{H+1} - 1$  nodes and  $H$  levels of cascading merges. All workflows are assumed to execute on all nodes of the aggregation tree. Each workflow comprises of a sequence of aperiodic jobs. Each job is assigned an end-to-end deadline equal to  $500 \times H \times 10^\alpha$  simulation time units, where  $H$  denotes the height of the aggregation tree and  $\alpha$  is chosen uniformly across the interval  $[0, DR]$ . The default height of the aggregation tree is 5 and the number of stages is 63. The parameter  $DR$ , called deadline ratio parameter, states the order of variability in job deadlines and is set to 2.0 in all experiments, unless otherwise specified. This choice allows the end-to-end deadlines of jobs to vary by a factor of  $10^2$ . The computation times of jobs are chosen proportional to their deadlines based on a uniform distribution with a mean of  $\frac{D}{H \times JR}$  in which  $D$  denotes the corresponding end-to-end deadline and  $JR$  denotes the job resolution parameter (the degree of variability in jobs' computations times). Unless otherwise specified, the  $JR$  factor is set to 20. Different computation times are within 12% of the mean value. Jobs are activated with offsets randomly chosen from the interval  $[0, OR \times 500 \times H]$ , wherein  $OR$  denotes the offset resolution factor, set to 0.5 throughout the experiments (note that  $500 \times H$  is the minimum end-to-end deadline of any job). Job priorities are assigned based on their end-to-end deadlines.

Our new delay composition framework and holistic analysis [8] are simulated and compared in terms of different system and workload parameters such as the number of stages in the system, job resolution, deadline ratio, offset resolution, and the number of real-time jobs. Every experiment is repeated for 50 times and the 95<sup>th</sup> percentile confidence level is within 1% of the plotted average value. For the sake of readability, we do not plot the confidence values.

### A. Resource Utilization

The delay composition framework can be employed as an admission controller, which admits only the part of the workload that it deems to be feasible (all admitted workflows are guaranteed to meet their end-to-end deadline). Guaranteeing the schedulability of all admitted workflows is a vital requirement for mission-critical real-time applications. This avoids system resources from being channeled into serving jobs that may not meet their end-to-end deadline, ensuring more efficient utilization of the system resources.

For each experiment, we employ the delay composition framework and holistic approach as admission controllers in independent executions that work on the same input. In each experiment, the same job set (which is sufficient to overload both approaches) is fed to both admission controllers and then the admitted jobs for each approach are executed through the aggregation tree. The overload region is of interest since it determines the capacity of the system under each admission controller.

In the first experiment, we varied the height of the aggregation tree (and hence the number of stages). Figure 6a presents the comparison of average per stage utilization with increasing values of the height of the tree from 1 to 5 (the number of stages varies from 3 to 63) in this experiment. The height of the tree is also the path length for the jobs. The delay composition framework and holistic approach are denoted as DCA and Holistic, respectively. As depicted in the figure, the resource utilization of delay composition framework remains nearly constant. This shows that the pessimism in analysis is independent of the number of stages in the system. On the other hand, the utilization of the holistic approach drops significantly as the system size increases. The reason is that the holistic approach analyzes one stage at a time and the pessimism grows as the path length increases. Hence, our analysis is especially useful for large systems.

In the next experiment, we varied the size of the real-time jobs in the system. The job resolution parameter,  $JR$ , is decreased from 80 to 5 (increasing the job sizes by a factor of 16) and the results for average per stage utilization are shown in Figure 6b. A job resolution of 80 denotes a large number of small jobs, while a value of 5 denotes a small number of large jobs. Under preemptive scheduling, the system utilization of both approaches remains nearly constant. Following a non-preemptive policy, the average per stage utilization of both DCA and holistic analysis techniques strictly decreases as jobs become larger. The reason is that the blocking delay component becomes significantly large as job sizes increase.

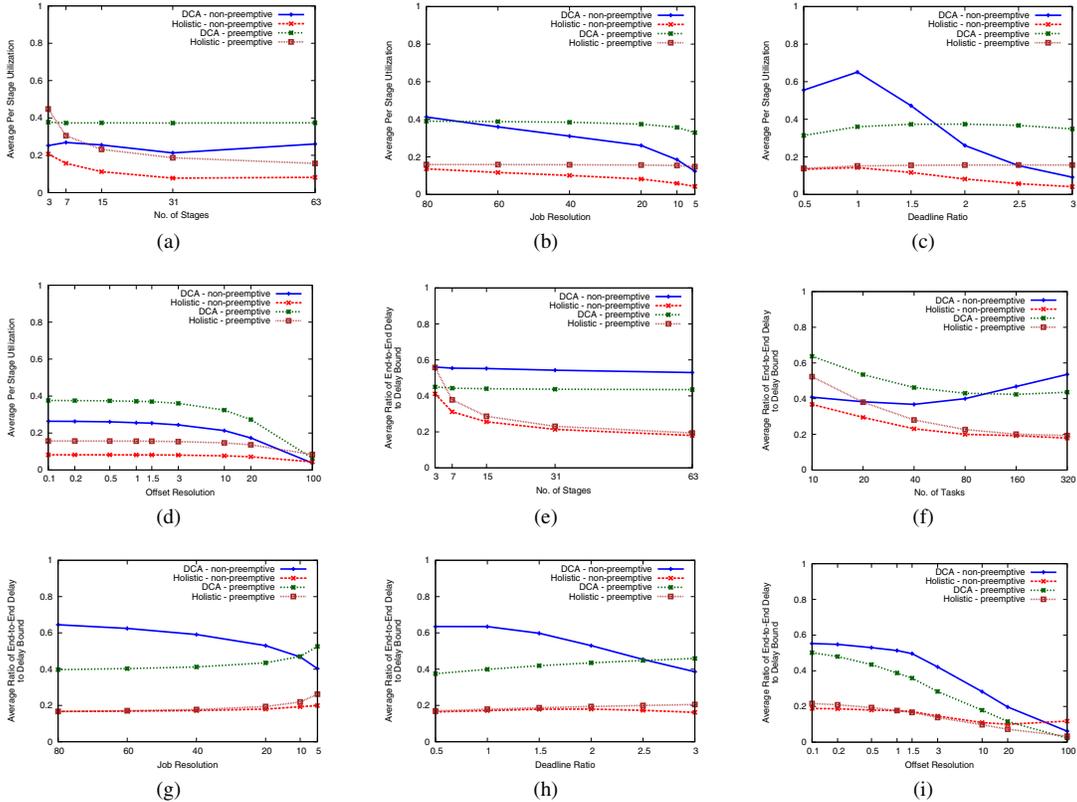


Figure 6: Resource utilization results and end-to-end delay bound accuracy

The delay composition framework consistently outperforms holistic analysis under non-preemptive as well as preemptive scheduling.

Next, we changed the deadline ratio parameter  $DR$ , increasing from 0.5 to 3.0 (recall that the deadlines vary by a factor of  $10^{DR}$ ). Figure 6c shows the average per stage utilization for this experiment. As job sizes are proportional to end-to-end deadlines, when the variability in deadlines increases, an increased number of larger jobs are introduced into the job set. When scheduling follows a non-preemptive policy, larger low priority jobs impose larger blocking delays on higher priority jobs. Therefore, as observed in the figure, under both DCA and holistic analysis, for larger deadline ratios (beyond 1) the effect of blocking delay becomes significant and the average utilization drops thereof. However, for the entire range of deadline ratios, the delay composition framework outperforms holistic analysis.

Figure 6d plots the effect of increasing the offset resolution parameter. We increased the offset resolution value from 0.1 to 100. Note that the  $x$  axis is drawn in logarithmic scale. While the average per stage utilization under holistic analysis is independent of arrival offsets, that of the DCA is a function of offset values. The results report that when offset resolution remains below 100 times of the minimum job deadline, the delay composition framework remains superior to the holistic approach.

### B. End-to-End Delay Bound Accuracy

In this section, we evaluate the tightness of the delay composition framework by comparing the actual delay of

jobs obtained through simulation to the analytically calculated worst-case delay bounds. As we are only interested in delay and not schedulability, we do not perform admission control. We compare our analysis framework with holistic analysis, and in each experiment retain the given workload constant in order to reflect merely the impact of parameters under study.

Figure 6e compares the average ratio of end-to-end delays obtained from simulation to the analytically computed delay bounds with respect to increase in the system size. As the results show, under preemptive scheduling, the holistic approach calculates tighter bounds than DCA for small-scale systems. However, as system scales, the delay bounds of DCA become significantly more accurate than those of holistic (achieving 6% and 24% improvement in accuracy respectively at 7 and 63 stages). Under non-preemptive scheduling, the delay composition framework again outperforms the holistic approach with 14% – 34% difference in the level of tightness.

We then conducted an experiment to evaluate the effect of input workload on delay bounds. Figure 6f shows the corresponding results for job sets of size 10, 20, 40, 80, 160, and 320. As depicted in the figure, the DCA delay bound accuracy is superior to that of the holistic analysis for all workload values. The reason is that the holistic approach relies on per-stage analyses and reflects the contribution of each higher priority job to the end-to-end delay at all stages along its path. In contrast, the delay composition framework takes an end-to-end, per job approach. The latter leads to more accurate estimation of the worst-case delay for DCA. Furthermore, under non-preemptive scheduling policy, while the tightness of holistic delay bound decreases, that of the DCA improves.

This is because for DCA, the difference between the delay bound and the actual delay remains more or less constant with increase in the number of jobs. Therefore, the delay ratio improves for DCA. On the other hand, for holistic analysis, the difference between the delay bound and the actual delay grows significantly with load and hence the delay ratio drops.

Under preemptive scheduling, the delay ratio of both analysis techniques decreases with increase in workload. For DCA, the reason lies in the fact that it accounts for a delay of two stage execution times for the contribution of each higher priority job in delay of other lower priority ones in the worst case (rather than one in non-preemptive scheduling). We encourage the reader to review our prior work [1] for an explanation of this source of pessimism in the analysis. Under both non-preemptive and preemptive scheduling, DCA consistently has a higher ratio of end-to-end delay to the computed delay bound as compared to holistic analysis (more than 20% when the number of jobs is over 80).

Figure 6g plots the results for the effect of job resolution on delay bounds. As depicted in the figures, when changing jobs sizes, the delay composition framework is superior to the holistic approach under both scheduling policies (up to 45% and 25% improvement in tightness, under non-preemptive and preemptive policies, respectively).

The impact of deadline ratio on the tightness is reported in Figure 6h. Recall that a value of  $x$  for deadline ratio implies a  $10^x$  variability in jobs deadlines. As indicated before, the decrease in delay ratio under a non-preemptive policy is related to the blocking delay component becoming more significant. Both analysis approaches predict a larger worst-case blocking factor leading to a reduced value for the delay ratio for both analysis techniques. Here again, the delay composition framework outperforms the holistic approach by about 45% and 25% under non-preemptive and preemptive scheduling schemes, respectively.

The impact of offset in job arrivals is presented in Figure 6i. Note that the  $x$  axis is in logarithmic scale. The worst-case delay bound estimation of the DCA is a function of offset values. The reduction in delay ratios of the holistic is due to the fact that increase in offset values reduces the interference in jobs execution which leads to decrease in the value of actual end-to-end delay. We observe that up to an offset parameter value of nearly 100, DCA outperforms holistic analysis in terms of the achieved average ratio of the end-to-end delay of jobs to the computed delay bound.

## VIII. CONCLUSIONS

In this paper, we investigate timing properties and delay composability of jobs in multisensor data aggregation systems. We propose a theoretical framework which extends previously proposed delay composition theory for a class of systems characterized by a “MERGE” primitive. We provide intuition as to why analyzing such systems can be difficult, and develop a framework to determine offline schedulability of multicriticality distributed workload in such networked systems. Our framework assists to prevent system resources from being

channeled into serving jobs which may not complete their end-to-end execution within prespecified deadlines. We provide delay composition rules for aggregation workflows under non-preemptive as well as preemptive scheduling policies based on characteristics of concurrent workflows and their corresponding arrival offsets. We extensively evaluate our framework through simulations and show that our theoretical framework is significantly more accurate than traditional analysis techniques and effectively utilizes distributed resources. Our framework is especially beneficial for large systems.

## ACKNOWLEDGEMENTS

This research was sponsored in part by ONR grant N00014-10-1-0172 and NSF grant CNS 07-20513 and in part by the Army Research Laboratory and was accomplished under Cooperative Agreement Number W911NF-09-2-0053. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation here on.

## REFERENCES

- [1] P. Jayachandran and T. Abdelzaher, “Transforming acyclic distributed systems into equivalent uniprocessors under preemptive and non-preemptive scheduling,” in *ECRTS*, July 2008, pp. 233–242.
- [2] X. Li and J. Cao., “Coordinated workload scheduling in hierarchical sensor networks for data fusion applications,” *Journal of Computer Science and Technology*, vol. 23, May 2008.
- [3] F. Iandola, F. Saremi, T. Abdelzaher, P. Jayachandran, and A. Yener, “Real-time capacity of networked data fusion,” in *International Conference on Information Fusion*, 2011.
- [4] J. Xu and D. Parnas, “On satisfying timing constraints in hard real-time systems,” *IEEE Transactions on Software Engineering*, vol. 19, no. 1, pp. 70–84, January 1993.
- [5] G. Fohler and K. Ramamritham, “Static scheduling of pipelined periodic tasks in distributed real-time systems,” in *Euromicro Workshop on Real-Time Systems*, June 1997, pp. 128–135.
- [6] B. Kao and H. Garcia-Molina, “Deadline assignment in a distributed soft real-time system,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 8, no. 12, pp. 1268–1274, 1997.
- [7] Y. Zhang, C. Lu, C. Gill, P. Lardieri, and G. Thaker, “End-to-end scheduling strategies for aperiodic tasks in middleware,” University of Washington at St. Louis, Tech. Rep. WUCSE-2005-57, December 2005.
- [8] K. Tindell and J. Clark, “Holistic schedulability analysis for distributed hard real-time systems,” *Elsevier Microprocessing and Microprogramming*, vol. 40, no. 2-3, pp. 117–134, 1994.
- [9] R. Pellizzoni and G. Lipari, “Improved schedulability analysis of real-time transactions with earliest deadline scheduling,” in *RTAS*, March 2005, pp. 66–75.
- [10] L. Thiele, S. Chakraborty, and M. Naedele, “Real-time calculus for scheduling hard real-time systems,” in *IEEE International Symposium on Circuits and Systems*, vol. 4, May 2000, pp. 101–104.
- [11] B. Jonsson, S. Perathoner, L. Thiele, and W. Yi, “Cyclic dependencies in modular performance analysis,” in *ACM EMSOFT*, Oct. 2008, pp. 179–188.
- [12] A. Koubaa and Y.-Q. Song, “Evaluation and improvement of response time bounds for real-time applications under non-preemptive fixed priority scheduling,” *International Journal of Production and Research*, vol. 42, no. 14, pp. 2899–2913, July 2004.
- [13] P. Jayachandran and T. Abdelzaher, “End-to-end delay analysis of distributed systems with cycles in the task graph,” in *ECRTS*, July 2009.
- [14] A. N. Audsley, A. Burns, M. Richardson, and K. Tindell, “Applying new scheduling theory to static priority pre-emptive scheduling,” *Software Engineering*, vol. 8, no. 5, pp. 284–292, 1993.