



Universiteit  
Leiden  
The Netherlands

## Improved hard real-time scheduling and transformations for embedded Streaming Applications

Spasic, J.

### Citation

Spasic, J. (2017, November 14). *Improved hard real-time scheduling and transformations for embedded Streaming Applications*. Retrieved from <https://hdl.handle.net/1887/59459>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/59459>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The following handle holds various files of this Leiden University dissertation:

<http://hdl.handle.net/1887/59459>

**Author:** Spasic, J.

**Title:** Improved hard real-time scheduling and transformations for embedded Streaming Applications

**Issue Date:** 2017-11-14

## Chapter 4

# Exploiting Parallelism in Hard Real-Time Systems to Maximize Performance

Jelena Spasic, Di Liu, Todor Stefanov, “Exploiting Resource-constrained Parallelism in Hard Real-Time Streaming Applications”, *In Proceedings of the International Conference on Design, Automation and Test in Europe (DATE’16)*, pp. 954–959, Dresden, Germany, March 14-18, 2016.

---

**T**HIS chapter presents our solution to the problem of exploiting the right amount of parallelism in a streaming application, **Problem 2** given in Section 1.3, according to an MPSoC platform such that performance is maximized and the timing guarantees are provided. That is, the chapter describes our solution approach consisting of an unfolding graph transformation and an algorithm that adapts the parallelism in the application according to the resources in an MPSoC by using the unfolding transformation.

The remainder of this chapter continues with the problem description in Section 4.1 and summarizes our contributions in Section 4.2. Then, we give an overview of the related work in Section 4.3. A motivational example is given in Section 4.4. It is followed by the description of our proposed solution approach given in Sections 4.5 and 4.6. The experimental evaluation of our proposed approach is presented in Section 4.7. The concluding discussion is given in Section 4.8.

## 4.1 Problem Statement

To meet the computational demands and timing requirements of modern streaming applications, the parallel processing power of MPSoC platforms has to be exploited efficiently. Exploiting the available parallelism in an MPSoC platform to guarantee performance and timing constraints is a challenging task. This is because it requires the designer to expose the right amount of parallelism available in the application and to decide how to allocate and schedule the tasks of the application on the available processing elements such that the platform is utilized efficiently and the timing constraints are met. However, as introduced in Section 1.3, the given initial parallel application specification often is not the most suitable one for the given MPSoC platform. To better utilize the underlying MPSoC platform, the initial specification of an application, that is, the initial task graph, should be transformed by an *unfolding graph transformation* to an alternative one that exposes more parallelism while preserving the same application behavior. The unfolding graph transformations proposed so far: 1) introduce additional tasks for managing data among tasks' replicas [KM08], [FKBS11], which introduces communication and scheduling overhead; 2) do data reordering or increase rates of data production/consumption on channels [ZBS13], [SLA12], which causes an increase of buffer sizes of data communication channels between the tasks and an increase of the application latency. Thus, special care should be taken during the unfolding transformation to avoid all the unnecessary overheads. Moreover, having more tasks' replicas than necessary results in an inefficient system due to overheads in code and data memory, scheduling and inter-tasks communication [FKBS11], [ZBS13]. Thus, the right amount of parallelism (tasks' replicas), that is, the proper values of unfolding factors, depending on the underlying MPSoC platform, should be determined in a parallel application specification to achieve maximum performance and timing guarantees.

Therefore, in this chapter, we investigate the following sub-problems: (1) How to efficiently **unfold a given initial acyclic SDF graph** of an application to avoid unnecessary communication/scheduling overheads and unnecessary increases in buffer sizes and the application latency?, and (2) How to **find a proper unfolding factor of each task in the initial graph**, such that the obtained alternative graph exposes the right amount of parallelism **that maximizes the utilization of the available processors in an MPSoC platform under hard real-time scheduling?**

## 4.2 Contributions

Our contributions to the solution of the research problem described in Section 4.1 are summarized as follows:

- We propose a new unfolding graph transformation for SDF graphs which results in graphs with shorter application latency and smaller buffer sizes compared to the related approaches [KM08], [FKBS11], [SLA12], [ZBS13], as shown in Section 4.7.
- We propose a new algorithm for finding a proper value for the unfolding factor of each task in a graph when mapping the graph on a platform such that the platform is utilized as much as possible under hard real-time scheduling.
- We show, on a set of real-life streaming applications, that in more than 98% of the experiments, our unfolding graph transformation and algorithm result in a solution with a shorter latency, smaller buffer sizes and smaller values for unfolding factors compared to the solution obtained from [ZBS13] while the same performance and timing requirements are satisfied.

*Scope of work.* We assume that a given SDF graph is acyclic. This limitation comes from the hard real-time scheduling framework, presented in Chapter 3, we use to schedule an SDF graph. However, as already mentioned earlier in Chapter 3, even with this limitation our approach is still applicable to many real-life streaming applications because a recent work [TA10] has shown that around 90% of streaming applications can be modeled as acyclic SDF graphs. In addition, our approach does not unfold *stateful* tasks and input/output tasks. A *stateful* task is a task which current execution depends on its previous execution, thus those executions cannot be run in parallel. Input and output tasks are the tasks connected to the environment, hence they are not unfolded.

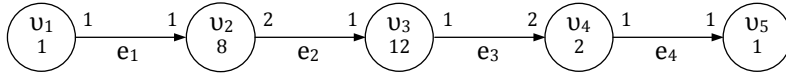
## 4.3 Related Work

[KM08] proposes an Integer Linear Programming (ILP) based approach for maximizing the throughput of an application modeled as an SDF graph by exploiting data parallelism when mapping the application on a platform with fixed number of processors. However, an ILP-based approach suffers from an exponential worst-case time complexity. To overcome the time complexity issue of the approach in [KM08], [FKBS11] separates the task replication and the allocation of replicas. However, decomposing the problem into two strongly related problems and solving them separately has a negative impact on the

solution quality. In addition, the maximum data-level parallelism is revealed in the application without considering the platform constraints. In contrast, in our approach, we solve the problem of task replication and the mapping of replicas simultaneously while taking into account the platform constraints. Both approaches [KM08] and [FKBS11] use *splitter* (S) and *joiner* (J) tasks to distribute and merge data streams processed by replicas, see Figure 4.2(a). Those tasks introduce additional communication overhead as data streams have to be sent to them and to the replicas. Moreover, the splitter/joiner tasks have to be considered in the process of mapping and scheduling of tasks. In contrast, in our approach, we do not introduce additional tasks for data management, but we propose a new transformation on an SDF graph in Section 4.5 where the data is sent by replicas of the original tasks only to replicas which need the data for computation. Thus, we avoid the overhead of scheduling splitter/joiner tasks and duplicated data transfers, as shown in Section 4.7.

[SLA12] proposes a throughput driven transformation of an application modeled as an SDF graph for mapping the application on a platform. The graph transformation method in [SLA12] increases the rates of data production/consumption and hence increases the buffer capacities needed to store the data, see Figure 4.2(b). In addition, to enable unfolding of tasks, multiple firings of a certain task in the initial graph are combined into one firing of the corresponding task in the transformed graph, see the increased execution times of tasks in Figure 4.2(b), which leads to an increase in latency. In contrast, our transformation technique does not increase the rates of data production/consumption on communication channels and does not combine multiple task firings into one firing which in turn leads to shorter application latency and smaller buffer sizes of the communication channels, as shown in Section 4.7.

The closest to our work, in terms of scope and methods proposed to efficiently utilize the parallelism of an application mapped onto resource-constrained platform, is the work in [ZBS13]. The authors in [ZBS13] propose an approach for exploiting *just-enough* parallelism when mapping a streaming application modeled as an SDF graph on a platform with fixed number of processing elements. The graph transformation method in [ZBS13] transforms an initial SDF graph to functionally equivalent CSDF graph while keeping the same rates of data production/consumption on communication channels, see Figure 4.2(c). However, the transformation approach in [ZBS13] is not efficient in terms of application latency and buffer sizes of the communication channels, as shown in Section 4.7. Moreover, the proposed algorithm in [ZBS13] for finding the values of unfolding factors and the mapping of task replicas does

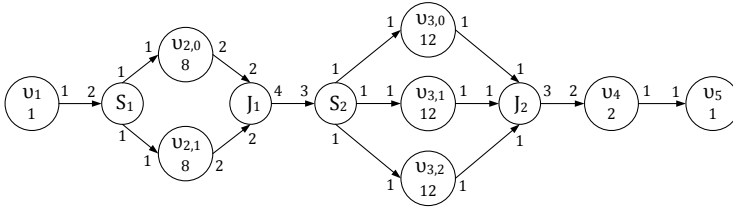
Figure 4.1: An SDF graph  $G$ .

not reveal the right amount of parallelism, but it reveals more parallelism than needed and hence the platform is unnecessarily overloaded, as shown in Section 4.7. In contrast, the approach we propose unfolds a graph by doing more aggressive token-flow analysis leading to shorter application latency and smaller buffer sizes. In addition, our approach finds smaller unfolding factors for tasks which leads to less memory needed to store the code of replicas and less memory to implement communication channels between the replicas.

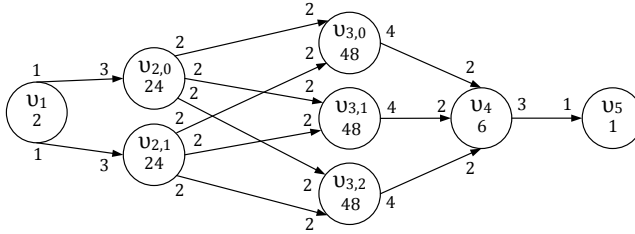
## 4.4 Motivational Example

In the first part of this section, we motivate the need for our new unfolding graph transformation. The throughput of graph  $G$  given in Figure 4.1 when scheduled under our ISPS presented in Chapter 3 is the same as the throughput obtained under *self-timed* scheduling [SGB08] and it is equal to  $\frac{1}{24}$ . Note that an unfolding graph transformation is used to increase the application throughput if it is allowed by the hardware platform on which the application is executed. Let us assume that actors  $v_2$  and  $v_3$  of graph  $G$  in Figure 4.1 are unfolded by factors 2 and 3, respectively, in order to increase the throughput of  $G$ . Figure 4.2 shows four functionally equivalent graphs obtained after applying the unfolding transformations proposed by [KM08], [FKBS11] – see Figure 4.2(a), by [SLA12] – see Figure 4.2(b), and by the transformation in [ZBS13] – see Figure 4.2(c), while the graph given in Figure 4.2(d) is obtained by applying our transformation described in Section 4.5. Our transformation method unfolds an SDF graph by doing more aggressive data token flow analysis with the aim to spread equally the workload of an actor during the hyperperiod and run in parallel as much replicas of the actor as possible.

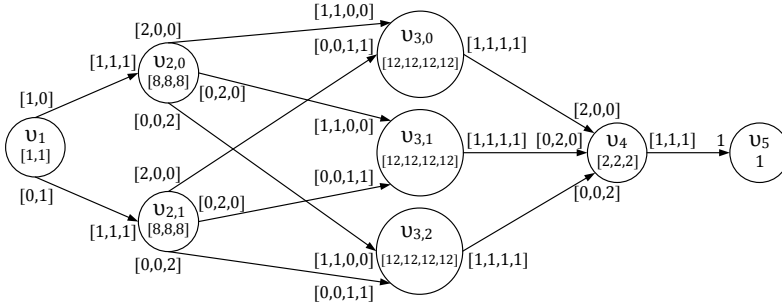
Table 4.1 gives for all four equivalent graphs of  $G$  the throughput  $\mathcal{R}_{\text{out}}$  of the output actor, actor  $v_5$ , the maximum latency  $\mathcal{L}_{\text{in} \rightarrow \text{out}}$  on an input-output path, the total size  $M$ , of the communication buffers, the total code size  $CS$ , and the total number of processors  $m$  needed to schedule the graphs under ISPS and the self-timed scheduling while achieving the same throughput  $\mathcal{R}_{\text{out}}$ . We can see from the table that by applying our unfolding transformation we can obtain, under ISPS, 2.29, 3.14, and 1.43 times shorter latency and 2.08, 2.75, and 1.33 times smaller buffers than the unfolding methods in [KM08]



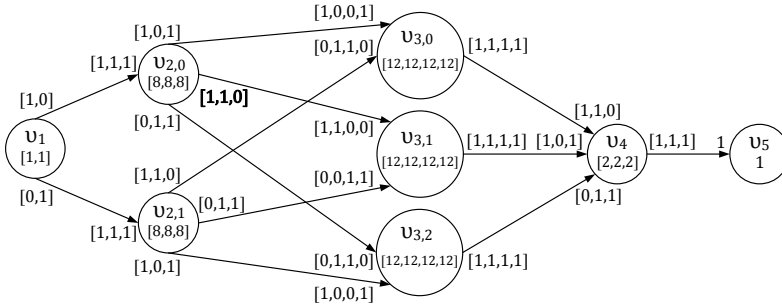
(a) Equivalent of G in Figure 4.1 after the transformation in [KM08], [FKBS11]



(b) Equivalent of G in Figure 4.1 after the transformation in [SLA12]



(c) Equivalent of G in Figure 4.1 after the transformation in [ZBS13]



(d) Equivalent of G in Figure 4.1 after our transformation

**Figure 4.2:** Equivalent graphs of the SDF graph in Figure 4.1 by unfolding actor  $v_2$  by factor 2 and  $v_3$  by factor 3.



**Table 4.1:** Results for  $G$  transformed by different transformation approaches.

Approach	ISPS					[SGB08]				
	$\mathcal{R}_{\text{out}}[\frac{1}{\mu\text{s}}]$	$\mathcal{L}_{\text{in} \rightarrow \text{out}}[\mu\text{s}]$	$M[\text{B}]$	$CS[\text{kB}]$	$m$	$\mathcal{R}_{\text{out}}[\frac{1}{\mu\text{s}}]$	$\mathcal{L}_{\text{in} \rightarrow \text{out}}[\mu\text{s}]$	$M[\text{B}]$	$CS[\text{kB}]$	$m$
[KM08], [FKBS11]	1/8	128	50	40	5	1/8	67	31	40	12
[SLA12]	1/8	176	66	36	5	1/8	93	57	36	8
[ZBS13]	1/8	80	32	36	5	1/8	76	24	36	8
our	1/8	56	24	36	5	1/8	62	21	36	8

**Table 4.2:** Results for  $G$  transformed and mapped on 2 processors by different approaches.

Approach	$\mathcal{R}_{\text{out}}[\frac{1}{\mu\text{s}}]$	$\mathcal{L}_{\text{in} \rightarrow \text{out}}[\mu\text{s}]$	$M[\text{B}]$	$CS[\text{kB}]$	$m$
[ZBS13]	1/18	180	31	44	2
our	1/18	108	16	32	2

and [FKBS11], [SLA12], and [ZBS13], respectively. The number of processors needed to schedule the graph obtained after the transformation under ISPS is equal for all the transformation methods. Under self-timed scheduling [SGB08] we obtain 1.08, 1.5, and 1.23 times shorter latency, while buffers are smaller 1.47, 2.71, and 1.14 times compared to the related approaches. Assuming one-to-one mapping for the self-timed scheduling, we need the same number of processors to schedule the unfolded graph obtained by the methods in [SLA12] and [ZBS13], and 1.5 times less processors than the unfolding methods in [KM08] and [FKBS11]. For both scheduling algorithms we obtain equal code size as the unfolding methods in [SLA12] and [ZBS13], and 1.11 times smaller code size than the methods in [KM08] and [FKBS11]. From Table 4.1, we see that our unfolding transformation approach presented in Section 4.5 is more efficient than the approaches in [KM08], [FKBS11], [SLA12], and [ZBS13].

So far, we considered only the unfolding transformation. Now, we would like to focus on the algorithm for finding the proper unfolding factors for actors when a graph is mapped onto resource-constrained platform and scheduled by a hard real-time scheduler such that the throughput of the graph is maximized. Here, we want to compare our algorithm in Section 4.6 with the approach in [ZBS13], because only that approach, among the related approaches, exploits the parallelism in an application under hard real-time scheduling. For example, in order to schedule graph  $G$  in Figure 4.1 on a platform with 2 processors while maximizing the throughput under hard real-time scheduling, the approach in [ZBS13] finds a vector of unfolding factors  $\vec{f} = [1, 2, 4, 1, 1]$ .

However, there exists a smaller vector of unfolding factors, such as  $\vec{f} = [1, 1, 3, 1, 1]$ , such that  $G$  is schedulable on 2 processors and the throughput is maximized. This smaller vector  $\vec{f}$  is found by our algorithm in Section 4.6. Table 4.2 gives the throughput  $\mathcal{R}_{\text{out}}$ , latency  $\mathcal{L}_{\text{in} \rightarrow \text{out}}$ , buffer sizes  $M$  and code

size  $CS$  when  $G$  is unfolded and mapped on  $m = 2$  processors by applying the approach in [ZBS13] and by applying our algorithm presented in Section 4.6. We can see from the table that by applying our algorithm we obtain under ISPS 1.67 times shorter latency, 1.94 times smaller buffers, and 1.38 smaller code size than the approach in [ZBS13]. From these results and the results given in Table 4.1, we clearly show the necessity and usefulness of the graph unfolding transformation presented in Section 4.5, and the algorithm for finding proper values for the unfolding factors presented in Section 4.6.

## 4.5 New Unfolding Transformation for SDF Graphs

Our new unfolding transformation method is given in Algorithm 3. The algorithm takes an SDF graph  $G$  and a vector of unfolding factors  $\vec{f}$  and produces an unfolded graph  $G'$ , which is a CSDF graph. The initial SDF graph and its unfolded version given in the form of a CSDF graph are functionally equivalent, meaning that both of them generate the same sequence of output data tokens for a given sequence of input data tokens. The algorithm consists of three phases. The first phase is given in lines 1 to 4 in Algorithm 3. Given that the execution semantics of the SDF model allows any integer multiple of the basic repetition vector also as a valid repetition vector, in line 1 of Algorithm 3 the basic repetition vector  $\vec{q}$  of  $G$  is replaced by  $\vec{q}^f = \text{lcm}(\vec{f}) \cdot \vec{q}$ , where  $\text{lcm}(\vec{f})$  is the least common multiple of all elements in  $\vec{f}$ . Then in lines 2 to 4, for each channel  $e_u$  in  $G$ , a matrix  $d$  is constructed containing as many columns as the number of tokens produced/consumed on the channel during one iteration of  $G$  with repetition vector  $\vec{q}^f$ . Each column in  $d$  contains in row 0 an index  $p$ ,  $d[0][t] = p$ , which is the index of the firing of the producer actor,  $p \geq 0$ , which produces the  $t^{\text{th}}$  token, and an index  $c$  in row 1,  $d[1][t] = c$ , representing the index of the firing of the consumer actor,  $c \geq 0$ , which consumes the  $t^{\text{th}}$  token on  $e_u$ . Constructing matrix  $d$  for channel  $e_2$  of graph  $G$  in Figure 4.1 when  $\vec{f} = [1, 2, 3, 1, 1]$  is given in Figure 4.3, lines 2 to 4.

In the second phase the topology of the equivalent CSDF graph  $G'$  is created, which is given in lines 5 to 14 in Algorithm 3. In the equivalent CSDF graph  $G'$ , every actor is replicated a certain number of times, as determined by the unfolding vector, lines 6 to 8. Then each channel  $e_u$  in the initial graph is replicated a certain number of times in the equivalent graph such that each replica of the producer on  $e_u$  is connected to each replica of the consumer on  $e_u$ , as given in lines 9 to 12 of Algorithm 3. The motivation behind unfolding is to equally distribute the workload of an actor in the initial graph by running in parallel replicas corresponding to that actor. The workload of an actor within

**Algorithm 3:** Procedure to unfold an SDF graph.

---

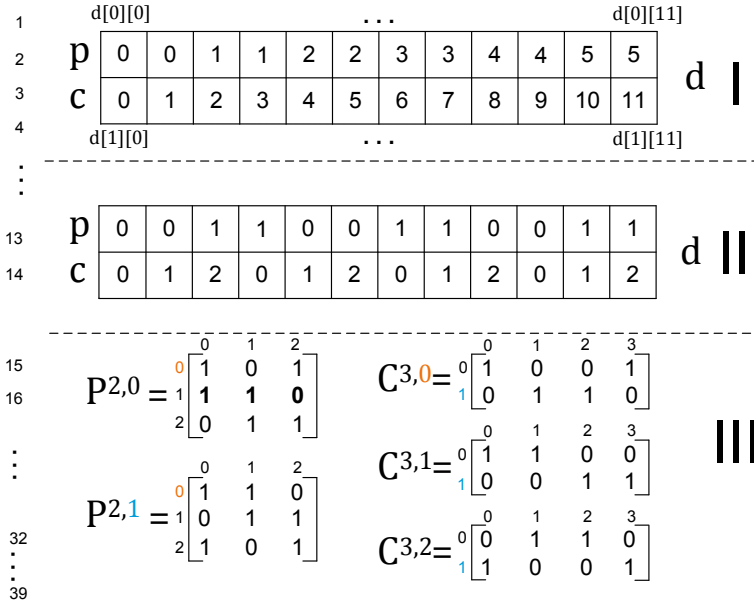
**Input:** An SDF graph  $G = (\mathcal{V}, \mathcal{E})$ , a vector of unfolding factors  $\vec{f}$ .  
**Output:** The equivalent CSDF graph  $G' = (\mathcal{V}', \mathcal{E}')$ .

```

1 Take  $\vec{q}^f = [\text{lcm}(\vec{f}) \cdot q_1, \dots, \text{lcm}(\vec{f}) \cdot q_N]$  as a repetition vector of  $G$ ;
2 for communication channel  $e_u = (v_i, v_j) \in \mathcal{E}$  do
3   Get production rate  $prd$  and consumption rate  $cns$  on  $e_u$ ;
4   Construct a matrix  $d$ ,  $d[0][t] = p$ ,  $d[1][t] = c$ ,  $t \in [0, prd \cdot q_i^f - 1]$ ,  $p$  is the index of  $v_i$  firing
   which produces  $t^{th}$  token,  $c$  is the index of  $v_j$  firing which consumes  $t^{th}$  token on  $e_u$ ;
5  $\mathcal{V}' \leftarrow \emptyset, \mathcal{E}' \leftarrow \emptyset$ ;
6 for actor  $v_i \in \mathcal{V}$  do
7   for  $k = 0$  to  $f_i - 1$  do
8     Add replica  $v_{i,k}$  to  $\mathcal{V}'$ ;
9 for communication channel  $e_u = (v_i, v_j) \in \mathcal{E}$  do
10   for replica  $v_{i,k}$  of  $v_i$  do
11     for replica  $v_{j,l}$  of  $v_j$  do
12       Add  $e'_u = (v_{i,k}, v_{j,l})$  to  $\mathcal{E}'$ ;
13   for  $t = 0$  to  $prd \cdot q_i^f - 1$  do
14      $d[0][t] = d[0][t] \bmod f_i$ ,  $d[1][t] = d[1][t] \bmod f_j$ ;
15 for communication channel  $e_u = (v_i, v_j) \in \mathcal{E}$  do
16   Get production rate  $prd$  and consumption rate  $cns$  on  $e_u$ ;
17   Create empty/zero matrices  $P^{i,k}$  with size  $f_j \times q_{i,k}$ ,  $k \in [0, f_i - 1]$ ;
18   Create empty/zero matrices  $C^{j,l}$  with size  $f_i \times q_{j,l}$ ,  $l \in [0, f_j - 1]$ ;
19   for  $h = 0$  to  $q_{i,0} - 1$  do
20     for  $k = 0$  to  $f_i - 1$  do
21       Initialize a prod. counter seq.  $cnt_{\text{prod}}$  of length  $f_j$  to 0;
22       for  $o = 0$  to  $prd - 1$  do
23          $cnt_{\text{prod}}[d[1][h \cdot k \cdot prd + o]] = cnt_{\text{prod}}[d[1][h \cdot k \cdot prd + o]] + 1$ ;
24       for  $l = 0$  to  $f_j - 1$  do
25          $P^{i,k}[l][h] = cnt_{\text{prod}}[l]$ ;
26   for  $h = 0$  to  $q_{j,0} - 1$  do
27     for  $l = 0$  to  $f_j - 1$  do
28       Initialize a cons. counter seq.  $cnt_{\text{cons}}$  of length  $f_i$  to 0;
29       for  $o = 0$  to  $cns - 1$  do
30          $cnt_{\text{cons}}[d[0][h \cdot l \cdot cns + o]] = cnt_{\text{cons}}[d[0][h \cdot l \cdot cns + o]] + 1$ ;
31       for  $k = 0$  to  $f_i - 1$  do
32          $C^{j,l}[k][h] = cnt_{\text{cons}}[k]$ ;
33   for  $k = 0$  to  $f_i - 1$  do
34     for  $l = 0$  to  $f_j - 1$  do
35       if all entries in row  $P^{i,k}[l][\ ]$  are 0 then
36         Delete a channel  $e'_u$  connecting replicas  $v_{i,k}$  and  $v_{j,l}$ ;
37       else
38         Associate production sequence  $P^{i,k}[l][\ ]$  and consumption sequence  $C^{j,l}[k][\ ]$ 
         with  $e'_u = (v_{i,k}, v_{j,l})$ ;
39 return  $G'$ ;

```

---



**Figure 4.3:** Unfolding channel  $e_2$  from the graph in Figure 4.1 by using Algorithm 3 when  $\vec{f} = [1, 2, 3, 1, 1]$ .

one graph iteration is determined by the corresponding repetition value of the actor. Thus, each replica  $v_{i,k} \in G'$  of an actor  $v_i \in G$  will have the repetition  $q_{i,k}$ :

$$q_{i,k} = \frac{q_i^f}{f_i} = \frac{q_i \cdot \text{lcm}(\vec{f})}{f_i}. \quad (4.1)$$

For example, after the unfolding of the SDF graph in Figure 4.1 with the unfolding vector  $\vec{f} = [1, 2, 3, 1, 1]$  we obtain the graph shown in Figure 4.2(d) with the repetition vector  $\vec{q}' = [6, 3, 3, 4, 4, 4, 6, 6]$ , where  $q_{2,0} = q_{2,1} = \frac{1 \cdot \text{lcm}(1, 2, 3, 1, 1)}{2} = 3$ . Lines 13 to 14 convert the firing production/consumption indexes in  $d[0][t]/d[1][t]$  for each token  $t$  produced/consumed on channel  $e_u$  into the indexes corresponding to the index  $k$  of the replica which produces/consumes  $t$ . This is illustrated for channel  $e_2$  in Figure 4.3, lines 13 and 14.

Lines 15 to 39 represent the third phase of Algorithm 3 and they derive the production and consumption sequences for new channels and perform final placement of the new channels between the corresponding actor replicas. More specifically, for each source replica  $v_{i,k}$  and destination replica  $v_{j,l}$  of a channel, a production matrix  $P^{i,k}$  and consumption matrix  $C^{j,l}$  is created from matrix  $d$  in lines 15 to 32. The index of each raw in a production matrix

$P^{i,k}$  corresponds to the index  $l$  of a destination replica  $v_{j,l}$ . The index of each column in a production matrix  $P^{i,k}$  corresponds to the firing index of source replica  $v_{i,k}$ . Elements in a production matrix of a source replica contain the number of tokens produced by a certain firing of that replica. Similar holds for elements in a consumption matrix. The created matrices  $P^{2,0}$ ,  $P^{2,1}$ ,  $C^{3,0}$ ,  $C^{3,1}$ ,  $C^{3,2}$  for the source replicas  $v_{2,0}$ ,  $v_{2,1}$  and destination replicas  $v_{3,0}$ ,  $v_{3,1}$ ,  $v_{3,2}$  on channel  $e_2$  are given in Figure 4.3. For example, the value 0 in element  $P^{2,0}[2][0]$  says that 0 tokens are produced by the 0 firing of source replica  $v_{2,0}$  for the destination replica  $v_{3,2}$ . Once these matrices are constructed, the production and consumption sequences on channel replicas are extracted from the corresponding rows in matrices, as given in lines 33 to 38 in Algorithm 3. For example, the production sequence on the channel between  $v_{2,0}$  and  $v_{3,1}$  in Figure 4.2(d) is extracted from row  $P^{2,0}[1][\ ]$  in matrix  $P^{2,0}$  and is equal to  $[1, 1, 0]$ . The extracted production/consumption sequences on replicas of channel  $e_2$  can be seen in Figure 4.2(d). The unfolded graph  $G'$  is returned in line 39 of Algorithm 3.

## 4.6 The Algorithm for Finding Proper Unfolding Factors

In order to efficiently utilize the parallelism available in an application when mapping the application on a resource-constrained platform under hard real-time scheduling, proper unfolding factors for actors of the application have to be determined. Therefore, in this section, we present an algorithm which derives the proper unfolding factors which maximize the utilization of the platform, that is, maximize the application throughput.

The algorithm is given in Algorithm 4. It takes an SDF graph  $G$ , where the actors are scheduled by ISPS presented in Chapter 3, a platform with  $m$  processors, a scheduling algorithm  $A$  [LL73], an allocation heuristic  $H$  [CGJ96] and a quality factor  $\rho$ . A quality factor  $\rho \in (0, 1]$  determines how much of the platform processing resources we want to utilize, with  $\rho = 1$  corresponding to full utilization. The algorithm returns the best solution vector of unfolding factors  $\vec{f}^{best}$ .

Line 1 in Algorithm 4 initializes each unfolding factor of an actor in  $G$  to 1 and  $G'$  to  $G$ . Then, the upper bound  $\hat{f}_i$  of unfolding factor  $f_i$  for each actor  $v_i$  in  $G$  is computed in line 2 in Algorithm 4 by using Equation (4.2) which is

similar to Equation (3) in [ZBS13]:

$$\hat{f}_i = \frac{\text{lcm}\{x_1, x_2, \dots, x_n\}}{x_i}, \quad (4.2)$$

where  $x_i = \frac{\text{lcm}\{W_1, W_2, \dots, W_n\}}{W_i}$  and  $W_i = \sum_{\varphi=1}^{\varphi=\phi_i} C_i(\varphi) \cdot r_i$  is the workload of actor  $v_i$  during one hyperperiod. The logic behind the computation of upper bounds on unfolding factors is the same as in [ZBS13]. That is, by unfolding every actor  $v_i$  in  $G$  by an upper bound  $\hat{f}_i$  we obtain a CSDF graph  $G'$ , for which the minimum period  $\check{T}_{i,k}$  of each replica  $v_{i,k}$  under ISPS is equal to  $\sum_{\varphi=1}^{\varphi=\phi_{i,k}} C_{i,k}(\varphi)$ , meaning that each actor in the unfolded graph fully utilizes the processor which it runs on, hence, leading to the maximum throughput. Line 3 finds the utilization of graph  $G'$  when  $G'$  is scheduled on  $m$  processors by invoking Algorithm 5. The best utilization of  $G'$  is initialized in line 4 to be the first schedulable solution on  $m$  processors found by Algorithm 5 in line 3. Line 5 finds the bottleneck actor in  $G'$ . The bottleneck actor  $v_{b,k}$  is the actor with the heaviest workload during one hyperperiod,  $W_{b,k} = \max_{v_{i,k} \in \mathcal{V}'} W_{i,k}$ . If multiple actors have the same maximum workload, then the one with the smallest code size is selected to be the bottleneck. If the current utilization  $u_{G'}$  does not meet the quality requirement checked in line 6, the unfolding factor  $f_b$  of the bottleneck actor  $v_{b,k}$  is increased in line 7 and the graph is unfolded by using Algorithm 3 in line 8. Note that stateful actors and input and output actors are not unfolded, that is, the upper bound on their unfolding factors is 1. The utilization  $u_{G'}$  of the unfolded graph  $G'$  mapped on  $m$  processors is calculated in line 9 by Algorithm 5. If the current utilization  $u_{G'}$  is higher than the best utilization in line 10, then in line 11 the best utilization becomes the one found in line 9 and the best solution vector of unfolding factors becomes the current vector of unfolding factors. Line 12 finds the bottleneck actor in the unfolded graph  $G'$ . Lines 6 to 12 are repeated and the algorithm terminates when either a pre-specified quality factor  $\rho$  is satisfied ( $u_{G'} \geq \rho \cdot m$ ) or the unfolding factor of a bottleneck actor exceeds its upper bound  $\hat{f}_b$  ( $f_b \geq \hat{f}_b$ ).

We see that Algorithm 4 uses Algorithm 5 for finding the utilization of the unfolded graph  $G'$  when mapped on a platform with  $m$  processors. Algorithm 5 takes the unfolded CSDF graph  $G'$ , a platform with  $m$  processors, a scheduling algorithm  $A$  [LL73] and an allocation heuristic  $H$  [CGJ96] as inputs. Line 1 calculates periods of actors in  $G'$  scheduled by ISPS presented in Chapter 3 by using Equation (3.5). Equation (3.5) can be written as Equation (4.3) and Equation (4.4):

$$T_i = \frac{\text{lcm}(\vec{r})}{r_i} \cdot s, \forall v_i \in \mathcal{V}, \quad (4.3)$$

---

**Algorithm 4:** Finding proper unfolding factors for an SDF graph mapped onto resource-constrained platform.

---

**Input:** An SDF graph  $G$ , the number of processors in a platform  $m$ , quality factor  $\rho$ , a scheduling algorithm  $A$ , an allocation heuristic  $H$ .

**Output:** Vector of unfolding factors  $\vec{f}^{best}$ .

```

1  $\vec{f} = [1, 1, \dots, 1]$ ;  $G' = G$ ;
2 Compute the upper bound  $\hat{\vec{f}}$  of  $\vec{f}$  by Equation (4.2);
3 Get  $u_{G'}$  of  $G'$  by Algorithm 5 when scheduled by  $A$  and  $H$  on  $m$ ;
4  $u_{G^{best}} = u_{G'}$ ;  $\vec{f}^{best} = \vec{f}$ ;
5 Find the bottleneck actor  $v_{b,k}$  in  $G'$ ;
6 while  $u_{G'} < \rho \cdot m$  and  $f_b < \hat{f}_b$  do
7    $f_b = \hat{f}_b + 1$ ;
8   Get  $G'$  by unfolding  $G$  by Algorithm 3;
9   Get  $u_{G'}$  of  $G'$  by Algorithm 5 when scheduled by  $A$  and  $H$  on  $m$ ;
10  if  $u_{G'} > u_{G^{best}}$  then
11     $u_{G^{best}} = u_{G'}$ ;  $\vec{f}^{best} = \vec{f}$ ;
12  Find the bottleneck actor  $v_{b,k}$  in  $G'$ ;
13 return  $\vec{f}^{best}$ .
```

---



---

**Algorithm 5:** Procedure to find the utilization of a CSDF graph mapped onto resource-constrained platform.

---

**Input:** A CSDF graph  $G'$ , the number of processors in a platform  $m$ , a scheduling algorithm  $A$ , an allocation heuristic  $H$ .

**Output:** Graph utilization  $u_{G'}$ .

```

1 Calculate  $s$  by Equation (4.4); calculate  $T_i$  by Equation (4.3) by using the calculated  $s$ ;
2 Calculate  $u_{G'}$  by Equation (4.5);
3 while  $G'$  is not schedulable on  $m$  by  $A$  and  $H$  do
4    $s = s + 1$ ;
5   Calculate  $T_i$  by using  $s$  in Equation (4.3); calculate  $u_{G'}$  by Equation (4.5);
6 return  $u_{G'}$ .
```

---

$$s = \left\lceil \frac{\hat{W}}{\text{lcm}(\vec{r})} \right\rceil, \quad (4.4)$$

where  $\text{lcm}(\vec{r})$  is the least common multiple of all repetition entries in  $\vec{r}$  and  $\hat{W} = \max_{v_i \in \mathcal{V}} W_i$  is the maximum workload during one hyperperiod. Note that periods computed by Equation (4.3) are the minimum periods for actors scheduled by ISPS and that there exist other larger valid periods for actors

by taking any integer  $s > \left\lceil \frac{\hat{W}}{\text{lcm}(\vec{r})} \right\rceil$ . Once the actor periods are computed, the

utilization of actor  $v_i$ , denoted as  $u_i$ , can be computed as  $u_i = \sum_{\varphi=1}^{\varphi=\phi_i} C_i(\varphi) / T_i$ , where  $u_i \in (0, 1]$ . For a graph  $G$ ,  $u_G$  is the total utilization of  $G$  given by:

$$u_G = \sum_{v_i \in \mathcal{V}} u_i = \sum_{v_i \in \mathcal{V}} \frac{\sum_{\varphi=1}^{\varphi=\phi_i} C_i(\varphi)}{T_i}. \quad (4.5)$$

The total utilization of a graph directly determines the minimum number of processors needed to schedule the graph, as explained earlier in Section 2.2.5. The utilization  $u_{G'}$  of  $G'$  is calculated in line 2 in Algorithm 5 by using Equation (4.5). Actor periods computed by Equation (4.3) and Equation (4.4) represent the minimum periods when the actors are scheduled under ISPS on a platform with unlimited number of processors. It may happen that these minimum periods lead to a graph which is not schedulable on a platform with only  $m$  processors. Hence, in line 3 by using the utilization  $u_{G'}$  calculated in line 2 we check if  $G'$  can be scheduled on  $m$  processors by using the corresponding schedulability test for  $A$  and  $H$  [DB11]. If  $G'$  is not schedulable on the platform, we decrease  $u_{G'}$  until  $G'$  becomes schedulable by increasing the actor periods  $T_i$ . This is done in lines 4 and 5 in Algorithm 5. Once the graph  $G'$  becomes schedulable on  $m$  processors by  $A$  and  $H$ , Algorithm 5 returns the utilization of the unfolded graph  $G'$  in line 6.

## 4.7 Evaluation

We present two experiments to evaluate the techniques proposed in Section 4.5 and Section 4.6. In the first experiment, we evaluate the efficiency of our unfolding transformation in comparison to the unfolding transformation methods in [KM08], [FKBS11], [SLA12], and [ZBS13]. In the second experiment, we evaluate the efficiency of Algorithm 4 presented in Section 4.6 in terms of performance and time complexity by comparing our approach to the related approach in [ZBS13]. The experiments were performed on the real-life applications from the StreamIt benchmarks suit [TA10], given in Table 4.3. These applications were modeled as SDF graphs and  $|\mathcal{V}|$  denotes the number of actors in an SDF graph, while  $|\mathcal{E}|$  denotes the number of communication channels. The results of the evaluations are shown in Figure 4.4, Figure 4.5, and Figure 4.6. In all these figures, each vertical line shows the variations in the corresponding results among all the applications. The upper and lower ends of a vertical line represent the maximum and minimum values of the corresponding result while the marker at the middle of each vertical line represents the geometric mean of the result. Note that the Y axis in Figure 4.4 to



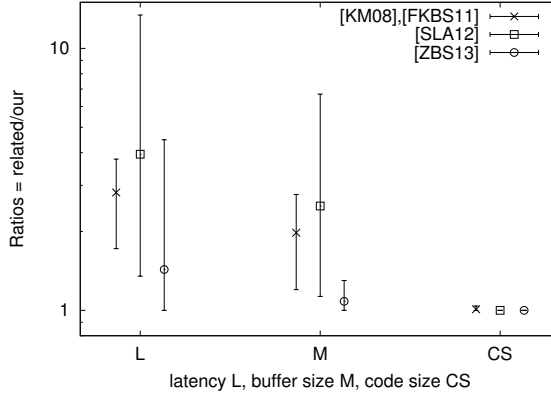
**Table 4.3:** *Benchmarks used for evaluation.*

Benchmark	$ \mathcal{V} $	$ \mathcal{E} $
Discrete cosine transform (DCT)	8	7
Fast Fourier transform (FFT)	17	16
Time delay equalization (TDE)	29	28
Data encryption standard (DES)	53	60
Bitonic Sorting	40	46
Channel Vocoder	55	70
Filterbank	85	99
Serpent	120	128
MPEG2	23	26
Vocoder	114	147
FMRadio	43	53

Figure 4.6 has a logarithmic scale. We run all the experiments on an Intel Core i7-2620M CPU running at 2.70 GHz with Linux Ubuntu 12.4.

#### 4.7.1 Efficiency of the Proposed Unfolding Transformation

In this section, we evaluate the performance of our unfolding transformation method proposed in Section 4.5 by comparison to the related unfolding transformation methods in [KM08], [FKBS11], [SLA12], and [ZBS13]. In this experiment, first we use Algorithm 4 to find a vector of unfolding factors for each application in Table 4.3 mapped on a platform with 64 processors with partitioned First-Fit Decreasing Earliest Deadline First (FFD-EDF) scheduler and quality factor  $\rho = 0.9$ . Then, for each application, we use the found vector of unfolding factors to unfold the application graph by applying our transformation method and the related transformation methods [KM08], [FKBS11], [SLA12], [ZBS13]. Finally, we use the ISPS framework presented in Chapter 3 to calculate the latency, buffer sizes and code size when the unfolded graphs are scheduled by FFD-EDF on 64-processor platform. The ratios between the results obtained by *related* transformation methods and *our* transformation in terms of application latency ( $\mathcal{L}$ ), buffer sizes ( $M$ ) and code size (CS) are given in Figure 4.4. We can see that our method outperforms all the related methods, and delivers on average 2.82, 3.95, and 1.43 times shorter latency and 1.98, 2.5, and 1.08 times smaller buffers than the method in [KM08] and [FKBS11], [SLA12], and [ZBS13], respectively. Although the methods in [KM08] and [FKBS11] introduce additional actors for data management, the average increase in the total code size is only 1%. The other two transformation methods, [SLA12]

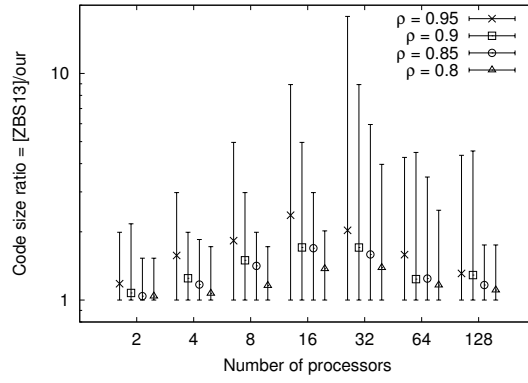


**Figure 4.4:** Comparison of our unfolding transformation to the approaches in [KM08], [FKBS11], [SLA12], [ZBS13].

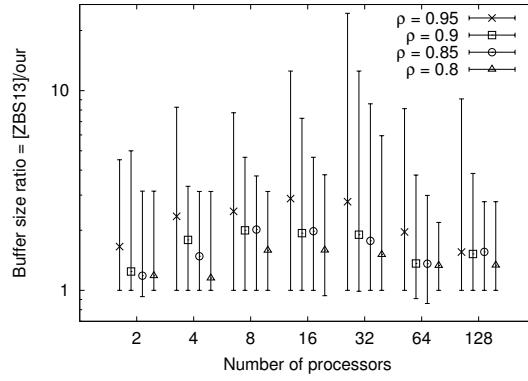
and [ZBS13], have the same code size as our method. Note that all the methods achieve the same application throughput.

#### 4.7.2 Performance of Algorithm 4

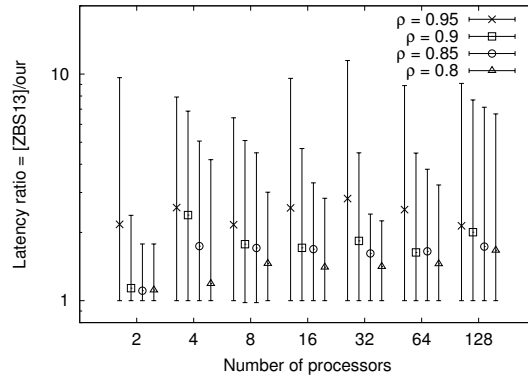
We evaluate the performance of Algorithm 4 by comparison to the related approach in [ZBS13]. For each application  $app$  in Table 4.3, we construct 28 system configurations  $(app, m, \rho)$  with number of processors  $m \in \{2, 4, 8, 16, 32, 64, 128\}$ , and utilization quality  $\rho \in \{0.8, 0.85, 0.9, 0.95\}$ . We run Algorithm 4 with FFD-EDF scheduler for each  $(app, m, \rho)$  configuration to obtain a vector of unfolding factors  $\vec{f}^{best}$ . Then, for each configuration, we unfold the corresponding application graph by the obtained vector  $\vec{f}^{best}$  by using Algorithm 3. Finally, we use the ISPS framework presented in Chapter 3 to calculate the latency of an application, buffer sizes and code size when the unfolded graphs are scheduled by FFD-EDF on the corresponding platform. We perform the same experiment by running the related algorithm proposed in [ZBS13] and using the ISPS framework in Chapter 3 for each  $(app, m, \rho)$ . The obtained ratios for the total code size, total buffer sizes, and latency between the approach in [ZBS13] and our approach are given in Figure 4.5(a), Figure 4.5(b), and Figure 4.5(c), respectively. We can see that by using Algorithm 4 we can achieve up to 17.85 times smaller code size (see Figure 4.5(a),  $\rho=0.95, m=32$ ), up to 24.4 times smaller buffers (see Fig. 4.5(b),  $\rho=0.95, m=32$ ) and up to 11.47 shorter latency (see Figure 4.5(c),  $\rho=0.95, m=32$ ) than the approach in [ZBS13]. Note that both approaches meet



(a) Code size ratio (higher is better)



(b) Buffer size ratio (higher is better)



(c) Latency ratio (higher is better)

**Figure 4.5:** Results of performance evaluation of our proposed approach in comparison to the approach in [ZBS13].

the same throughput requirements. Regarding the buffer sizes, we obtain in 5 experiments out of 308 experiments larger buffer sizes by up to 1.16 times than the approach in [ZBS13] (see for example Figure 4.5(b),  $\rho=0.85$ ,  $m=64$ ). However, for all these experiments we do less unfolding, so we obtain smaller code size. In the case of latency, in 2 experiments out of 308 we get latency which is by 2% larger than the corresponding latency when the approach in [ZBS13] is applied (see Figure 4.5(c),  $m=8$ ). However, in these two cases, we obtain smaller code size and smaller buffer sizes than the approach in [ZBS13].

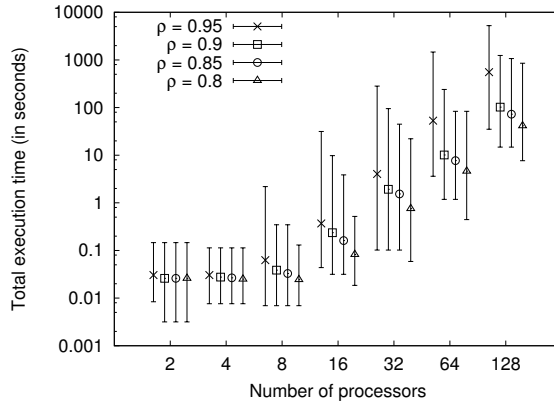
### 4.7.3 Time Complexity of Algorithm 4

We evaluate the efficiency of our algorithm for finding proper values of unfolding factors in terms of the execution time of our Algorithm 4 to find a solution. The execution times for different quality factors and different number of processors in a platform are given in Figure 4.6(a). We compare these execution times with the corresponding execution times of the related approach in [ZBS13]. The comparison is given in Figure 4.6(b).

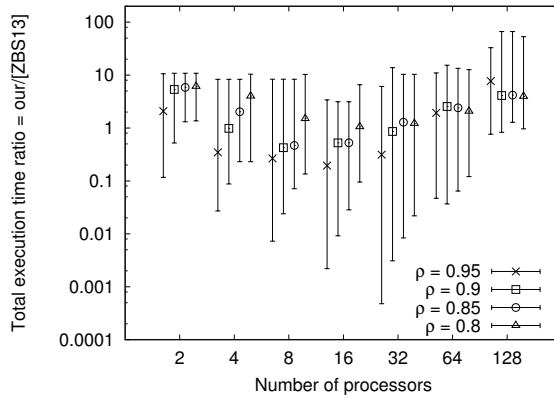
As can be seen from Figure 4.6(a), for platforms containing up to 16 processors, our Algorithm 4 takes in the worst case 32 seconds to find a solution, and less than 1 second on average for all values of quality factor  $\rho$ . For a platform with 32 processors, the execution time of our algorithm is 5 minutes in the worst case, and up to 4 seconds on average. In the case of a 64-processor platform our algorithm needs 25 minutes in the worst case to find a solution, and up to 53 seconds on average. Finally, for a platform with 128 processors Algorithm 4 takes 88 minutes in the worst case and up to 9 minutes on average to find a solution. In addition, it can be seen in Figure 4.6(b) that our approach is on average up to 8 times slower than the approach in [ZBS13] which is acceptable given that our approach delivers solutions of better quality, as shown in Section 4.7.2, within a matter of minutes.

## 4.8 Discussion

As a solution to a problem of exploiting the right amount of parallelism with the aim to achieve the maximum achievable throughput when mapping a streaming application modeled by an SDF graph on a resource-constrained platform under hard real-time scheduling, we presented in this chapter a new unfolding graph transformation and an algorithm which uses the transformation to adapt the parallelism in the application when mapping the application on the platform. Experiments on a set of real-life streaming applications



(a) Running time (in seconds) of Algorithm 4



(b) Execution time ratio (lower is better)

**Figure 4.6:** Results of time evaluation of our proposed approach in comparison to the approach in [ZBS13]

demonstrate that: 1) our unfolding transformation gives shorter latency and smaller buffer sizes when compared to the related approaches; and 2) our algorithm finds, in a matter of minutes, a solution with smaller code size, smaller buffer sizes and shorter latency in 98% of the experiments, while meeting the same performance and timing requirements when compared to an existing approach.

