

Buffer Merging — A Powerful Technique for Reducing Memory Requirements of Synchronous Dataflow Specifications

Praveen K. Murthy
Fujitsu Labs of America, Sunnyvale California
pmurthy@fla.fujitsu.com

Shuvra S. Bhattacharyya
University of Maryland, College Park
ssb@eng.umd.edu

Categories and subject descriptors: D.3.2 [**Language Classifications**]: Dataflow languages; D.3.4 [**Programming Languages**]: Processors—code generation, compilers, optimization; D.2.2 [**Tools and Techniques**]: Structured programming, Petri nets

General Terms: Design Methodology, DSP and Embedded Systems, Block diagram compiler, Dataflow

Additional Keywords: Synchronous dataflow, memory optimization, lifetime analysis, graph coloring, buffer overlaying, path covering, array lifetime

Abstract

We develop a new technique called buffer merging for reducing memory requirements of synchronous dataflow (SDF) specifications. SDF has proven to be an attractive model for specifying DSP systems, and is used in many commercial tools like System Canvas, SPW, and Cocentric. Good synthesis from an SDF specification depends crucially on scheduling, and memory is an important metric for generating efficient schedules. Previous techniques on memory minimization have either not considered buffer sharing at all, or have done so at a fairly coarse level (the meaning of this will be made more precise in the paper). In this paper, we develop a buffer overlaying strategy that works at the level of an input/output edge pair of an actor. It works by algebraically encapsulating the lifetimes of the tokens on the input/output edge pair, and determines the maximum amount of the input buffer space that can be reused by the output. We develop the mathematical basis for performing merging operations, and develop several algorithms and heuristics for using the merging technique

for generating efficient implementations. We show improvements of up to 48% over previous techniques.

1. INTRODUCTION

Memory is an important metric for generating efficient code for DSPs used in embedded applications. This is because most DSPs have very limited amounts of on-chip memory, and adding off-chip memory is frequently not a viable option due to the speed, power, and cost penalty this entails. High-level language compilers, like C compilers have been ineffective for generating good DSP code [23]; this is why most DSPs are still programmed manually in assembly language. However, this is a tedious, error-prone task at best, and the increasing complexity of the systems being implemented, with shorter design cycles, will require design development from a higher level of abstraction.

One potential approach is to do software synthesis from block-diagram languages. Block diagram environments for DSPs have proliferated recently, with industrial tools like System Canvas from Angeles Design Systems [17], SPW from Cadence design systems, the ADS tool from Hewlett Packard, and the Cocentric system studio [5] from Synopsys, and academic tools like Ptolemy [4] from UC Berkeley, and GRAPE from K. U. Leuven [8]. Reasons for their popularity include ease-of-use, intuitive semantics, modularity, and strong formal properties of the underlying dataflow models.

Most block diagram environments for DSPs that allow software synthesis, use the technique of threading for constructing software implementations. In this method, the block diagram is scheduled first. Then the code-generator steps through the schedule, and pieces together code for each actor that appears in the schedule by taking it from a predefined library. The code generator also performs memory allocation, and expands the macros for memory references in the generated code.

Clearly, the quality of the code will be heavily dependent on the schedule used. Hence, we consider in this paper scheduling strategies for minimizing memory usage. Since the scheduling techniques we develop operate on the coarse-grain, system level description, these techniques are somewhat orthogonal to the optimizations that might be employed by tools lower in the flow. For example, a general purpose compiler cannot make usually use of the global control and dataflow that our scheduler can exploit. Thus, the techniques we develop in this paper are complimentary to the work being done on developing better procedural language compilers for DSPs [10][11]. Since the individual actors are programmed in procedural languages like 'C', the output of our SDF compiler is sent to a procedural language compiler to optimize the internals of each actor, and to possibly further optimize the code at a global level (for example, by performing global register allocation.) In particular, the techniques we develop operate on the graphs at a high enough level that particular architectural features of the target processor are largely irrelevant.

2. PROBLEM STATEMENT AND ORGANIZATION OF THE PAPER

The specific problem addressed by this paper is the following. Given a schedule for an SDF graph, there are several strategies that can be used for implementing the buffers needed on the edges of the graph. Previous work on minimizing these buffer sizes has used two models: implementing each buffer separately (for example, in [1][2][21]), or using lifetime analysis techniques for sharing buffers (for example, in [6][16][20]). In this paper, we present a third strategy—buffer merging. This strategy allows sharing of input and output buffers systematically, something that the lifetime-based approaches of [16][20] are unable to do, and something that the separate buffer approaches of [1][2][21] do not even attempt to do. The reason that lifetime-based approaches break down when input/output edges are considered is because they make the conservative assumption that an output buffer becomes live as soon as an actor begins firing, and that an input buffer does not die until the actor has finished execution. Hence, the lifetimes of the input and output buffers overlap, and they cannot be shared. However, as we will show in this paper, relaxing this assumption by analyzing the production and consumption pattern of individual tokens results in significant reuse opportunities that can be efficiently exploited. However, the merging approach of this paper is complimentary to lifetime-based approaches because the merging technique is not able to exploit global sharing opportunities based on the topology of the graph and the schedule. It can only exploit sharing opportunities at the input/output level, based on a fine-grained analysis of token traffic during a single, atomic execution of an actor. Thus, we give a hybrid algorithm that combines both of these techniques and show that dramatic reductions in memory usage are possible compared to either technique used by itself.

In a synthesis tool called ATOMIUM, De Greef, Catthoor, and De Man have developed lifetime analysis and memory allocation techniques for single-assignment, static control-flow specifications that involve explicit looping constructs, such as for loops [6]. While the techniques in [6] are able to reuse and overlay variables and arrays very effectively, the worst case complexity of the algorithms used is exponential. Our algorithms used in this paper, in contrast, provably run in polynomial-time because we are able to exploit the particular, restricted structure of SDF programs and single-appearance schedules for these programs. The algorithms we develop are also purely graph-theoretic techniques, and do not use ILP formulations or array subscript analysis, problems that can have prohibitive complexity in general. It is important to emphasize this last point about complexity: single appearance schedules for SDF graphs allow the formulation of this fine-grained sharing to be performed efficiently whereas a general, non-single appearance, non-SDF model generally induces formulations having non-polynomial complexity. Perhaps this is one reason why previous work for SDF has not considered buffer sharing at the finer level as we do in this paper; we use the insights gained in previous work on single appearance scheduling to derive efficient formulations for buffer sharing at a more fine-grained level. Finally, block-diagram based languages are being used extensively in industry, as has been pointed out earlier, and all of these tools use the SDF model of computation (or a close variant of it). Hence techniques that directly apply to SDF are necessary and useful.

The CBP parameter that we develop in section 5.1, and more completely in [3], plays a role that is somewhat similar to the array index distances derived in the in-place memory

management strategies of Cathedral [22], which applies to nested loop constructs in Silage. The merging approach presented in this paper is different from the approach of [22] in that it is specifically targeted to the high regularity and modularity present in single appearance schedule implementations (at the expense of decreased generality). In particular, the CBP-based overlapping of SDF input/output buffers by shifting actor read and write pointers does not emerge in any straightforward way from the more general techniques developed in [22]. Our form of buffer merging is especially well-suited for incorporation with the SDF vectorization techniques (for minimizing context-switch overhead) developed at the Aachen University of Technology [19] since the absence of nested loops in the vectorized schedules allows for more flexible merging of input/output buffers.

Ritz et. al. [20] give an enumerative method for reducing buffer memory in SDF graphs. Their approach operates only on flat single appearance schedules since buffer memory reduction is tertiary to their goal of reducing code size and context-switch overhead (for which flat schedules are better). However, it's been shown in [2] that on practical applications, their method yields buffering requirements that can be much larger than using nested schedules with each buffer implemented separately.

In [21], Sung et. al. explore an optimization technique that combines procedure calls with inline code for single appearance schedules; this is beneficial whenever the graph has many different instantiations of the same basic actor. Thus, using parametrized procedure calls enables efficient code sharing and reduces code size even further. Clearly, all of the scheduling techniques mentioned in this paper can use this code-sharing technique also, and our work is complementary to this optimization.

This paper is organized as follows. In Sections 3 and 4, we define the notation and review the dataflow scheduling background on which the merging algorithms are based. In Section 5, we develop the theory of merging an input/output buffer pair, and show how the combined buffer size can be accurately calculated. We then show how chains of buffers can be merged. In Section 6, we develop a scheduling algorithm that determines schedules optimized for merge buffer usage for chains of SDF actors. In Section 7, we develop two scheduling algorithms for general acyclic SDF graphs; these scheduling algorithms both generate schedules optimized for merged buffer usage. One of the algorithms is based strictly on the merging technique, while the other is a hybrid algorithm that uses merging and lifetime analysis in alternating phases to get the benefit of both techniques. Section 8 presents the experimental results of running these algorithms on several practical SDF systems, and we conclude in Section 9. We omit all proofs and supporting mathematical machinery here in the interest of readability and space constraints, and refer the reader to [15] for these details.

3. NOTATION AND BACKGROUND

Dataflow is a natural model of computation to use as the underlying model for a block-diagram language for designing DSP systems. The blocks in the language correspond to actors in a dataflow graph, and the connections correspond to directed edges between the actors. These edges not only represent communication channels, conceptually implemented as FIFO queues, but also establish precedence constraints. An actor fires in a dataflow graph by removing tokens from its input edges and producing tokens on its output

edges. The stream of tokens produced this way corresponds naturally to a discrete time signal in a DSP system. In this paper, we consider a subset of dataflow called synchronous dataflow (SDF) [9]. In SDF, each actor produces and consumes a fixed number of tokens, and these numbers are known at compile time. In addition, each edge has a fixed initial number of tokens, called delays.

Fig. 1(a) shows a simple SDF graph. Each edge is annotated with the number of tokens produced (consumed) by its source (sink) actor. Given an SDF edge e , we denote the source actor, sink actor, and delay (initial tokens) of e by $src(e)$, $snk(e)$, and $del(e)$. Also, $prd(e)$ and $cns(e)$ denote the number of tokens produced onto e by $src(e)$ and consumed from e by $snk(e)$. If $prd(e) = cns(e) = 1$ for all edges e , the graph is called **homogenous**. In general, each edge has a FIFO buffer; the number of tokens in this buffer defines the state of the edge. Initial tokens on an edge are just initial tokens in the buffer. The size of this buffer can be determined at compile time, as shown below. The state of the graph is defined by the states of all edges.

A **schedule** is a sequence of actor firings. We compile an SDF graph by first constructing a **valid schedule** — a finite schedule that fires each actor at least once, does not deadlock, and produces no net change in the number of tokens queued on each edge (i.e. returns the graph to its initial state). We represent the minimum number of times each actor must be fired in a valid schedule by a vector q_G , indexed by the actors in G (we often suppress the subscript if G is understood). These minimum numbers of firings can be derived by finding the minimum positive integer solution to the **balance equations** for G , which specify that q must satisfy $prd(e)q(src(e)) = cns(e)q(snk(e))$, for all edges e in G .

The vector q , when it exists, is called the **repetitions vector** of G , and can be computed efficiently [2].

4. CONSTRUCTING MEMORY-EFFICIENT LOOP STRUCTURES

In [2], the concept and motivation behind **single appearance schedules (SAS)** has been defined and shown to yield an optimally compact inline implementation of an SDF graph with regard to code size (neglecting the code size overhead associated with the loop control). An SAS is one where each actor appears only once when loop notation is used. Figure 1 shows an SDF graph, and valid schedules for it. The notation $2B$ represents the firing sequence BB . Similarly, $2(B(2C))$ represents the schedule loop with firing sequence $BCCBCC$. We say that the **iteration count** of this loop is 2, and the body of this loop is $B(2C)$. Schedules 2 and 3 in figure 1 are single appearance schedules since actors A, B, C appear only once. An SAS like the third one in Figure 1(b) is called **flat** since it does not have any nested loops. In general, there can be exponentially many ways of nesting loops in a flat SAS.

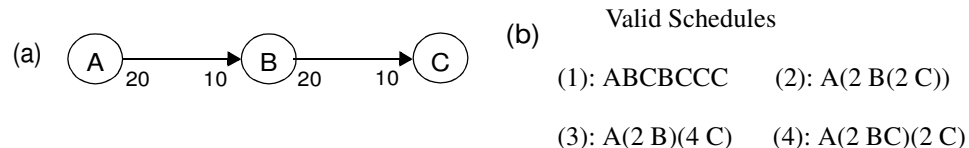


FIGURE 1. An example used to illustrate the interaction between scheduling SDF graphs and the memory requirements of the generated code.

Scheduling can also have a significant impact on the amount of memory required to implement the buffers on the edges in an SDF graph. For example, in Figure 1(b), the buffering requirements for the four schedules, assuming that one separate buffer is implemented for each edge, are 50, 40, 60, and 50 respectively.

4.1 OPTIMIZING FOR BUFFER MEMORY

We give priority to code-size minimization over buffer memory minimization; the importance of addressing this prioritization is explained in [2][14]. Hence, the problem we tackle is one of finding buffer-memory-optimal SAS, since this will give us the best schedule in terms of buffer-memory consumption amongst the schedules that have minimum code size. Following [2] and [14], we also concentrate on acyclic SDF graphs since algorithms for acyclic graphs can be used in the general SAS framework developed in [2] for SDF graphs that are not necessarily acyclic.

For an acyclic SDF graph, any topological sort $a b c \dots$ immediately leads to a valid flat SAS given by $(q(a)a) (q(b)b) \dots$. Each such flat SAS leads to a set of SASs corresponding to different nesting orders.

In [14] and [2], we define the buffering cost as the *sum* of the buffer sizes on each edge, assuming that each buffer is implemented *separately*, without any sharing. With this cost function, we give a post-processing algorithm called dynamic programming post optimization (**DPPO**) that organizes a buffer-optimal nested looped schedule for any given flat SAS. We also develop two heuristics for generating good topological orderings, called APGAN and RPMC.

In this paper, we use an alternative cost for implementing buffers. Our cost is based on overlaying buffers so that spaces can be re-used when the data is no longer needed. This technique is called **buffer merging**, since, as we will show, merging an input buffer with an output buffer will result in significantly less space required than their sums.

5. MERGING AN INPUT/OUTPUT BUFFER PAIR

Example 1: Consider the second schedule in figure 1(b). If each buffer is implemented separately for this schedule, the required buffers on edges AB and BC will be of sizes 20 and 20, giving a total requirement of 40. Suppose, however, that it is known that B consumes its 10 tokens per firing *before* it writes any of the 20 tokens. Then, when B fires for the first time, it will read 10 tokens from the buffer on AB , leaving 10 tokens there. Now it will write 20 tokens. At this point, there are 30 live tokens. If we continue observing the token traffic as this schedule evolves, it will be seen that 30 is the maximum number that are live at any given time. Hence, we see that in reality, we only need a buffer of size 30 to implement AB and BC . Indeed, the diagram shown in figure 2 shows how the read and write pointers for actor B would be overlaid, with the pointers moving right as tokens are

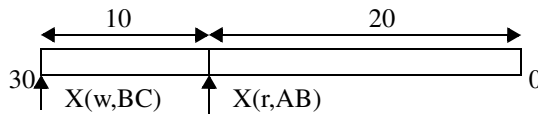


FIGURE 2. The merged buffer for implementing edges AB and BC in figure 1.

read and written. As can be seen, the write pointer, $X(w,BC)$ never overtakes the read pointer $X(r,AB)$, and the size of 30 suffices. Hence, we have merged the input buffer (of size 20) with the output buffer (of size 20) by overlapping a certain amount that is not needed because of the lifetimes of the tokens.

In order to merge buffers in this manner systematically, we introduce several new concepts, notation, and theorems. We assume for the rest of the paper that our SDF graphs are delayless because initial tokens on edges may have lifetimes much greater than tokens that are produced and consumed during the schedule, thus rendering the merging incorrect. This is not a big restriction, since if there are delays on edges, we can divide the graph into regions that are delayless, apply the merging techniques in those portions, and allocate the edges with delays separately. In practical systems, the number of edges having initial tokens is usually a small percentage of the total number of edges, especially in acyclic SDF systems or subsystems. We could even use retiming techniques to move delays around and try to concentrate them on a few edges so that the delayless region becomes as big as possible. Retiming to concentrate delays in this manner has been studied in a different context in [24]. The objective there is to facilitate more extensive vectorization of the input SDF graph. For example, FIGURE 3. (a), the SDF graph has two edges marked with one initial token (delay) each: edges AE and AB. This will be allocated separately and only the other edges will be considered for merging and sharing. In FIGURE 3.(b), the SDF graph can be retimed by firing a preamble schedule BCD to move the delays onto one edge, and thus allow more delay-free edges to be shared. We do not consider these techniques in this paper.

5.1 THE CBP PARAMETER

We define a parameter called the **consumed-before-produced (CBP)** value; this parameter is a property of the SDF actor and a particular input/output edge pair of that actor [3]. Informally, it gives the best known lower bound on the difference between the number of tokens consumed and number of tokens produced over the entire time that the actor is in the process of firing. Formally, let X be an SDF actor, let e_i be an input edge of X , and e_o be an output edge of X . Let the firing of X begin at time 0 and end at time T . Define $c(t)$ ($p(t)$) to be the number of tokens that have been consumed (produced) from (on) e_i (e_o) by time $t \in [0, T]$. The quantities $c(t)$ and $p(t)$ are monotonically non decreasing functions that increase from 0 to $cons(e_i)$ and $prd(e_o)$ respectively. Then, we define

$$CBP(X, e_i, e_o) = \text{MIN}_t \{c(t) - p(t)\}. \quad (\text{EQ 1})$$

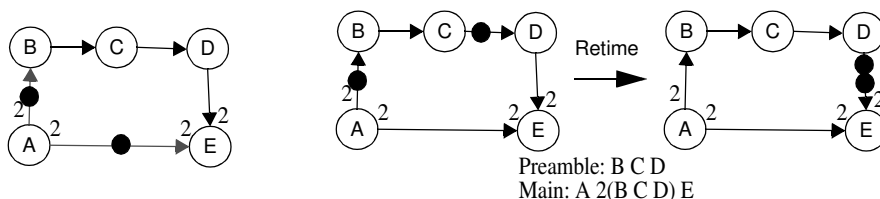


FIGURE 3. (a) An SDF graph with delays and edges that will be allocated separately and not merged. (b) Retiming can reduce number of separate edges, but requires a preamble.

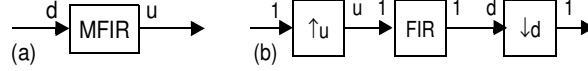


FIGURE 4. Polyphase FIR filter. The structure in (a) implements the graph (b) efficiently.

Note that at $t = 0$, nothing has been consumed or produced, so $c(0) - p(0) = 0$. At T , $p = \text{prd}(e_o)$ tokens have been produced and $c = \text{cns}(e_i)$ tokens have been consumed; hence, $c(T) - p(T) = c - p$. So, we immediately have

$$-p \leq \text{CBP}(X, e_i, e_o) \leq \text{MIN}(0, c - p). \quad (\text{EQ } 2)$$

There are several ways in which the CBP parameter could be determined. The simplest would be for the programmer of the actor to state it based on analyzing the written code inside the actor. This analysis is quite simple in many cases that occur commonly; a study of this type of analysis is reported in [3]. Automatic deduction by source code analysis could also be done, but is beyond the scope of this paper. An analysis of optimized assembly language code written for the polyphase FIR filter in the Ptolemy library (figure 4) shows that [3]

$$\text{CBP}(\text{MFIR}) = \begin{cases} 0 & \text{if } (u \leq d) \\ (d - u) & \text{if } (u > d) \end{cases}. \quad (\text{EQ } 3)$$

Such filters are commonly used in DSP and communication systems. For small, homogeneous actors like adders and multipliers, $\text{CBP} = 0$. If it is not possible to determine a good (meaning largest) lower bound for the CBP parameter, then the worst-case bound of $-p$ is assumed. As we will show, better bounds for CBP will enable smaller merged buffers, but even with the most pessimistic estimate, we can achieve smaller buffers with the merging technique; hence, the merging technique does not necessarily depend on good estimates for the CBP parameter, but certainly benefits from it if available.

5.2 R-SCHEDULES AND THE SCHEDULE TREE

As shown in [14], it is always possible to represent any single appearance schedule for an acyclic graph as

$$(i_L S_L)(i_R S_R), \quad (\text{EQ } 4)$$

where S_L and S_R are SASs for the subgraph consisting of the actors in S_L and in S_R , and i_L and i_R are iteration counts for iterating these schedules. In other words, the graph can be partitioned into a left subset and a right subset so that the schedule for the graph can be represented as in equation 4. SASs having this form at all levels of the loop hierarchy are called R-schedules [14].

Given an R-schedule, we can represent it naturally as a binary tree. The internal nodes of this tree will contain the iteration count of the subschedule rooted at that node. Subschedules S_L or S_R that only have one actor become leaf nodes containing the iteration count and the actor. Figure 5 shows schedule trees for the SAS in figure 1. Note that a schedule tree is not unique since if there are loop factors of 1, then the split into left and right subgraphs can be made at multiple places. In figure 5, the schedule tree for the flat SAS in fig-

ure 1(b)(3) is based on the split $\{A\}\{B, C\}$. However, we could also take the split to be $\{A, B\}\{C\}$. As we will show below, the cost function will not be sensitive to which split is used as they both represent the same schedule.

Define $lf(v)$ to be the iteration count of the node v in the schedule tree. If v is a node of the schedule tree, then $subtree(v)$ is the (sub)tree rooted at node v . If T is a subtree, define $root(T)$ to be the root node of T . A subtree S is a **subset** of a subtree T , $S \subseteq T$ if there is a node v in T such that $S = subtree(v)$. A subtree S is a **strict subset** of a subtree T , $S \subset T$ if there is a node $v \neq root(T)$ in T such that $S = subtree(v)$.

Consider a pair of input/output edges e_i, e_o for an actor Y . Let $X = src(e_i)$, $Z = snk(e_o)$, $snk(e_i) = Y = src(e_o)$. Let T_{XYZ} be the smallest subtree of the schedule tree that contains the actors X, Y, Z . Similarly, let T_{XYZ} be the largest subtree of T_{XYZ} containing actors X, Y , but not containing Z . In figure 5, T_{ABC} is the entire tree, and $T_{A'BC}$ is the tree rooted at the node marked T_{BC} . Largest simply means the following: for every tree $T \subseteq T_{XYZ}$ that contains X, Y and not Z , $T_{XYZ} \supseteq T$. Smallest is defined similarly. Let G be an SDF graph, S be an SAS, and $T(G, S)$ be the schedule tree representing S .

Definition 1: The edge pair $\{e_i, e_o\}$ is said to be **output dominant (OD)** with respect to $T(G, S)$ if $T_{XYZ} \subset T_{XYZ}$ (note that \subset denotes the strict subset).

Definition 2: The edge pair $\{e_i, e_o\}$ is said to be **input dominant (ID)** with respect to $T(G, S)$ if $T_{XYZ} \subset T_{XYZ}$.

The edge pair $\{AB, BC\}$ is ID with respect to both the schedule trees depicted in figure 5. Intuitively, an OD edge pair results from X, Y being more deeply nested together in the SAS than Z .

Fact 1: For any edge pair $\{e_i, e_o\}$, and actors X, Y, Z as defined above, $\{e_i, e_o\}$ is either OD or ID with respect to $T(G, S)$.

Definition 3: For an OD (ID) edge pair $\{e_i, e_o\}$, and actor $snk(e_i) = Y = src(e_o)$, let I_1 be the product of the loop factors in all nodes on the path from the leaf node containing Y to the root node of T_{XYZ} (T_{XYZ}). I_1 is simply the total number of invocations of Y in the largest subschedule not containing Z (X). Similarly, let I_2 be the product of all the loop factors on the path from the leaf node containing Y to the root node of the subtree $T_{XY} \subseteq T_{XYZ}$ ($T_{YZ} \subseteq T_{XYZ}$), where T_{XY} (T_{YZ}) is taken to be the largest tree containing Y but not X (Z).

Note that $I_1 = lf(root(T))I_2$ where $T = T_{XYZ}$ for OD edges pairs and $T = T_{XYZ}$ for ID edge pairs.



FIGURE 5. Schedule trees for schedules in figure 1(b)(2) and (3)

In figure 5, for the schedule tree on the left, we have $I_1 = 2, I_2 = 1$ for B , and $I_1 = 2, I_2 = 2$ for B in the tree on the right.

5.3 BUFFER MERGING FORMULAE

5.3.1 Merging an input/output buffer pair

Given these definitions, we can prove the following theorem about the size of the merged buffer.

Theorem 1: [15] *Let an input-output edge pair $\{e_i, e_o\}$ of an actor Y , and a SAS S for the SDF graph G be given. Define $p = \text{prd}(e_o)$ and $c = \text{cns}(e_i)$. The total size of the buffer required to implement this edge pair is given by the following table:*

TABLE 1. Size of the merged buffer

	OD	ID
$c - p < 0$	$I_1 p + c - p + CBP $	$I_1 c + I_2(p - c) + c - p + CBP $
$c - p \geq 0$	$I_1 p + I_2(c - p) + CBP $	$I_1 c + CBP $

If the edge pair can be regarded as either OD or ID (this happens if $I_1 = I_2$), then the expressions in the 3rd column equal those in the 2nd column. Similarly, if $c = p$, then the expressions in the second row coincide with the expressions in the 3rd row. This verifies our assertion that it does not matter where the split is taken in the SAS when there are multiple choices. Note also that better lower bounds for the CBP make it less negative, reducing $|CBP|$, and thus the size of the merged buffer.

Lemma 1: [15] The size of the merged buffer is no greater than the sum of the buffer sizes implemented separately.

Note that the above lemma holds even with the worst case CBP estimate, where it is not known at all, and we have $|CBP| = p$. Hence, this confirms our assertion that while good CBP estimates benefit the buffer merging technique, the technique is useful even without these CBP estimates. For example, consider the graph and schedule in Figure 6. The left-most part shows how the buffers would be shared using lifetime analysis; in this case only the buffers AB and CD would be shared as their lifetimes are disjoint. Systematic techniques of optimizing the schedule for lifetime analysis, and performing the lifetime analysis and allocation efficiently, are presented in [16]. The middle part of the figure shows how the three buffers are overlaid to yield a total cost of 35, 10 less than the cost achieved via lifetime analysis. We only show the middle part of the full schedule: $5A\ 3B\ 5C\ 3D$. This merge assumes that B has a CBP of 0, meaning that it consumes all of its tokens before it writes any. The right part of the figure shows the case where the CBP is not known for B; it has to be assumed to be -5 for the worst case where B writes all 5 tokens before consuming any. In this case, the merged cost becomes 40 since an extra allowance has to be made for the 5 tokens that B will write before it consumes any. However, this merged cost is still less than that achieved via lifetime analysis; hence, the merging model is useful even with pessimistic estimates of CBP. Basically, the CBP captures token traffic for one firing, whereas the merging model captures the token traffic over several firings in a nested loop. Hence, when there are multiple iterations, the CBP can at worst offset the

buffer by the number of tokens produced on one firing; the other firings can still share the buffer. This figure also illustrates how a fine-grained buffering model leads to lifetime profiles that are non-rectangular and jagged. While one can use packing techniques to model these non-rectangular shapes, as done in [6], packing is NP-complete even when the shapes are rectangular [7]. In contrast, our model is algebraic and efficient, and overlays the buffers without needing to model the jagged shapes explicitly and without relying on heuristics for packing such shapes. We are able to this because we exploit the compact representation of SAS for SDF graphs; such an approach may not be possible in a more general model that may not even have static schedules, let alone a compact one such as an SAS.

Observation 1: For the MFIR of fig. 4, table 1 becomes

TABLE 2. Merged buffer size for MFIR

	OD	ID
$c - p < 0$	$I_1 p$	$I_1 c + I_2(p - c)$
$c - p \geq 0$	$I_1 p + I_2(c - p)$	$I_1 c$

Observation 2: For an input-output edge pair $\{e_i, e_o\}$ of an actor Y with $cons(e_i) = c = p = prd(e_o)$, and $CBP = 0$, table 1 simplifies to $I_1 p = I_1 c$ for all cases.

Homogenous actors are a common special case where observation 2 holds. In the remainder of the paper, we will assume that for illustrative examples, the CBP is equal to the upper bound in equation 2 unless otherwise specified; we do this for clarity of exposition since it avoids having to also list the CBP value for each input/output pair of edges. This means that the size of the merged buffer will be assumed to be taken from table 2 for all the illustrative examples we use (since for the MFIR, the CBP is equal to the upper bound in equation 2), unless specified otherwise. Note that none of our results are affected by this assumption; the assumption only applies to examples we use for illustrative purposes.

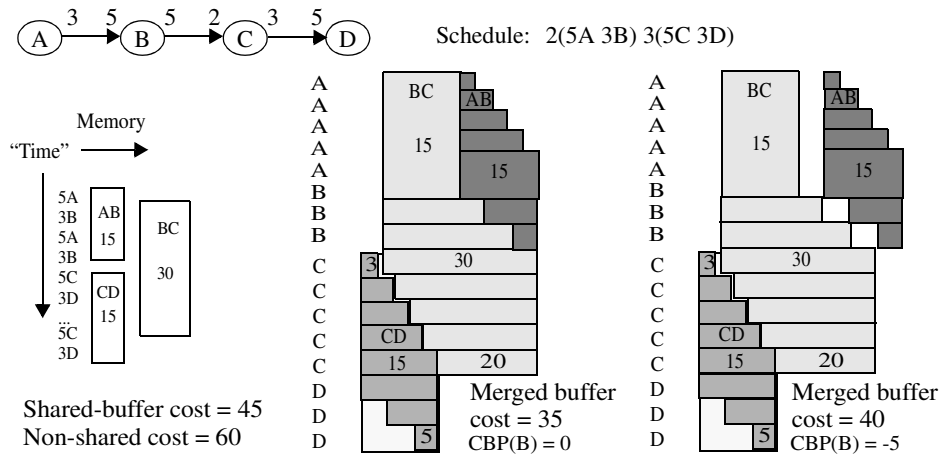


FIGURE 6. Merged buffers versus lifetime analysis based sharing.

5.3.2 Merging a chain of buffers

Let $b_i \oplus b_o$ denote the buffer resulting from merging the buffers b_i and b_o , on edges e_i and e_o respectively. Define $|b|$ to be the size of a buffer b . Define the **augmentation function** $A(b_i \oplus b_o)$ to be the amount by which the output buffer b_o has to be augmented due to the merge $b_i \oplus b_o$. That is,

$$A(b_i \oplus b_o) = |b_i \oplus b_o| - |b_o|. \quad (\text{EQ 5})$$

For OD edge pairs, $|b_o| = I_1 p$ and for ID edge pairs, $|b_o| = I_2 p$. Hence, table 2, can be rewritten in terms of the augmentation as

TABLE 3. Augmentation function for MFIR

	OD	ID
$c - p < 0$	0	$I_1 c - I_2 c$
$c - p \geq 0$	$I_2(c - p)$	$I_1 c - I_2 p$

Lemma 2: [15] The merge operator is associative; i.e, if $v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow v_4$ is a chain of four actors, $e_i = (v_i, v_{i+1}), i = 1, 2, 3$, and b_i are the respective buffers, then $|(b_1 \oplus b_2) \oplus b_3| = |b_1 \oplus (b_2 \oplus b_3)|$.

Theorem 2: [15] Let $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_k, k > 2$, be a path (a chain of actors and edges) in the SDF graph. Let b_i be the buffer on the output edge of actor v_i , and let S be a given SAS (according to which the b_i are determined). Then,

$$|b_1 \oplus \dots \oplus b_{k-1}| = \sum_{i=2}^{k-1} A(b_{i-1} \oplus b_i) + |b_{k-1}| \quad (\text{EQ 6})$$

6. A HEURISTIC FOR MERGED COST-OPTIMAL SAS

Until now, we have assumed that a SAS was given; we computed the merged costs based on this SAS. In this section, we develop an algorithm to generate the SAS so that the merged cost is minimized. In [13], a DPPO formulation is given for chain-structured SDF graphs that organizes the optimal loop hierarchy for any SAS based on the cost function where every buffer is implemented separately. In this section, we give a DPPO formulation that uses the new, buffer merging cost function developed in the previous section for organizing a good loop hierarchy for a chain-structured graph. However, unlike the result in [13], our formulation for this new cost function is not optimal for reasons we will show below; however, it is still a good heuristic technique to use.

6.1 FACTORING

In [14], we show that factoring a SAS by merging loops (in other words, generating nested loops) by the greatest extent possible is not harmful to buffer memory reduction, and that the buffering requirements in a fully factored looped schedule are less than or equal to the requirements in the non-factored loop. Of-course, this result depends on the buffering cost function being used. For example, the result does not, in general, hold under the shared

buffer model used in [16]. Happily, this result does hold for the merging cost function, as shown by the following theorem:

Theorem 3: [15] *Suppose that $S = (i_L S_{ik})(i_R S_{k+1j})$ is a valid SAS for a chain-structured SDF graph G . Define $\text{cost}(S)$ to be the size of the buffer obtained by merging all the buffers on the edges in S . Then, for any positive integer γ that divides i_L and i_R , the schedule*

$$S' = \gamma \left\{ \left(\frac{i_L}{\gamma} S_{ik} \right) \left(\frac{i_R}{\gamma} S_{k+1j} \right) \right\}$$

satisfies $\text{cost}(S') \leq \text{cost}(S)$.

6.2 DPPO FORMULATION

Let $v_i \rightarrow v_{i+1} \rightarrow \dots \rightarrow v_j$ be a sub-chain of actors in the chain-structured SDF graph. The basic idea behind the DPPO formulation is to determine where the split should occur in this chain, so that the SAS S_{ij} for it may be represented as

$$S_{ij} = (i_L S_{ik})(i_R S_{k+1j}).$$

If S_{ik} and S_{k+1j} are known to be optimal for those subchains, then all we have to do to compute S_{ij} is to determine the $i \leq k < j$ where the split should occur; this is done by examining the cost for each of these k . In order for the resulting S_{ij} to be optimal, the problem must have the optimum substructure property: the cost computed at the interfaces (at the split points) should be independent of the schedules S_{ik} and S_{k+1j} . Now, if each buffer is implemented separately, then the cost at the split point is simply the size of the buffer on the edge crossing the split, and this does not depend on what schedule was chosen for the left half (S_{ik}). Hence, the algorithm would be optimal then [13]. However, for the merging cost function, it turns out that the interface cost does depend on what S_{ik} and S_{k+1j} are, and hence this DPPO formulation is not optimal, as shown in [15]. It is a greedy heuristic that attempts to give a good approximation to the minimum. In section 8, we show that on practical SDF systems, this heuristic can give better results than the technique of [13]. In order to compute the interface costs, let the buffers on the edges be $b_i, \dots, b_k, b_{k+1}, \dots, b_{j-1}$. Now suppose that the split occurs at k . That is, actors v_i, \dots, v_k are on the left side of the split. Since we know $\text{cost}(S_{ik})$ and $\text{cost}(S_{k+1j})$ (these are memoized, or stored in the dynamic programming table), we have (by theorem 2) that

$$\text{cost}(S_{ik}) = |b_{k-1}| + A_{ik}, \text{ and } \text{cost}(S_{k+1j}) = |b_{j-1}| + A_{k+1j},$$

where A is the augmentation term. Hence, in order to determine the cost of splitting at k , we have to determine $b_{k-1} \oplus b_k$ and $b_k \oplus b_{k+1}$. Using theorem 2, the total cost is thus given by

$$c_{ij}(k) = \frac{\text{cost}(S_{ik}) - |b_{k-1}| + A(b_{k-1} \oplus b_k) + A(b_k \oplus b_{k+1}) + \text{cost}(S_{k+1j})}{1}. \quad (\text{EQ 7})$$

We then choose the k that minimizes the above cost:

$$\text{cost}(S_{ij}) = \text{MIN}_{i \leq k < j} \{c_{ij}(k)\}.$$

In order to compute $A(b_{k-1} \oplus b_k)$ and $A(b_k \oplus b_{k+1})$, we need the appropriate I_1 and I_2 factors. We can compute I_1, I_2 efficiently, as is shown in [15], and omitted here due to space constraints. The entire DPPO algorithm then has a running time of $O(n^3)$ where n is the number of actors in the chain.

Unfortunately, we cannot prove that the DPPO formulation of section 6.2 is optimal; details can be found in [15].

7. ACYCLIC GRAPHS

In this section, we extend the merging techniques to arbitrary, delayless, acyclic SDF graphs. The techniques we develop here can easily be extended to handle graphs that have delays, as discussed in Section 5. SASs for SDF graphs that contain cycles can be constructed in an efficient and general manner by using the loose interdependence scheduling framework (LISF) [2]. The LISF operates by decomposing the input graph G into a hierarchy of acyclic SDF graphs. Once this hierarchy is constructed, any algorithm for scheduling acyclic SDF graphs can be applied to each acyclic graph in the hierarchy, and the LISF combines the resulting schedules to construct a valid SAS for G . Thus, the LISF provides an efficient mechanism by which the techniques developed here can be applied to general (not necessarily acyclic) topologies.

When acyclic graphs are considered, there are two other dimensions that come into play for designing merging algorithms. The first dimension is the choice of the topological ordering of the actors; each topological ordering leads to a set of SASs. This dimension has been extensively dealt with before in [2], where we devised two heuristic approaches for determining good topological orderings. While these heuristics were optimized for minimizing the buffer memory cost function where each buffer is implemented separately, they can be used with the new merged cost function as well. We leave for future work to design better heuristics for the merged cost function, if it is possible.

The second dimension is unique to the merge cost function, and is the issue of the set of paths that buffers should be merged on. In other words, given a topological sort of the graph, and a nested SAS for this graph, there still remains the issue of what paths buffers should be merged on. For the chain-structured graphs of the previous section, there is only one path, and hence this is not an issue. Since an acyclic graph can have an exponential number of paths, it does become an issue when acyclic graphs are considered. In the following sections, we develop two approaches for determining these paths. The first approach determines the optimum set of paths along which buffers should be merged, and then performs lifetime analysis (using techniques from [16]) once to exploit further sharing opportunities not captured by the merging. Note that buffer merging only captures sharing opportunities among buffers that lie on some path of actors; buffers not related this way still have to be shared by lifetime analysis techniques. The algorithm we give is optimal in the sense that for a given topological ordering and SAS, our algorithm will determine the lowest merge cost implementation when buffers are merged in a linear order along the paths. The second approach is a bottom up approach that combines lifetime analysis techniques and the merging approach in a more fine-grained level, where merging and lifetime-analysis are performed iteratively. However, a drawback of this approach is that it

is of high complexity and is slow. Yet another dimension can be introduced by not merging the buffers linearly; this is captured by clustering as we show later.

7.1 PATH COVERING

Determining the best set of paths to merge buffers on can be formulated as a path covering problem. Essentially, we want a disjoint set of paths Ψ such that each edge in the graph is in exactly one path in Ψ . The total buffering cost is then determined by merging the buffers on the edges in each path, and summing the resulting costs.

Example 2: Consider the graph shown in figure 7. The schedule tree is shown on the right, and represents the SAS $5A \ 2(3(2B \ 3C) \ 2D)$. There are two possible ways of merging buffers in this graph: $b_1 \oplus b_2 \oplus b_3 + b_4$ and $b_1 \oplus b_2 + b_4 \oplus b_3$. These correspond to the paths $\{(AB, BC, CD), (AC)\}$ and $\{(AB, BC), (AC, CD)\}$. The non-merged costs for the buffers in each edge, for the schedule shown, are given by $b_1 = 60$, $b_2 = 6$, $b_3 = 36$, and $b_4 = 90$. Thus, if each of these were to be implemented separately, the total buffering cost would be $60 + 6 + 36 + 90 = 192$. It can be verified that $b_1 \oplus b_2 \oplus b_3 + b_4 = 180$, and $b_1 \oplus b_2 + b_4 \oplus b_3 = 168$. Hence, the better set of paths to use for this example is $\{(AB, BC), (AC, CD)\}$.

Given a directed graph (digraph) G , an **edge-oriented path** is a sequence of edges e_1, e_2, \dots, e_n such that $src(e_i) = snk(e_{i-1})$ for each $i = 2, 3, \dots, n$. A **node-oriented path** is a sequence of nodes v_1, \dots, v_n such that (v_i, v_{i+1}) is an edge in the graph for each $i = 1, \dots, n-1$. The **edge-set of a node-oriented path** v_1, \dots, v_n is the set of edges (v_i, v_{i+1}) . An **edge-oriented path cover** χ is defined as a set of edge-disjoint edge-oriented paths whose union is the entire edge set of G . A **node-oriented path cover** Ψ is defined as a set of node-disjoint node-oriented paths whose union is the entire node set of G . In other words, each node in G appears in exactly one node-oriented path $p \in \Psi$. The **edge-set of a node-oriented path cover** is the union of the edge-sets of each node-oriented path in the cover.

For an SDF graph G , define the **buffer cost of an edge-oriented path** e_1, e_2, \dots, e_n as $|b_1 \oplus \dots \oplus b_n|$, where b_i is the buffer on edge e_i . Define the **buffer cost of an edge-oriented path cover** χ as the sum of the buffer costs of the paths in the cover.

Definition 4: The PATH SELECTION PROBLEM FOR BUFFER MERGING (**PSPBM**) in an acyclic SDF graph is to find an edge-oriented path cover of minimum buffer cost.

In order to solve the path selection problem, we first derive a weighted, directed MERGE GRAPH from the SDF graph $G = (V, E)$. The **MERGE GRAPH** MG is defined as $MG = (V_{MG}, E_{MG}, w)$, where

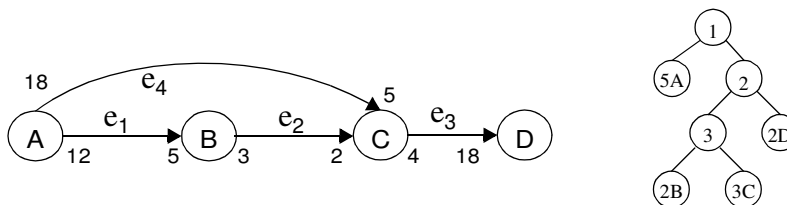


FIGURE 7. An example to show the variation of the merge cost with path selection.

$$\begin{aligned}
V_{MG} &= \{v_e \mid e \in E\} \cup \{s_e \mid e \in E\}, \\
E_{MG} &= \{(v_{e_1}, v_{e_2}) \mid e_1 \in E, e_2 \in E, \text{snk}(e_1) = \text{src}(e_2)\} \cup \{(v_e, s_e) \mid e \in E\}, \\
w((v_{e_1}, v_{e_2})) &= A(b(e_1) \oplus b(e_2)), \text{ and} \\
w((v_e, s_e)) &= |b(e)|.
\end{aligned}$$

The buffer on an edge e in the SDF graph G is denoted by $b(e)$. Figure 8 shows the MERGE GRAPH for the SDF graph in figure 7. The nodes of type s_e are called **S-type nodes**.

Fact 2: If the SDF graph is acyclic, the associated MERGE GRAPH is also acyclic.

Given a weighted digraph G , the **weight of a path** p in G , denoted as $w(p)$ is the sum of the weights on the edges in the path. The **weight of a path cover** Ψ , denoted as $w(\Psi)$ is the sum of the path weights of the paths in Ψ .

We define a **maximal path cover** for the MERGE GRAPH as a node-oriented path cover Ψ such that each path $p \in \Psi$ ends in an S-type node. A **minimum weight maximal path cover (MWMPC)** Ψ^* is a maximal path cover of minimum weight.

Definition 5: The **MWMPC problem** for MERGE GRAPHS is to find an MWMPC.

Given an MWMPC Ψ for the MERGE GRAPH MG , for each path p in the MWMPC, replace each (non S-type) node v_e in p by the corresponding edge e in the SDF graph G to get an edge-oriented path q in G . Let χ be the set of paths q . Note that we do not have any edges corresponding to S-type nodes in MG . Then we have the following obvious result:

Lemma 3: [15] The set of edge-oriented paths χ constructed above is a solution to the PSPBM problem; that is, χ is an edge-oriented path cover of minimum buffer cost for the SDF graph G .

For example, in figure 8, the MWMPC is given by $\{s_{e_1}, (v_{e_1} \rightarrow v_{e_2} \rightarrow s_{e_2}), (v_{e_4} \rightarrow v_{e_3} \rightarrow s_{e_3}), s_{e_4}\}$, and it can be verified easily that this corresponds to the optimal buffer merge paths shown in example 2.

In [12], Moran et al. give a technique for finding maximum weight path covers in digraphs. We modify this technique slightly to give an optimum, polynomial time algorithm for finding an MWMPC in a MERGE GRAPH. Given a weighted, directed acyclic graph $G = (V, E, w)$, with $V = \{v_1, \dots, v_n\}$, we first derive the following weighted bipartite graph $G_B = (X, Y, E_B, w_B)$:

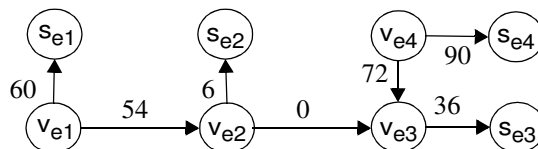


FIGURE 8. The MERGE GRAPH for the SDF graph in figure 7.

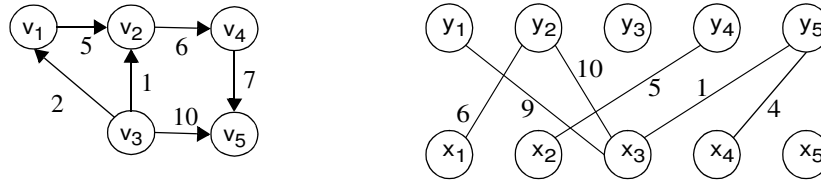


FIGURE 9. The bipartite graph derived from a weighted digraph.

$$X = \{x_i \mid v_i \in V\}, Y = \{y_i \mid v_i \in V\}, E_B = \{(x_i, y_j) \mid (v_i, v_j) \in E\}$$

$$w_B((x_i, y_j)) = W^* - w((v_i, v_j)) \forall (v_i, v_j) \in E, \text{ where}$$

$$W^* = \text{MAX}_{(v_i, v_j) \in E} (w((v_i, v_j))) + 1$$

Figure 9 shows a weighted digraph and the corresponding bipartite graph.

A **matching** M is a set of edges such that no two edges in M share an endpoint. The weight of a matching is the sum of the weights of the edges in the matching. A **maximum weight matching** is a matching of maximum weight. Maximum weight matchings in bipartite graphs can be found in polynomial time ($O(|V_B|^2 \cdot \log(|V_B|))$); for example, using the ‘‘Hungarian’’ algorithm [18].

Given a matching M in G_B , define $E_\Psi = \{(v_i, v_j) \mid (x_i, y_j) \in M\}$. That is,

$$(v_i, v_j) \in E_\Psi \Leftrightarrow (x_i, y_j) \in M \quad (\text{EQ 8})$$

Theorem 4: [15] *If M^* is a maximum weight matching in G_B , and E_Ψ is as defined above, then $\Psi^* = \text{Co}(V, E_\Psi)$ is an MWMPC for the MERGE GRAPH G , where $\text{Co}(V, E_\Psi)$ denotes the connected components of (V, E_Ψ) .*

The algorithm for finding an optimal set of paths along which to merge buffers is summarized in figure 10. Once all the buffers along the paths have been merged, we perform an overall lifetime analysis on the set of merged buffers to exploit any sharing opportunity left between the buffers.

7.1.1 Running time

As already mentioned, the matching step takes $O(|V_b|^2 \cdot \log(|V_b|))$, where V_b is the set of nodes in the bipartite graph. Since $|V_b| = |X| + |Y| = 2 \cdot |V_{MG}|$, where X, Y are the node sets in the bipartite graph, and V_{MG} is the node set in the MERGE GRAPH, we have

Procedure `determineMergePaths(SDF Graph G)`

```

 $G_E \leftarrow \text{MergeGraph}(G)$ 
 $G_B \leftarrow \text{BipartiteGraph}(G_E)$ 
 $M^* \leftarrow \text{maxMatching}(G_B)$ 
 $(v_i, v_j) \in \Psi^* \Leftrightarrow (x_i, y_j) \in M^*$ 
return  $\Psi^*$ 

```

FIGURE 10. Procedure for computing the optimum set of paths along which buffers should be merged.

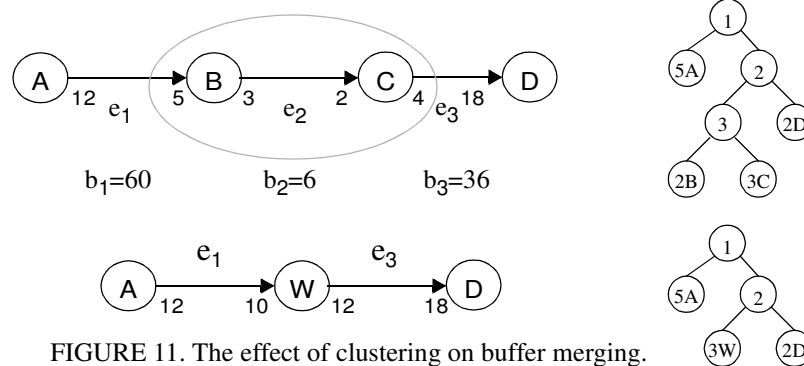


FIGURE 11. The effect of clustering on buffer merging.

$|V_{MG}| = 2 \cdot |E|$ and $|V_b| = 4 \cdot |E| = O(|E|)$, where E is the set of edges in the SDF graph. The construction of the MERGE GRAPH takes time $O(|V| \cdot \log(|V|) + |E|)$ if the SDF graph has actors with constant in-degree and out-degree. The $\log(|V|)$ comes from computing the buffer merges to determine the weights in the MERGE GRAPH; recall that since the SAS is given now, computing a merge requires traversal of the schedule tree, and can be done in time $\log(|V|)$ on average if the tree is balanced. If the tree is not balanced, then the merge computation could take $O(|V|)$ time, meaning that the MERGE GRAPH construction takes $O(|V|^2 + |E|)$ time. The bipartite graph also takes time $O(|E|)$. Hence, the overall running time is dominated by the matching step, and takes time $O(|E|^2 \cdot \log(|E|))$, or $O(|V|^2 \cdot \log(|V|))$ if the SDF graph is sparse.

7.1.2 Clustering

While the buffer merging technique as developed results in significant reductions in memory requirements, even more reduction can be obtained if other optimizations are considered. The first of these is the use of clustering. Until now, we have implicitly assumed that the buffers that are merged along a chain are overlaid in sequence. However, this may be a suboptimal strategy since it may result in a fragmented buffer where lot of storage is wasted. Hence, the optimization is to determine the sub-chains along a chain where buffers should be profitably merged, and not to blindly merge all buffers in a chain. This can be captured via clustering, where the cluster will determine the buffers that are merged. For instance, consider the SDF graph in figure 11. If we merge the buffers in the top graph, we get a merged buffer of size 90. However, if we merge the two edges in the clustered graph at the bottom, where actors B and C have been clustered together into actor W , we get a merged buffer of size 66. The edge between B and C is implemented separately, and it requires 6 storage units. Hence the total buffering cost goes down to 72. The reason that this happens is shown in figure 12. The buffer of size 6 between the two larger buffers fragments the overall buffer and results in some space being wasted. The clustering removes this small buffer and merges only the two larger ones, enabling more efficient use of storage.

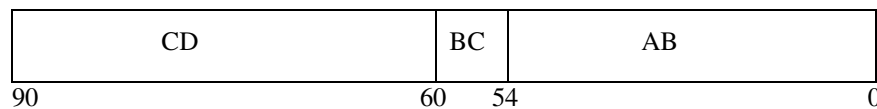


FIGURE 12. Fragmentation due to merging all buffers in series.

The above optimization can be incorporated into the path covering algorithm by introducing transitive edges in the MERGE GRAPH construction. Instead of just having edges between nodes that correspond to input/output edges of the same actor in the SDF graph, we introduce edges between two nodes if they correspond to edges on some directed path in the SDF graph. Figure 13 shows the MERGE GRAPH for the SDF graph in figure 11.

7.2 A BOTTOM-UP APPROACH

Now we describe another technique for determining merge paths that also combines lifetime analysis techniques from [16]. Briefly, the lifetime analysis techniques developed in [16] construct an SAS optimized using a particular shared-buffer model that exploits temporal disjointedness of the buffer lifetimes. The method then constructs an intersection graph that models buffer lifetimes by nodes and edges between nodes if the lifetimes intersect in time. FirstFit allocation heuristics [15] are then used to perform memory allocation on the intersection graph. The shared buffer model used in [16] is useful for modeling the sharing opportunities that are present in the SDF graph as a whole, but is unable to model the sharing opportunities that are present at the input/output buffers of a single actor. The model has to make the conservative assumption that all input buffers are simultaneously live with all output buffers of an actor while the actor has not fired the requisite number of times in the periodic schedule. This means that input/output buffers of a single actor cannot be shared under this model. However, the buffer merging technique developed in this paper models the input/output edge case very well, and is able to exploit the maximum amount of sharing opportunities. However, the merging process is not well suited for exploiting the overall sharing opportunities present in the graph, as that is better modeled by lifetime analysis. Hence, the bottom-up approach we give here combines both these techniques, and allows maximum exploitation of sharing opportunities at both the global level of the overall graph, and the local level of an individual input/output buffer pair of an actor.

The algorithm is stated in figure 14. It basically makes several passes through the graph, each time merging a suitable pair of input/output buffers. For each merge, a global memory allocation is performed using the combined lifetime of the merged buffer. That is, the start time of the merged buffer is the start time of the input buffer, and the end time is the end time of the output buffer (the procedure `changeIntersectionGraph` performs this). If the allocation improves, then the merge is recorded (procedure `recordMerge`). After examining each node and each pair of input/output edge pairs, we determine whether the best recorded merge improved the allocation. If it did, then the merge is performed (procedure `mergeRecorded`), and another pass is made through the graph where every node and its input/output edge pairs is examined. The algorithm stops when there is no further improvement.

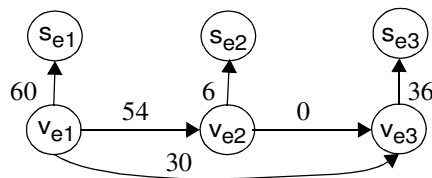


FIGURE 13. Adding transitive edges to the MERGE GRAPH for nonlinear merging.

7.2.1 Running time analysis

The loops labelled (1), (2), and (3) take

$$\sum_{i=1}^{|V|} \text{indeg}(v_i) \cdot \text{outdeg}(v_i) \quad (\text{EQ 9})$$

steps, where $\text{indeg}(v)$ is the in-degree of actor v and $\text{outdeg}(v)$ is the out-degree of actor v . In the worst possible case, we can show that this sum is $O(|V|^3)$, assuming a dense, acyclic graph. If the graph is not dense, and the in- and out- degrees of the actors are bounded by pre-defined constants, as they usually are in most SDF specifications, then equation 9 would be $O(|V|)$. The merging step in line (4) can be precomputed and stored in a matrix since merging a buffer with a chain of merged buffers just involves merging the buffer at the end of the chain and summing the augmentation. This precomputation would store the results in an $|E| \times |E|$ matrix, and would take time $O(|E|^2 \cdot \log |V|)$ in the

Procedure `mergeBottomUp(SDF Graph G , SAS S)`

```

 $I \leftarrow \text{computeIntersectionGraph}(G, S)$ 
 $cur_{best} \leftarrow \text{allocate}(I)$ 
 $iter_{best} \leftarrow cur_{best}$ 

while (true)
  for each node  $v \in G$  (1)
    for each input edge  $e_i$  of  $v$  (2)
      for each output edge  $e_o$  of  $v$  (3)
         $b \leftarrow b_i \oplus b_o$  (4)
         $I \leftarrow \text{changeIntersectionGraph}(S, b)$ 
         $m \leftarrow \text{allocate}(I)$ 
        if ( $m < iter_{best}$ )
          recordMerge( $b, I, m$ )
           $iter_{best} \leftarrow m$ 
        fi
      restore( $I, b_i, b_o, S$ )
    end for
  end for
end for
if ( $iter_{best} < cur_{best}$ )
   $cur_{best} \leftarrow iter_{best}$ 
  mergeRecorded()
else
  break
fi
end while

```

FIGURE 14. A bottom-up approach that combines buffer merging with lifetime analysis.

average case, and $O(|E|^2 \cdot |V|)$ time in the worst case. So line (4) would end up taking a constant amount of time since the precomputation would occur before the loops. The `intersectionGraph` procedure can take $O(|E|^2)$ time in the worst case. While this could be improved by recognizing the incremental change that actually occurs to the lifetimes, it is still hampered by the fact that the actual allocation heuristic still takes time $O(|E|^2)$. The overall while loop can take $O(|E|)$ steps since each edge could end up being merged. Hence, the overall running time, for practical systems, is $O(|V| \cdot |E|^3)$ which is $O(|V|^4)$ for sparse graphs. Improvement, if any, can be achieved by exploring ways of implementing the FirstFit heuristic to work incrementally (so that it does not take $O(|E|^2)$); however, this is unlikely to be possible as the merged buffer will have a different lifetime and size, and the allocation has to be redone from scratch each time.

8. EXPERIMENTAL RESULTS

8.1 CD-DAT

Consider the SDF representation of the CD-DAT sample rate conversion example from [13], shown in figure 15. The best schedule obtained for this graph in [13], using the non-merged buffering model, has a cost of 260. If we take this SAS, and merge the buffers, then the cost goes down to 226. Applying the new DPPO formulation of Section 6.2 based on the merging cost, gives a different SAS, having a merged cost of 205. This represents a reduction of more than 20% from previous techniques.

8.2 HOMOGENOUS SDF GRAPHS

Unlike the techniques in [13][2], the buffer merging technique is useful even if there are no rate changes in the graph. For instance, consider a simple, generic image-processing system implemented using SDF shown in figure 16. This graph has a number of pixelwise operators that can be considered to have a *CBP* of 0 for any input-output edge pair. The graph is homogenous because one token is exchanged on all edges; however, the token can be a large image. Most previous techniques, and indeed many current block-diagram code-generators (SPW, Ptolemy, DSPCanvas) will generate a separate buffer for each edge, requiring storage for 8 image tokens; this is clearly highly wasteful since it can be almost seen by inspection that 3 image buffers would suffice. Our buffer merging technique gives an allocation of 3 buffers as expected. In particular, for the example below, we can choose the path to be from *A* to *Disp* and apply the merge along that path. Applying the bottom-up approach will reduce this further to 2 buffers since the lifetime analysis shows that *C*'s output can reuse the location that *B* used.

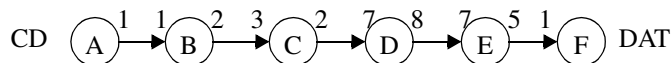


FIGURE 15. The CD-DAT sample rate converter.

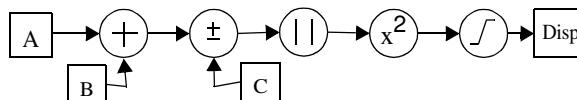


FIGURE 16. An image processing flow with pixelwise operators.

8.3 A NUMBER OF PRACTICAL SYSTEMS

Table 4 shows the results of applying the bottom-up and top-down techniques to a number of practical SDF systems, and compares these techniques to lifetime based approaches. These are multirate filterbank systems (*qmf*), a satellite receiver (satrec) implementation from [20], and a number of other systems taken from the Ptolemy library: a 16 QAM modem (16qamModem), a pulse amplitude modulation system (4pamxmitrec), a phased array receiver (phasedArray), a overlap-add FFT system (overAddFFT), and block vocoder (blockVox). The filterbank examples are denoted using the following notation: “qmf23_3d” means that the system is a complete construction-reconstruction system of depth 3; that is, 8 channels. The “23” denotes that a 1/3-2/3 split is used for the spectrum; that is, the signal is recursively divided by taking 1/3 of the spectrum (low-pass) and passing 2/3 of it as the high-pass component. The “qmf12_xd” denote filter banks where a 1/2-1/2 split is used, and the “qmf235_xd” systems denote filterbanks where a 2/5-3/5 split is used. Figure 17(b) shows the “qmf12_3d” system. The rate-changing actors in this system are polyphase FIR filters for which the CBP parameter is obtained using equation 3. The largest of these examples have around 50 actors. The examples used in these experiments are consistent with their use in [16]. We have also tried to use examples from the public domain such as the Ptolemy library, and previously published examples such as the satellite receiver. Finally, it should be noted that even a moderately sized SDF graph with tens of nodes and lots of rate changes has a huge design space of single appearance schedules, and hence is a good test of these scheduling techniques.

The columns named BU(R) and BU(A) give the results of applying the bottom-up algorithm on topological orderings generated by the RPMC and APGAN heuristics respectively, and the columns labelled TD(A) and TD(R) do the same for the top-down algorithm. The column name “bestShrd” gives the best result of applying the lifetime-analysis algorithms from [16]. The “best NonShrd” column contains the best results obtained under the non-shared buffer models, using the algorithms in [2]. The last column gives the percentage improvement of the best of the TD(R), TD(A), BU(R), BU(A) columns (marked in bold in each of these columns) compared to the “bestShrd” column; that is, the improvement of the combined buffer merging and lifetime analysis approach compared to the pure lifetime approach. As can be seen, the improvements in memory requirements averages 12% over these examples, and is as high as 49% in one case.

The experiment also shows that for this set of examples, the bottom-up approach generally outperforms the top-down approach; however, when the top-down approach does outperform the bottom-up approach, it does so significantly as two of the examples show.

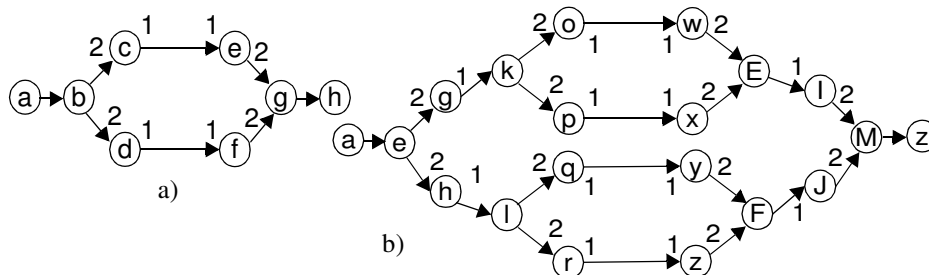


FIGURE 17. SDF graph for a two-sided filterbank. a) Depth 1 filterbank, b) Depth 2 filterbank. The produced/consumed numbers not specified are all unity.

TABLE 4. Buffer sizes for several SDF systems under various scheduling algorithms

System	bestNon Shrd	best Shrd	TD (R)	TD (A)	BU (R)	BU (A)	best Merged	% impr
16qamModem	35	9	15	15	8	8	8	11.1
4pamxmitrec	49	35	35	20	49	18	18	48.6
aqmf12_2d	34	9	13	16	9	13	9	0.0
aqmf12_3d	78	16	29	36	16	27	16	0.0
aqmf235_2d	122	55	63	91	52	65	52	5.5
aqmf235_3d	492	240	289	553	237	367	237	1.3
aqmf23_2d	60	24	26	35	22	26	22	8.3
aqmf23_3d	173	63	74	126	61	78	61	3.2
blockVox	409	135	198	199	132	129	129	4.4
nqmf23	209	132	156	251	129	230	129	2.3
overAddFFT	1222	514	514	386	514	514	386	24.9
phasedArray	2496	2071	1672	1672	1672	1672	1672	19.3
satrec	1542	991	1233	773	1441	972	773	22.0

9. CONCLUSION

Earlier work on SDF buffer optimization has focused on the separate buffer model, and the lifetime model, in which buffers cannot share memory space if any part of the buffers simultaneously contain live data. Our work on buffer merging in this paper has formally introduced a third model of buffer implementation in which input and output buffers can be overlaid in memory if subsets of the buffers have disjoint lifetimes. The technique of buffer merging is able to encapsulate the lifetimes of tokens on edges algebraically, and use that information to develop near-optimal overlaying strategies. While the mathematical sophistication of this technique is especially useful for multirate DSP applications that involve numerous input/output buffer accesses per actor invocation, a side benefit is that it is highly useful for homogenous SDF graphs as well, particularly those involving image and video processing systems since the savings can be dramatic. We have given an analytic framework for performing buffer merging operations, and developed a dynamic programming algorithm that is able to generate loop hierarchies that minimize this merge cost function for chains of actors.

For general acyclic graphs, we have developed two algorithms for determining the optimal set of buffers to merge. The first of these techniques is an innovative formulation using path-covering for determining a provably optimal set of paths (under certain assumptions) on which buffers should be merged. Since this technique is a pure buffer-merging technique, and does not use lifetime analysis, it is faster and might be useful in cases where fast compile times are especially important. The second of these techniques, the bottom-up merging algorithm, combines merging and lifetime analysis. Our experiments show improvements over the separate-buffer and lifetime-based implementations of the best of the merging techniques of 12% on average on a number of practical systems, with some systems exhibiting upto a 49% improvement.

As mentioned before, lifetime-based approaches break down when input/output edges are considered because they make the conservative assumption that an output buffer becomes

live as soon as an actor begins firing, and that an input buffer does not die until the actor has finished execution. This conservative assumption is made in [16] primarily to avoid having to pack arrays with non-rectangular lifetime profiles; if the assumption is relaxed, we would get a jagged, non-rectangular lifetime profile, and this could in theory be packed to yield the same memory consumption requirements as the buffer merging technique. However, packing these non-rectangular patterns efficiently is significantly more difficult (as shown by the exponential worst-case complexity of the techniques in [6]), and moreover, it still does not take into account the very fine-grained production and consumption pattern modeled by the CBP parameter. Hence, the buffer merging technique finesses the problem of packing arrays that have non-rectangular lifetime profiles by providing an exact, algebraic framework that exploits the particular structure of SDF graphs and single appearance looped schedules. This framework can then be used with the lifetime-based approach of [16] efficiently to get significant reductions in buffer memory usage.

Buffer merging does not render separate-buffers or lifetime-based buffer sharing obsolete. Separate buffers are useful for implementing edges that contain delays efficiently. Furthermore, they provide a tractable cost function with which one can rigorously prove useful results on upper bound memory requirements [2]. Lifetime-based sharing is a dual of the merging approach, as mentioned already, and can be fruitfully combined with the merging technique to develop a powerful hybrid approach that is better than either technique used alone, as we have demonstrated with the algorithm of Section 7.2.

10. REFERENCES

- [1] M. Ade, R. Lauwereins, J. Peperstraete, "Implementing DSP Applications on Heterogeneous Targets Using Minimal Size Data Buffers," IEEE Wkshp. on Rapid Sys. Prototyping, 1996.
- [2] S. S. Bhattacharyya, P. K. Murthy, E. A. Lee, *Software Synthesis from Dataflow Graphs*, Kluwer, 1996.
- [3] S. S. Bhattacharyya, P. K. Murthy, "The CBP Parameter—a Useful Annotation for SDF Compilers," Proceedings of the ISCAS, Geneva, Switzerland, May 2000.
- [4] J. Buck, S. Ha, E. A. Lee, D. G. Messerschmitt, "Ptolemy: a Framework for Simulating and Prototyping Heterogeneous Systems," *Intl. J. of Computer Simulation*, Jan. 1995.
- [5] J. Buck, R. Vaidyanathan, "Heterogeneous Modeling and Simulation of Embedded Systems in El Greco," CODES, May 2000.
- [6] E. De Greef, F. Catthoor, H. De Man, "Array Placement for Storage Size Reduction in Embedded Multimedia Systems," Intl. Conf. on Application Specific Systems, Architectures, and Processors, 1997.
- [7] M. R. Garey, D. S. Johnson, *Computers and Intractability*, Freeman, 1979.
- [8] R. Lauwereins, P. Wauters, M. Ade, and J. A. Peperstraete, "Geometric Parallelism and Cyclo-static Data Flow in GRAPE-II," Proc. IEEE Wkshp Rapid Sys. Proto., 1994.

- [9] E. A. Lee, D. G. Messerschmitt, "Static Scheduling of Synchronous Dataflow Programs for Digital Signal Processing," *IEEE Trans. on Computers*, Feb., 1987.
- [10] S. Liao, S. Devadas, K. Keutzer, S. Tjiang, A. Wang, "Code Optimization Techniques in Embedded DSP Microprocessors," *DAES*, vol.3, (no.1), Kluwer, Jan. 1998.
- [11] P. Marwedel, G. Goossens, editors, *Code Generation for Embedded Processors*, Kluwer, 1995.
- [12] S. Moran, I. Newman, Y. Wolfstahl, "Approximation Algorithms for Covering a Graph by Vertex-Disjoint Paths of Maximum Total Weight," *Networks*, Vol. 20, pp55-64, 1990.
- [13] P. K. Murthy, S. S. Bhattacharyya, E. A. Lee, "Minimizing Memory Requirements for Chain-Structured SDF Graphs," Proc. of ICASSP, Australia, 1994.
- [14] P. K. Murthy, S. S. Bhattacharyya, E. A. Lee, "Joint Code and Data Minimization for Synchronous Dataflow Graphs," *Journal on Formal Methods in System Design*, July 1997.
- [15] P. K. Murthy, S. S. Bhattacharyya, "Buffer Merging - A Powerful Technique for Reducing Memory Requirements of Synchronous Dataflow Specifications," Tech report UMIACS-TR-2000-20, University of Maryland Institute for Advanced Computer Studies, College Park, MD 20742, <http://www.cs.umd.edu/TRs/TRumiacs.html>, April 2000.
- [16] P. K. Murthy, S. S. Bhattacharyya, "Shared Buffer Memory Implementations of Synchronous Dataflow Specifications Using Lifetime Analysis Techniques," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, Vol. 20, No. 2, pp. 177-198, February 2001.
- [17] P. K. Murthy, E. G. Cohen, S. Rowland, "System Canvas: A New Design Environment for Embedded DSP and Telecommunication Systems," Proceedings of the Ninth International Symposium on Hardware/Software Codesign, Copenhagen, Denmark, April 2001.
- [18] C. Papadimitriou, K. Steiglitz, *Combinatorial Optimization*, Dover, 1998.
- [19] S. Ritz, M. Pankert, and H. Meyr, "Optimum Vectorization of Scalable Synchronous Dataflow Graphs," Proc. of the Intl. Conf. on ASAP, Oct. 1993.
- [20] S. Ritz, M. Willems, H. Meyr, "Scheduling for Optimum Data Memory Compaction in Block Diagram Oriented Software Synthesis," Proc. ICASSP 95, May 1995.
- [21] W. Sung, J. Kim, S. Ha, "Memory Efficient Synthesis from Dataflow Graphs," ISSS, Hinschu, Taiwan, 1998.
- [22] I. Verbauwhede, F. Catthoor, J. Vandewalle, and H. De Man, "In-place Memory Management of Algebraic Algorithms on Application Specific ICs," *J. VLSI SP*, 1991.
- [23] V. Zivojinovic, J. M. Velarde, C. Schlager, H. Meyr, "DSPStone — A DSP-oriented Benchmarking Methodology," ICSPAT, 1994.

- [24] V. Zivojinovic, S. Ritz, H. Meyr, "Retiming of DSP Programs for Optimum Vectorization," Proceedings of the ICASSP, April, 1994.