

Generating Compact Code from Dataflow Specifications of Multirate Signal Processing Algorithms

Shuvra S. Bhattacharyya, *Member, IEEE*, Joseph T. Buck,
Soonhoi Ha, *Member, IEEE*, and Edward A. Lee, *Fellow, IEEE*

Abstract—Synchronous dataflow (SDF) semantics are well-suited to representing and compiling multirate signal processing algorithms. A key to this match is the ability to cleanly express iteration without overspecifying the execution order of computations, thereby allowing efficient schedules to be constructed. Due to limited program memory, it is often desirable to translate the iteration in an SDF graph into groups of repetitive firing patterns so that loops can be constructed in the target code. This paper establishes fundamental topological relationships between iteration and looping in SDF graphs, and presents a scheduling framework that provably synthesizes the most compact looping structures for a large class of practical SDF graphs. By modularizing different components of the scheduling framework, and establishing their independence, we show how other scheduling objectives, such as minimizing data buffering requirements or increasing the number of data transfers that occur in registers, can be incorporated in a manner that does not conflict with the goal of code compactness.

I. INTRODUCTION

IN THE dataflow model of computation, pioneered by Dennis [6], a program is represented as a directed graph in which the nodes represent computations and the arcs specify the passage of data. Synchronous dataflow (SDF) [15] is a restricted form of dataflow in which the nodes, called *actors*, consume a fixed number of data items, called *tokens* or *samples*, per invocation and produce a fixed number of output samples per invocation. SDF and related models have been

Manuscript received May 25, 1993; revised December 1, 1994. This work was part of the Ptolemy project, supported by the Advanced Research Projects Agency and U. S. Air Force (RASSP program, Contract F33615-93-C-1317), Semiconductor Research Corporation (Project 94-DC-008), National Science Foundation (MIP-9201605), Office of Naval Technology (Naval Research Laboratories), State of California MICRO program, and the following companies: Bell Northern Research, Cadence, Dolby, Hitachi, Mentor Graphics, Mitsubishi, NEC, Pacific Bell, Philips, Rockwell, Sony, and Synopsys. This paper was recommended by Associate Editor D. Mlynski.

S. S. Bhattacharyya was with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA. He is now with the Semiconductor Research Laboratory, Hitachi America, Ltd., San Jose, CA 95134 USA.

J. T. Buck was with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA. He is now with Synopsys, Inc., Mountain View, CA 94043 USA.

S. Ha was with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA. He is now with the Department of Computer Engineering, Seoul National University, Sinlim-Dong, Gwanak-Ku, Seoul 151-742 Korea.

E. A. Lee is with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, CA 94720 USA.

IEEE Log Number 9409315.

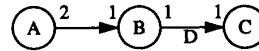


Fig. 1. A simple SDF graph.

studied extensively in the context of synthesizing assembly code for signal processing applications, for example [8]–[11], [17], [19]–[21].

Fig. 1 shows a simple SDF graph with three actors, labeled A, B and C. Each arc is annotated with the number of samples produced by its source and the number of samples consumed by its sink. Thus, actor A produces two samples on its output arc each time it is invoked and B consumes one sample from its input arc. The “D” on the arc directed from B to C designates a unit delay, which we implement as an initial token on the arc.

In SDF, *iteration* is induced whenever the number of samples produced on an arc (per invocation of the source actor) does not match the number of samples consumed (per sink invocation) [13]. For example, in Fig. 1, actor B must be invoked two times for every invocation of actor A. Multirate applications often involve a large amount of iteration and thus subroutine calls must be used extensively, code must be replicated, or loops must be organized in the target program. The use of subroutine calls to implement repetition may reduce throughput significantly however, particularly for graphs involving small granularity. On the other hand, we have found that code duplication can quickly exhaust on-chip program memory [12]. Thus, it may be essential that we arrange loops in the target code. In this paper we develop topological relationships between iteration and looping in SDF graphs.

We emphasize that in this paper, we view dataflow as a programming model, not as a form of computer architecture [2]. Several programming languages used for DSP, such as Lucid [25], SISAL [16], and Silage [10] are based on, or include dataflow semantics. The developments in this paper are applicable to this class of languages. Compilers for such languages can easily construct a representation of the input program as a hierarchy of dataflow graphs. It is important for a compiler to recognize SDF components of this hierarchy, since in DSP applications, usually a large fraction of the computation can be expressed with SDF semantics. For example, in [7] Dennis shows how to convert recursive stream functions in SISAL-2 into SDF graphs.

In [12] How showed that we can often greatly improve looping by clustering subgraphs that operate at the same repetition rate, and scheduling such subgraphs as a single unit. Fig. 1 shows how this technique can improve looping. A naive scheduler might schedule this SDF graph as CABCB, which offers no looping possibility within the schedule period. However, if we first group the subgraph $\{B, C\}$ into a hierarchical “supernode” Ω , a scheduler will generate the schedule $A\Omega\Omega$. To highlight the repetition in a schedule, we let the notation $(nX_1X_2\cdots X_m)$ designate n successive repetitions of the firing sequence $X_1X_2\cdots X_m$. We refer to a schedule expressed with this notation as a **looped schedule**. Using this notation, and substituting each occurrence of Ω with a subschedule for the corresponding subgraph, our clustering of the uniform-rate set $\{B, C\}$ leads to either $A(2BC)$ or $A(2CB)$, both of which expose the full potential for looping in the SDF graph of Fig. 1.

We explored the looping problem further in [5]. First, we generalized How’s scheme to exploit looping opportunities that occur across sample-rate changes. Our approach involved constructing the subgraph hierarchy in a pairwise fashion by clustering exactly two nodes at each step. Our subgraph selection was based on frequency of occurrence—we selected the pair of adjacent nodes whose associated subgraph had the largest repetition count. The “repetition count” of a subgraph can be viewed as the number of times that a minimal schedule for the subgraph is repeated in a minimal schedule for the overall graph. We will define this concept precisely in the next section.

By not discriminating against sample-rate boundaries, our approach exposed looping more thoroughly than How’s scheme. Furthermore, by selecting subgraphs based on repetition count, we reduced data memory requirements, an aspect that How’s scheme did not address.

Clustering a subgraph must be done with care since certain groupings cause deadlock. Thus, for each candidate subgraph, we must first verify that its consolidation does not result in an unschedulable graph. One way to perform this check is to attempt to schedule the new SDF graph [14], but this approach is extremely time consuming if a large number of clustering candidates must be considered. In [5], we employed a computationally more efficient method in which we maintained the subgraph hierarchy on the acyclic precedence graph rather than the SDF graph. Thus we could verify whether or not a grouping introduced deadlock by checking whether or not it introduced a cycle in the precedence graph. Furthermore, we showed that this check can be performed quickly by applying a *reachability matrix*, which indicates for any two precedence graph nodes (invocations) P_1 and P_2 , whether there is a precedence path from P_1 to P_2 .

Two limitations surfaced in the approach of [5]. First, the storage cost of the reachability matrix proved prohibitive for multirate applications involving very large sample rate changes. Observe that this cost is quadratic in the number of distinct actor *invocations* (precedence graph nodes). For example, a rasterization actor that decomposes an image into component pixels may involve a sample-rate change on the order of 250 000 to 1. If the rasterization output is connected to

a homogenous block (for example, a gamma level correction), this block alone will produce on the order of $(250\,000)^2 = 6.25 \times 10^{10}$ entries in the reachability matrix! Thus very large rate changes preclude straightforward application of the reachability matrix; this is unfortunate because looping is most important precisely for such cases. The second limitation in [5] is its failure to process cyclic paths in the graph optimally. Since cyclic paths limit looping, first priority should be given to preserving the full amount of looping available within the strongly connected components [1] of the graph. As Fig. 2 illustrates, clustering subgraphs based on repetition count alone does not fully carry out this goal.

In this paper, we develop a class of uniprocessor scheduling algorithms that extract the most compact looping structure from the cyclic paths in the SDF graph. This scheduling *framework* is based on a topological quality that we call “tight interdependence.” We show that for SDF graphs that contain no tightly interdependent subgraphs, our framework always synthesizes the most compact looping structures. Interestingly and fortunately, a large majority of practical SDF graphs seem to fall into this category. Furthermore, for this class of graphs, our technique does not require use of the reachability matrix, the precedence graph, or any other unreasonably large data structure. For graphs that contain tightly interdependent subgraphs, we show that our scheduling framework naturally isolates the minimal subgraphs that require special care. Only when analyzing these “tightly interdependent components,” do we need to apply reachability matrix-based analysis, or some other explicit deadlock-detection scheme.

An important aspect of our scheduling framework is its flexibility. By modularizing the framework into “sub-algorithms,” we allow other scheduling objectives to be integrated in a manner that does not conflict with code compactness objectives. Also, we show how decisions that a scheduler makes about grouping, or “clustering,” computations together can be formally evaluated in terms of their effects on program compactness. As an example, we demonstrate a very efficient clustering technique for increasing the amount of buffering that is done in machine registers, as opposed to memory, and we prove that this clustering strategy preserves codes space compactness for a large class of SDF graphs.

II. BACKGROUND

An SDF program is normally translated into a loop, where each iteration of the loop executes one cycle of a periodic schedule for the graph. In this section we summarize important properties of such periodic schedules. Most of the terminology introduced in this and subsequent sections is summarized in the glossary at the end of the paper.

For an SDF graph G , we denote the set of nodes in G by $N(G)$ and the set of arcs in G by $A(G)$. For an SDF arc α , we let $source(\alpha)$ and $sink(\alpha)$ denote the nodes at the source and the sink of α ; we let $p(\alpha)$ denote the number of samples produced by $source(\alpha)$, $c(\alpha)$ denote the number of samples consumed by $sink(\alpha)$, and we denote the delay on α by $delay(\alpha)$. We define a **subgraph** of G to be that SDF graph formed by any $Z \subseteq N(G)$ together with the set of arcs

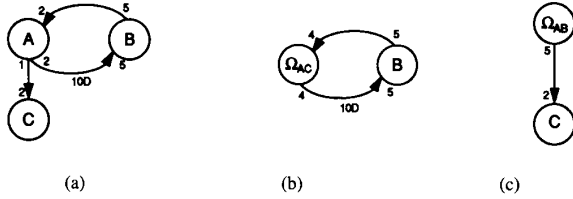


Fig. 2. This example illustrates how clustering based on repetition count alone can conceal looping opportunities within cyclic paths. Part (a) depicts a multirate SDF graph. Two pairwise clusterings lead to graphs that have schedules—{A, B}, having repetition count 2, and {A, C}, having repetition count 5 (clustering B and C results in deadlock). Clustering the subgraph with the highest repetition count yields the hierarchical topology in (b), for which the most compact schedule is $(2B)(2\Omega_{AC})B\Omega_{AC}B(2\Omega_{AC}) \Rightarrow (2B)(2(2A)C)B(2A)CB(2(2A)C)$. Clustering the subgraph {A, B} of lower repetition count, as depicted in part (c), yields the more compact schedule $(2\Omega_{AB})(5C) \Rightarrow (2(2B)(5A))(5C)$.

$\{\alpha \in A(G) | source(\alpha), sink(\alpha) \in Z\}$. We denote the subgraph associated with the subset of nodes Z by $subgraph(Z, G)$; if G is understood, we may simply write $subgraph(Z)$. If N_1 and N_2 are two nodes in an SDF graph, we say that N_1 is a *successor* of N_2 if there is an arc directed from N_2 to N_1 ; we say that N_1 is a *predecessor* of N_2 if N_2 is a successor of N_1 ; and we say that N_1 and N_2 are *adjacent* if N_1 is a predecessor or successor of N_2 . A sequence of nodes (N_1, N_2, \dots, N_k) is a *path* from N_1 to N_k if N_{i+1} is a successor of N_i for $i = 1, 2, \dots, (k-1)$. A sequence of nodes (N_1, N_2, \dots, N_k) is a *chain* that joins N_1 and N_k if N_{i+1} is adjacent to N_i for $i = 1, 2, \dots, (k-1)$.

We can think of each arc in G as having a FIFO queue that buffers the tokens that pass through the arc. Each FIFO contains an initial number of samples equal to the delay on the associated arc. Firing a node in G corresponds to removing $c(\alpha)$ tokens from the head of the FIFO for each input arc α , and appending $p(\beta)$ tokens to the FIFO for each output arc β . After a sequence of 0 or more firings, we say that a node is *fireable* if there are enough tokens on each input FIFO to fire the node. An *admissible sequential schedule* (“sequential” is used to distinguish this type of schedule from a parallel schedule) for G is a finite sequence $S = S_1 S_2 \dots S_N$ of nodes in G such that each S_i is fireable immediately after S_1, S_2, \dots, S_{i-1} have fired in succession.

We say that a sequential schedule S is a *periodic schedule* if it invokes each node at least once and produces no net change in the number of tokens on any arc’s FIFO—for each arc α , (the number of times $source(\alpha)$ is fired in S) $\times p(\alpha) =$ (the number of times $sink(\alpha)$ is fired in S) $\times c(\alpha)$. A *periodic admissible sequential schedule* (PASS) is a schedule that is both periodic and admissible. We will use the term *valid schedule* to describe a schedule that is a PASS, and the term *consistent* to describe an SDF graph that has a PASS. Except where otherwise stated, we deal only with consistent SDF graphs in this paper.

In [14], it is shown that for each connected SDF graph G , there is a unique minimum number of times that each node needs to be invoked in a periodic schedule. We specify these minimum numbers of firings by a vector of positive integers q_G , which is indexed by the nodes in G , and we denote the component of q_G corresponding to a node N by

$q_G(N)$. Every PASS for G invokes each node N a multiple of $q_G(N)$ times, and corresponding to each PASS S , there is a positive integer $J(S)$ called the *blocking factor* of S , such that S invokes each $N \in N(G)$ exactly $J(S)q_G(N)$ times. We call q_G the *repetitions vector* of G . For example in Fig. 2(a), $q_G(A) = 10, q_G(B) = 4$, and $q_G(C) = 5$. The following properties of repetitions vectors are established in [14]:

Fact 1: The components of a repetitions vector are collectively coprime.

Fact 2: The *balance equation* $q_G(source(\alpha)) \times p(\alpha) = q_G(sink(\alpha)) \times c(\alpha)$ is satisfied for each arc α in G .

Given a subset Z of nodes in a connected SDF graph G , we define $q_G(Z) = gcd(\{q_G(N) | N \in Z\})$, where gcd denotes the greatest common divisor. We can interpret $q_G(Z)$ as the number of times that G invokes the “subsystem” Z . We will use the following property of connected subsystems which is derived in [4].

Fact 3: If G is a connected SDF graph, and Z is a connected subset of $N(G)$, then for each $N \in Z$, $q_G(N) = q_G(Z)q_{subgraph(Z)}(N)$.

For our hierarchical scheduling approach, we will apply the concept of *clustering* a subgraph. This process is illustrated in Fig. 2. Here $subgraph(\{A, C\})$ of (a) is clustered into the hierarchical node Ω_{AC} , and the resulting SDF graph is shown in (b). Similarly, clustering $subgraph(\{A, B\})$ results in the graph of (c). Each input arc α to a clustered subgraph P is replaced by an arc α' having $p(\alpha') = p(\alpha)$, and $c(\alpha') = c(\alpha) \times q_G(sink(\alpha)) / q_G(N(P))$, the number of samples consumed from α in one *invocation* of *subgraph* P . Similarly we replace each output arc β with β' such that $c(\beta') = c(\beta)$, and

$$p(\beta') = p(\beta) \times q_G(source(\alpha)) / q_G(N(P)).$$

The following properties of clustered subgraphs are proven in [4].

Fact 4: Suppose G is a connected SDF graph, Z is a subset of nodes in G , G' is the SDF graph that results from clustering $subgraph(Z)$ into the hierarchical node Ω , and S' is a PASS for G' . Suppose that S_Z is a PASS for $subgraph(Z)$ such that for each $N \in Z, S_Z$ invokes N ($q_G(N) / q_G(Z)$) times. Let S^* denote the schedule that results from replacing each appearance of Ω in S' with S_Z . Then S^* is a PASS for G .

Fact 5: Suppose G is a connected SDF graph, Z is a subset of nodes in G , and G' is the SDF graph that results from clustering $subgraph(Z)$ into the node Ω . Then $q_{G'}(\Omega) = q_G(Z)$; and for any node N in G' other than $\Omega, q_{G'}(N) = q_G(N)$.

Given a directed graph G , we say that G is **strongly connected** if for any pair of distinct nodes A, B in G , there is a path from A to B and a path from B to A . We say that a strongly connected graph is *nontrivial* if it contains more than one node. Finally, a *strongly connected component* of G is a subset of nodes Z such that $subgraph(Z, G)$ is strongly connected, and there is no strongly connected subset of $N(G)$ that properly contains Z . For example $\{A, B\}$ and $\{C\}$ are the strongly connected components of Fig. 2(a).

Similarly, we define a *connected component* of a directed graph G to be a maximal subset of nodes Z such that for any

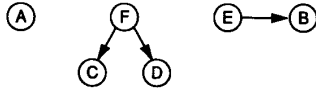


Fig. 3. A directed graph that has three connected components.

pair of distinct members A, B of Z, there is a chain that joins A and B. For example in Fig. 3, the connected components are {A}, {C, D, F}, and {B, E}.

Given a connected SDF graph G, and an arc α in G, we define $total_consumed(\alpha, G)$ to be the total number of samples consumed from α in a minimal schedule period for G. Thus $total_consumed(\alpha, G) = q_G(sink(\alpha))c(\alpha)$. Finally, given an SDF graph G, a looped schedule S for G and a node N in G, we define $appearances(N, S)$ to be the number of times that N appears in S, and we say that S is a **single appearance schedule** if for each $N \in N(G)$, $appearances(N, S) = 1$. For example, consider the two schedules $S_1 = CA(2B)C$ and $S_2 = A(2B)(2C)$ for Fig. 1. We have $appearances(C, S_1) = 2$; $appearances(C, S_2) = 1$; S_1 is not a single appearance schedule because C appears more than once; and S_2 is a single appearance schedule. Single appearance schedules form the class of schedules that allow in line code generation without any code space or subroutine penalty.

III. SUBINDEPENDENCE

Our scheduling framework for synthesizing compact nested loop structures is based on a form of precedence independence, which we call *subindependence*.

Definition 1: Suppose that G is a connected SDF graph. If Z_1 and Z_2 are disjoint, nonempty subsets of $N(G)$ we say that “ Z_1 is **subindependent** of Z_2 in G” if for every arc α in G such that $source(\alpha) \in Z_2$ and $sink(\alpha) \in Z_1$, we have $delay(\alpha) \geq total_consumed(\alpha, G)$. We occasionally drop the “in G” qualification if G is understood from context. If (Z_1 is subindependent of Z_2) and ($Z_1 \cup Z_2 = N(G)$), then we write ($Z_1 |G Z_2$), and we say that Z_1 is *subindependent* in G.

Thus Z_1 is subindependent of Z_2 if no samples produced from Z_2 are consumed by Z_1 in the same schedule period that they are produced; and $Z_1 |G Z_2$ if Z_1 is subindependent of Z_2 , and Z_1 and Z_2 form a partition of the nodes in G. For example, consider Fig. 2(a). Here $q_G(A, B, C) = (10, 4, 5)$, and the complete set of subindependence relationships is (1) {A} is subindependent of {C}; (2) {B} is subindependent of {C}; (3) {A, B} |G {C}; and {C} is subindependent of {B}.

The following property of subindependence follows immediately from definition 1.

Fact 6: If G is a strongly connected SDF graph and X, Y, and Z are disjoint subsets of $N(G)$, then (a) (X is subindependent of Z) and (Y is subindependent of Z) \Rightarrow (X \cup Y) is subindependent of Z. (b) (X is subindependent of Y) and (X is subindependent of Z) \Rightarrow X is subindependent of (Y \cup Z).

Our scheduling framework is based on the following condition for the existence of a single appearance schedule, which is developed in [4].

Fact 7: An SDF graph has a valid single appearance schedule iff for each nontrivial strongly connected component Z,

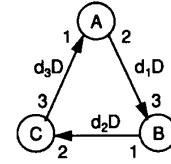


Fig. 4. An illustration of loose and tight interdependence. Here d_1 , d_2 , and d_3 represent the number of delays on the associated arcs. This SDF graph is tightly interdependent if and only if ($d_1 < 6$), ($d_2 < 2$), and ($d_3 < 3$).

there exists a partition X, Y of Z such that X *subgraph*(Z) Y, and *subgraph*(X) and *subgraph*(Y) each have single appearance schedules.

A related condition was developed independently by Ritz *et al.* in [22], which discusses single appearance schedules in the context of *minimum activation schedules*. For example, the schedule A(2CB) for Fig. 1 results in 5 activations since invocations of C and B are interleaved. In contrast, the schedule A(2B)(2C) requires only one activation per actor, for a total of 3 activations. In the objectives of [22], the latter schedule is preferable because in that code generation framework there is a large overhead associated with each activation. However such overhead can often be avoided with careful instruction scheduling and register allocation, as [19] demonstrates. We prefer the former schedule, which has less looping overhead and requires less memory for buffering.

Fact 7 implies that for an SDF graph to have a single appearance schedule, we must be able to decompose each nontrivial strongly connected component into two subsets in such a way that one subset is subindependent of the other. Another implication of fact 7 is that every acyclic SDF graph has a single appearance schedule. We can easily construct a single appearance schedule for an acyclic SDF graph. We simply pick a root node N_1 ; schedule all of its invocations in succession; remove N_1 from the graph and pick a root node N_2 of the remaining graph; schedule all of N_2 's invocations in succession; and so on until we have scheduled all of the nodes. By this procedure, we get a cascade of loops ($q_G(N_1)N_1)(q_G(N_2)N_2) \cdots (q_G(N_k)N_k)$, which gives us a single appearance schedule.

Definition 2: Suppose that G is a nontrivial strongly connected SDF graph. Then we say that G is **loosely interdependent** if $N(G)$ can be partitioned into Z_1 and Z_2 such that $Z_1 |G Z_2$. We say that G is **tightly interdependent** if it is not loosely interdependent.

For example, consider the strongly connected SDF graph in Fig. 4. The repetitions vector for this graph is $q_G(A, B, C) = (3, 2, 1)$. Thus the graph is loosely interdependent if and only if ($d_1 \geq 6$) or ($d_2 \geq 2$) or ($d_3 \geq 3$).

In this section we have introduced topological properties of SDF graphs that are related to the existence of single appearance schedules. In the following section we use these properties to develop our scheduling framework and to demonstrate some of its useful qualities.

IV. THE CLASS OF LOOSE INTERDEPENDENCE ALGORITHMS

The properties of loose/tight interdependence are important for organizing loops because, as we will show, the existence of

a single appearance schedule is equivalent to the absence of tightly interdependent subgraphs. However, these properties are useful even when tightly interdependent subgraphs are present. The following definition specifies how to use loose interdependence to guide the looping process.

Definition 3: Let A_1 be any algorithm that takes as input a nontrivial strongly connected SDF graph G , determines whether G is loosely interdependent, and if so, finds a subindependent subset of $N(G)$. Let A_2 be any algorithm that finds the strongly connected components of a directed graph. Let A_3 be any algorithm that takes an acyclic SDF graph and generates a valid single appearance schedule. Finally, let A_4 be any algorithm that takes a tightly interdependent SDF graph, and generates a valid looped schedule of blocking factor 1. We define the algorithm $L(A_1, A_2, A_3, A_4)$ as follows:

Input: a connected SDF graph G .

Output: a valid unit-blocking-factor looped schedule $S_L(G)$ for G .

Step 1: Use A_2 to determine the nontrivial strongly connected components Z_1, Z_2, \dots, Z_s of G .

Step 2: Cluster Z_1, Z_2, \dots, Z_s into nodes $\Omega_1, \Omega_2, \dots, \Omega_s$ respectively, and call the resulting graph G' . This is an acyclic SDF graph.

Step 3: Apply A_3 to G' ; denote the resulting schedule $S'(G)$.

Step 4: (Shown at the bottom of the page.)

The *for*-loop replaces each " Ω_i " in $S'(G)$ with a valid looped schedule for $subgraph(Z_i)$. From repeated application of fact 4, we know that these replacements yield a valid looped schedule S_L for G . We output S_L . ■

Remark 1: Observe that step 4 does not insert or delete appearances of actors that are not contained in a nontrivial strongly connected component Z_i . Since A_3 generates a single appearance schedule for G' , we have that for every node N that is not contained in a nontrivial strongly connected component of G , $appearances(N, S_L(G)) = 1$.

Remark 2: If C is a nontrivial strongly connected component of G and $N \in C$, then since $S_L(G)$ is

derived from $S'(G)$ by replacing the single appearance of each Ω_i , we have $appearances(N, S_L(G)) = appearances(N, S_L(subgraph(C)))$.

Remark 3: For each strongly connected component Z_k whose subgraph is loosely interdependent, L partitions Z_k into X and Y such that $X | subgraph(Z_k) Y$, and replaces the single appearance of Ω_k in $S'(G)$ with $S_x S_y$. If N is a member of the connected component X_i , then $N \neq Y$, so $appearances(N, S_x S_y) = appearances(N, S_L(subgraph(X_i)))$. Also since N cannot be in any other strongly connected component besides Z_k , and since $S'(G)$ contains only one appearance of Ω_k , we have $appearances(N, S_L(G)) = appearances(N, S_x S_y)$. Thus, for $i = 1, 2, \dots, v, N \in X_i \Rightarrow appearances(N, S_L(G)) = appearances(N, S_L(subgraph(X_i)))$. By the same argument, we can show that for $i = 1, 2, \dots, w, N \in Y_i \Rightarrow appearances(N, S_L(G)) = appearances(N, S_L(subgraph(Y_i)))$.

$L(\bullet, \bullet, \bullet, \bullet)$ defines a family of algorithms, which we call **loose interdependence algorithms** because they exploit loose interdependence to decompose the input SDF graph. Since nested recursive calls decompose a graph into finer and finer strongly connected components, it is easy to verify that any loose interdependence algorithm always terminates. Each loose interdependence algorithm $\lambda = L(A_1, A_2, A_3, A_4)$ involves the "sub-algorithms" A_1, A_2, A_3 , and A_4 , which we call, respectively, the *subindependence partitioning algorithm* of λ , the *strongly connected components algorithm* of λ , the *acyclic scheduling algorithm* of λ , and the *tight scheduling algorithm* of λ .

We will apply a loose interdependence algorithm to derive a *nonrecursive* necessary and sufficient condition for the existence of a single appearance schedule. First, we introduce two lemmas.

Lemma 1: Suppose G is a connected SDF graph; N is a node in G that is not contained in any tightly interdependent subgraph of G ; and λ is a loose interdependence algorithm. Then N appears only once in $S_\lambda(G)$, the schedule generated by λ .

Step 4:

for $i = 1, 2, \dots, s$

Let SZ denote $subgraph(Z_i)$.

Apply A_1 to SZ .

if $X, Y \subseteq Z_i$ are found such that $X | SZ Y$,

then

- Determine the connected components X_1, X_2, \dots, X_v of $subgraph(X)$, and the connected components Y_1, Y_2, \dots, Y_w of $subgraph(Y)$.

- Recursively apply algorithm L to construct the schedules

$$S_x = (q_{SZ}(X_1)S_L(subgraph(X_1)) \cdots (q_{SZ}(X_v)S_L(subgraph(X_v))),$$

$$S_y = (q_{SZ}(Y_1)S_L(subgraph(Y_1)) \cdots (q_{SZ}(Y_w)S_L(subgraph(Y_w))).$$

- Replace the (single) appearance of Ω_i in $S'(G)$ with $S_x S_y$.

else (SZ is tightly interdependent)

- Apply A_4 to obtain a valid schedule S_i for SZ .

- Replace the single appearance of Ω_i in $S'(G)$ with S_i .

end - if

end - for

The proof of lemma 1 can be found in the appendix.

Lemma 2: Suppose that G is a strongly connected SDF graph, $P \subseteq N(G)$ is subindependent in G , and C is a strongly connected subset of $N(G)$ such that $C \cap P \neq C$ and $C \cap P \neq \emptyset$. Then $C \cap P$ is subindependent in $subgraph(C)$.

Proof: Suppose that α is an arc directed from a member of $(C - (C \cap P))$ to a member of $(C \cap P)$. By the subindependence of P in G , $delay(\alpha) \geq c(\alpha) \times q_G(sink(\alpha))$, and by fact 3, $q_G(sink(\alpha)) \geq q_{subgraph(C)}(sink(\alpha))$. Thus, $delay(\alpha) \geq c(\alpha) \times q_{subgraph(C)}(sink(\alpha))$. Since this holds for any α directed from $(C - (C \cap P))$ to $(C \cap P)$, we conclude that $(C \cap P)$ is subindependent in C . QED

Corollary 1: Suppose that G is a strongly connected SDF graph, Z_1 and Z_2 are subsets of $N(G)$ such that $Z_1 | G Z_2$, and T is a tightly interdependent subgraph of G . Then $N(T) \subseteq Z_1$ or $N(T) \subseteq Z_2$.

Proof: (By contraposition.) If $N(T)$ has nonempty intersection with both Z_1 and Z_2 , then from lemma 2, $N(T) \cap Z_1$ is subindependent in T , so T is loosely interdependent. QED

Theorem 1: Suppose that G is a strongly connected SDF graph. Then G has a single appearance schedule iff every nontrivial strongly connected subgraph of G is loosely interdependent.

Proof: \Leftarrow Suppose every nontrivial strongly connected subgraph of G is loosely interdependent, and let λ be any loose interdependence algorithm. Since no node in G is contained in a tightly interdependent subgraph, lemma 1 guarantees that $S_\lambda(G)$ is a single appearance schedule for G .

\Rightarrow Suppose that G has a single appearance schedule and that C is a strongly connected subset of $N(G)$. Set $Z_0 = G$. From fact 7, there exist $X_0, Y_0 \subseteq Z_0$ such that $X_0 | subgraph(Z_0) Y_0$, and $subgraph(X_0)$ and $subgraph(Y_0)$ both have single appearance schedules. If X_0 and Y_0 do not both intersect C then C is completely contained in some strongly connected component Z_1 of $subgraph(X_0)$ or $subgraph(Y_0)$. We can then apply fact 7 to partition Z_1 into X_1, Y_1 , and continue recursively in this manner until we obtain a strongly connected $Z_k \subseteq N(G)$, with the following properties: Z_k can be partitioned into X_k and Y_k such that $X_k | subgraph(Z_k) Y_k$; $C \subseteq Z_k$; and $(X_k \cap C)$ and $(Y_k \cap C)$ are both nonempty. From lemma 2, $(X_k \cap C)$ is subindependent in $subgraph(C)$, so C must be loosely interdependent. QED

Corollary 2: Given a connected SDF graph G , any loose interdependence algorithm will obtain a single appearance schedule if one exists.

Proof: If a single appearance schedule for G exists, then from theorem 1, G contains no tightly interdependent subgraphs. In other words, no node in G is contained in a tightly interdependent subgraph of G . From lemma 1, the schedule resulting from any loose interdependence algorithm contains only one appearance for each actor in G . QED

Thus, a loose interdependence algorithm always obtains an optimally compact solution when a single appearance schedule exists. When a single appearance schedule does not exist, strongly connected graphs are repeatedly decomposed until tightly interdependent subgraphs are found. In general, however, there may be more than one way to decompose

$N(G)$ into two parts so that one of the parts is subindependent of the other. Thus, it is natural to ask the following question: Given two distinct partitions $\{Z_1, Z_2\}$ and $\{Z'_1, Z'_2\}$ such that $Z_1 | G Z_2$ and $Z'_1 | G Z'_2$, is it possible that one of these partitions leads to a more compact schedule than the other? Fortunately, as we will show in the remainder of this section, the answer to this question is “No”. In other words, any two loose interdependence algorithms that use the same tight scheduling algorithm always lead to equally compact schedules. The key reason is that tight interdependence is an additive property.

Lemma 3: Suppose that G is a connected SDF graph, Y and Z are subsets of $N(G)$ such that $(Y \cap Z) \neq \emptyset$, and $subgraph(Y)$ and $subgraph(Z)$ are both tightly interdependent. Then $subgraph(Y \cup Z)$ is tightly interdependent.

Proof: (By contraposition.) Let $H = Y \cup Z$, and suppose that $subgraph(H)$ is loosely interdependent. Then there exist H_1 and H_2 such that $H = H_1 \cup H_2$ and $H_1 | subgraph(H) H_2$. From $H_1 \cup H_2 = Y \cup Z$, and $Y \cap Z \neq \emptyset$, it is easily seen that H_1 and H_2 both have a nonempty intersection with Y , or they both have a nonempty intersection with Z . Without loss of generality, assume that $H_1 \cap Y \neq \emptyset$ and $H_2 \cap Y \neq \emptyset$. From lemma 2, $(H_1 \cap Y)$ is subindependent in $subgraph(Y)$, and thus $subgraph(Y)$ is not tightly interdependent. QED

Lemma 3 implies that each SDF graph G has a *unique* set $\{C_1, C_2, \dots, C_n\}$ of maximal tightly interdependent subgraphs such that $i \neq j \Rightarrow N(C_i) \cap N(C_j) = \emptyset$, and every tightly interdependent subgraph in G is contained in some C_i . We call each $N(C_i)$ a *tightly interdependent component* of G . It follows from theorem 1 that G has a single appearance schedule iff G has no tightly interdependent components. Furthermore, since the tightly interdependent components are unique, the performance of a loose interdependence algorithm, with regards to schedule compactness, is not dependent on the particular subindependence partitioning algorithm, the sub-algorithm used to partition the loosely interdependent components. The following theorem develops this result.

Theorem 2: Suppose G is an SDF graph that has a PASS, N is a node in G , and λ is a loose interdependence algorithm. If N is not contained in a tightly interdependent component of G , then N appears only once in $S_\lambda(G)$. On the other hand, if N is contained in a tightly interdependent component T then $appearances(N, S_\lambda(G)) = appearances(N, S_\lambda(subgraph(T)))$ —the number of appearances of N is determined entirely by the tight scheduling algorithm of λ .

Proof: If N is not contained in a tightly interdependent component of G , then N is not contained in any tightly interdependent subgraph. Then from lemma 1, $appearances(N, S_\lambda(G)) = 1$.

Now suppose that N is contained in some tightly interdependent component T of G . If $T = N(G)$ we are done. Otherwise we set $M_0 = N(G)$, and thus $T \neq M_0$; by definition, tightly interdependent graphs are strongly connected, so T is contained in some strongly connected component C of $subgraph(M_0)$.

If T is a proper subset of C , then $subgraph(C)$ must be loosely interdependent, since otherwise $subgraph(T)$ would not be a maximal tightly interdependent subgraph.

Thus, λ partitions $subgraph(C)$ into X and Y such that $X \subseteq subgraph(C) \setminus Y$. We set M_1 to be that connected component of $subgraph(X)$ or $subgraph(Y)$ that contains N . Since X, Y partition C , M_1 is a proper subset of M_0 . Also, from remark 3, $appearances(N, S_\lambda(subgraph(M_0))) = appearances(N, S_\lambda(subgraph(M_1)))$, and from corollary 1, $N(T) \subseteq M_1$.

On the other hand, if $T = C$, then we set $M_1 = T$. Since $T \neq M_0$, M_1 is a proper subset of M_0 ; from remark 2, $appearances(N, S_\lambda(subgraph(M_0))) = appearances(N, S_\lambda(subgraph(M_1)))$; and trivially, $T \subseteq M_1$.

If $T \neq M_1$, then we can repeat the above procedure to obtain a proper subset M_2 of M_1 such that $appearances(N, S_\lambda(subgraph(M_1))) = appearances(N, S_\lambda(subgraph(M_2)))$, and $N(T) \subseteq M_2$. Continuing this process, we get a sequence M_1, M_2, \dots . Since each M_i is a proper subset of its predecessor, we cannot repeat this process indefinitely—eventually, for some $k \geq 0$, we will have $N(T) = M_k$. But, by construction, $appearances(N, S_\lambda(G)) = appearances(N, S_\lambda(subgraph(M_0))) = appearances(N, S_\lambda(subgraph(M_1))) = \dots = appearances(N, S_\lambda(subgraph(M_k)))$; and thus $appearances(N, S_\lambda(G)) = appearances(N, S_\lambda(subgraph(T)))$. QED

Theorem 2 states that the tight scheduling algorithm is independent of the subindependence partitioning algorithm, and vice-versa. Any subindependence partitioning algorithm makes sure that there is only one appearance for each actor outside the tightly interdependent components, and the tight scheduling algorithm completely determines the number of appearances for actors inside the tightly interdependent components. For example, if we develop a new subindependence partitioning algorithm that is more efficient in some way (e.g., it is faster or minimizes data memory requirements), we can replace it for any existing subindependence partitioning algorithm without changing the “compactness” of the resulting schedules—we don’t need to analyze its interaction with the rest of the loose interdependence algorithm. Similarly, if we develop a new tight scheduling algorithm that schedules any tightly interdependent graph more compactly than the existing tight scheduling algorithm, we are guaranteed that using the new algorithm instead of the old one will lead to more compact schedules overall.

V. COMPUTATIONAL EFFICIENCY

The complexity of a loose interdependence algorithm λ depends on its subindependence partitioning algorithm λ_{sp} , strongly connected components algorithm λ_{sc} , acyclic scheduling algorithm λ_{as} , and tight scheduling algorithm λ_{ts} . From the proof of theorem 2, we see that λ_{ts} is applied exactly once for each tightly interdependent component. For example, the simplest solution for a tight scheduling algorithm would be to apply an algorithm from the family of class-S scheduling algorithms that are defined in [14]; class-S algorithms exist whose complexity is linear in the number of actor firings (assuming that the number of input and output edges for a given actor is bounded) [3]. Alternatively, a more elaborate technique such as that presented in [5] can be employed.

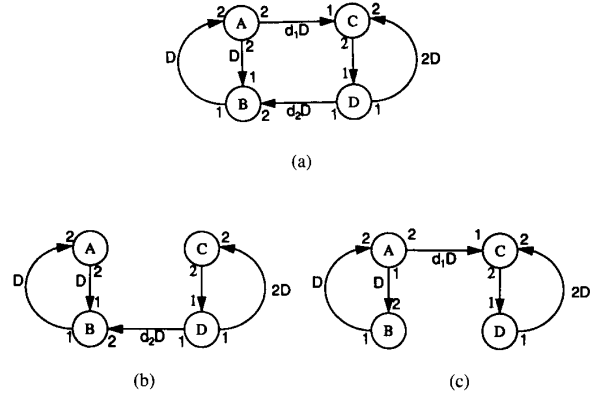


Fig. 5. An illustration of Theorem 3.

As mentioned earlier, one drawback of the technique of [5] is that it requires a reachability matrix, which has a storage cost that is quadratic in the number of actor firings. However, we greatly reduce this drawback by restricting application of the algorithm to only the tightly interdependent components. We are currently investigating other alternatives to scheduling tightly interdependent SDF graphs.

The other subalgorithms, λ_{sc} , λ_{as} , and λ_{sp} , are successively applied to decompose an SDF graph, and the process is repeated until all tightly interdependent components are found. In the worst case, each decomposition step isolates a single node from the current n -node subgraph, and the decomposition must be recursively applied to the remaining $(n - 1)$ -node subgraph. Thus, if the original program has n nodes, n decomposition steps are required in the worst case. Tarjan [24] first showed that the strongly connected components of a graph can be found in $O(m)$ time, where $m = \max(\text{number of nodes, number of arcs})$. Hence λ_{sc} can be chosen to be linear, and since at most $n \leq m$ decomposition steps are required, the total time that such a λ_{sc} accounts for in λ is $O(m^2)$. In Section III we presented a simple linear-time algorithm that constructs a single appearance schedule for an acyclic SDF graph. Thus λ_{as} can be chosen such that its total time is also $O(m^2)$.

The following theorem presents a simple topological condition for loose interdependence that leads to a linear subindependence partitioning algorithm λ_{sp} .

Theorem 3: Suppose that G is a nontrivial strongly connected SDF graph. From G , remove all arcs α for which $delay(\alpha) \geq c(\alpha) \times q_G(sink(\alpha))$, and call the resulting SDF graph G' . Then G is tightly interdependent if and only if G' is strongly connected.

For example, suppose that G is the strongly connected SDF graph in Fig. 5(a). The repetitions vector for G is $q_G(A, B, C, D) = (1, 2, 2, 4)$. This graph is loosely interdependent if $d_1 \geq 2$, which corresponds to $\{C, D\} | G \{A, B\}$, or if $d_2 \geq 4$, which corresponds to $\{A, B\} | G \{C, D\}$. The corresponding G' 's are depicted at the bottom of Fig. 5: Fig. 5(b) shows G' when $d_1 \geq 2$ and $d_2 < 4$, and Fig. 5(c) shows G' when $d_2 \geq 4$ and $d_1 < 2$. Observe that in both of these cases, G' is not strongly connected.

Proof: We prove both directions by contraposition.

\Rightarrow Suppose that G' is not strongly connected. Then $N(G')$ can be partitioned into Z_1 and Z_2 such that there is no arc directed from a member of Z_2 to a member of Z_1 in G' . Since no nodes were removed in constructing G' , Z_1 and Z_2 partition $N(G)$. Also, none of the arcs directed from Z_2 to Z_1 in G occur in G' . Thus, by the construction of G' , for each arc α in G directed from a member of Z_2 to a member of Z_1 , we have $\text{delay}(\alpha) \geq c(\alpha) \times q_G(\text{sink}(\alpha))$. It follows that $Z_1|GZ_2$, so G is loosely interdependent.

\Leftarrow Suppose that G is loosely interdependent. Then $N(G)$ can be partitioned into Z_1 and Z_2 such that $Z_1|GZ_2$. By construction of G' , there are no arcs in G' directed from a member of Z_2 to a member of Z_1 , so G' is not strongly connected. QED

Thus, λ_{sp} can be constructed as follows: (1) Determine $q_G(N)$ for each node N ; (2) Remove each arc α whose delay is at least $c(\alpha) \times q_G(\text{sink}(\alpha))$; (3) Determine the strongly connected components of the resulting graph; (4) If the entire graph is the only strongly connected component, then G is tightly interdependent; otherwise (5) cluster the strongly connected components—the resulting graph is acyclic and has at least two nodes. The strongly connected component corresponding to any root node of this graph is subindependent of the rest of the graph. An algorithm (first used in the Gabriel system [11]) that performs (1) in time $O(m)$ is described in [3]; it is obvious that (2) is $O(m)$; Tarjan's algorithm allows $O(m)$ for (3); and the checks in (4) and (5) are clearly $O(m)$ as well. Thus, we have a linear λ_{sp} , and the total time that λ spends in λ_{sp} is $O(m^2)$.

We have specified λ_{sp} , λ_{sc} , λ_{as} , and λ_{ts} such that the time complexity of the corresponding loose interdependence algorithm is $O(m^2 + f)$, where m is $\max(\text{number of nodes, number of arcs})$, and f is the number of actor firings. Note that our worst case estimate is conservative—in practice only a few decomposition steps are required to fully schedule a strongly connected subgraph, while our estimate assumes n steps, where n is the number of nodes in the input graph.

VI. CLUSTERING TO MAKE DATA TRANSFERS MORE EFFICIENT

In this section, we present a useful clustering technique for increasing the frequency of data transfers that occur through machine registers rather than memory, and we prove that this technique does not interfere with the code compactness potential of a loose interdependence algorithm—this clustering preserves the properties of loose interdependence algorithms discussed in Section IV.

Fig. 6 illustrates two ways in which arbitrary clustering decisions can conflict with code compactness objectives. Observe that Fig. 6(a) is an acyclic graph so it must have a single appearance schedule. Fig. 6(b) is the hierarchical SDF graph that results from clustering A and B in Fig. 6(a). It is easy to verify that this is a tightly interdependent graph. In fact, the only minimal periodic schedule for Fig. 6(a) that we can derive from this clustering is $C\Omega C \Rightarrow CAB C$. Thus, the

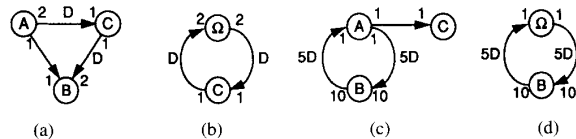


Fig. 6. Examples of how clustering can conflict with the goal of code compactness.

clustering of A and B in Fig. 6(a) cancels the existence of a single appearance schedule.

In Fig. 6(c), $\{A, B\}$ forms a tightly interdependent component and C is not contained in any tightly interdependent subgraph. From theorem 2, we know that any loose interdependence algorithm will schedule Fig. 6(c) in such a way that C appears only once. Now observe that the graph that results from clustering A and C, shown in Fig. 6(d), is tightly interdependent. It can be verified that the most compact minimal periodic schedule for this graph is $(5\Omega)B(5\Omega)$, which leads to the schedule $(5 AC)B(5 AC)$ for Fig. 6(c). By increasing the “extent” of the tightly interdependent component $\{A, B\}$ to subsume C, this clustering decision increases the minimum number of appearances of C in the final schedule.

Thus we see that a clustering decision can conflict with optimal code compactness if it introduces a new tightly interdependent component or extends an existing tightly interdependent component. In this section we present a clustering technique of great practical use and prove that it neither extends nor introduces tight interdependence. Our clustering technique and its compatibility with loose interdependence algorithms is summarized by the following claim: *Clustering two adjacent nodes A and B in an SDF graph does not introduce or extend a tightly interdependent component if (a) Neither A nor B is contained in a tightly interdependent component; (b) At least one arc directed from A to B has zero delay; (c) A and B are invoked the same number of times in a periodic schedule; and (d) B has no predecessors other than A or B.* The remainder of this section is devoted to proving this claim and explaining the corresponding clustering technique.

We motivate our clustering technique with the example shown in Fig. 7. One possible single appearance schedule for Fig. 7(a) is $(10 X)(10 Y)ZV(10 W)$. This is the *minimum activation* schedule preferred by Ritz *et al.* [22]; however, it is inefficient with respect to buffering. Due to the loop that specifies ten successive invocations of X, the data transfers between X and Y cannot take place in machine registers and 10 words of data-memory are required to implement the arc connecting X and Y. However, observe that conditions (a)–(d) of our above claim all hold for the adjacent pairs $\{X, Y\}$ and $\{Z, V\}$. Thus, we can cluster these pairs without cancelling the existence of a single appearance schedule. The hierarchical graph that results from this clustering is shown in Fig. 7(d); this graph leads to the single appearance schedule $(10\Omega_2)\Omega_1(10W) \Rightarrow (10 XY)ZV(10 W)$. In this second schedule, each sample produced by X is consumed by Y in the same loop iteration, so all of the transfers between X and Y can occur through a single machine register. Thus, the clustering of X and Y saves 10 words of buffer space for the

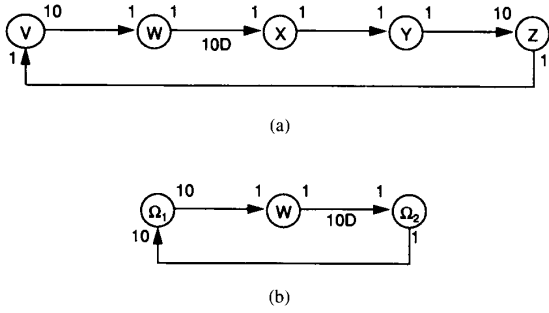


Fig. 7. An example of clustering to increase the amount of buffering that occurs through registers.

data transfers between X and Y, and it allows these transfers to be performed through registers rather than memory, which will usually result in faster code.

We will use the following additional notation in the development of this section.

Notation: Let G be an SDF graph and suppose that we cluster a subset W of nodes in G . We will refer to the resulting hierarchical graph as G' , and we will refer to the node in G' into which W has been clustered as Ω . For each arc α in G that is not contained in $\text{subgraph}(W)$, we denote the corresponding arc in G' by α' . Finally, if $X \subseteq N(G)$, we refer to the “corresponding” subset of $N(G')$ as X' . That is, X' consists of all members of X that are not in W ; and if X contains a member of W , then X' also contains Ω .

For example, if G is the SDF graph in Fig. 6(a), $W = \{A, B\}$, and α and β respectively denote the arc directed from A to C and the arc directed from C to B , then we denote the graph in Fig. 6(b) by G' , and in G' we denote the arc directed from Ω to C by α' and the arc denoted from C to Ω by β' . Also, if $X = \{A, C\}$, then $X' = \{\Omega, C\}$.

Lemma 4: Suppose that G is a strongly connected SDF graph and X_1, X_2 partition $N(G)$ such that $X_1 | G X_2$. Also suppose that A, B are nodes in G such that $A, B \in X_1$ or $A, B \in X_2$. If we cluster $W = \{A, B\}$ then the resulting SDF graph G' is loosely interdependent.¹

The proof of lemma 4 can be found in the appendix.

Definition 4: We say that two SDF graphs G_1 and G_2 are *isomorphic* if there exist bijective mappings $f_1: N(G_1) \rightarrow N(G_2)$ and $f_2: A(G_1) \rightarrow A(G_2)$ such that for each $\alpha \in A(G_1)$, $\text{source}(f_2(\alpha)) = f_1(\text{source}(\alpha))$, $\text{sink}(f_2(\alpha)) = f_1(\text{sink}(\alpha))$, $\text{delay}(f_2(\alpha)) = \text{delay}(\alpha)$, $p(f_2(\alpha)) = p(\alpha)$, and $c(f_2(\alpha)) = c(\alpha)$. Intuitively, two SDF graphs are isomorphic if they differ only by a relabeling of the nodes. For example, the SDF graph in Fig. 6(d) is isomorphic to $\text{subgraph}(\{A, B\})$ in Fig. 6(c).

We will use the following obvious fact about isomorphic SDF graphs.

Fact 8: If G_1 and G_2 are two isomorphic SDF graphs and G_1 is loosely interdependent then G_2 is loosely interdependent.

¹However, G' may be deadlocked even if G is not. This will not be a problem in our application of lemma 4.

Lemma 5: Suppose that G is an SDF graph, $M \subseteq N(G)$, $A_1 \in M$, and A_2 is an SDF node that is contained in $N(G)$ but not in M such that

- 1) A_2 is not adjacent to any member of $(M - \{A_1\})$, and
- 2) for some positive integer k , $q(A_2) = kq(A_1)$.

Then if we cluster $W = \{A_1, A_2\}$ in G , then $\text{subgraph}(M - \{A_1\} + \{\Omega\}, G')$ is isomorphic to $\text{subgraph}(M, G)$.

As a simple illustration, consider again the clustering example of Fig. 6(c) and (d). Let G and G' respectively denote the graphs of Fig. 6(c) and (d), and in Fig. 6(c), let $M = \{A, B\}$, $A_1 = A$, and $A_2 = C$. Then $(M - \{A_1\} + \{\Omega\}) = \{B, \Omega\}$, and clearly, $\text{subgraph}(\{B, \Omega\}, G')$ is isomorphic to $\text{subgraph}(\{A, B\}, G)$. The proof of lemma 5 can be found in the appendix.

Lemma 6: Suppose that G is a strongly connected SDF graph, and Z is a strongly connected subset of nodes in G such that $q_G(Z) = 1$. Suppose Z_1 and Z_2 are disjoint subsets of Z such that Z_1 is subindependent of Z_2 in $\text{subgraph}(Z)$. Then Z_1 is subindependent of Z_2 in G .

Proof: For each arc α directed from a member of Z_2 to a member of Z_1 , we have $\text{delay}(\alpha) \geq \text{total_consumed}(\alpha, \text{subgraph}(Z))$. From fact 3, $q_{\text{subgraph}(Z)}(N) = q_G(N)$ for all $N \in Z$. Thus, for all arcs α in $\text{subgraph}(Z)$, $\text{total_consumed}(\alpha, \text{subgraph}(Z)) = \text{total_consumed}(\alpha, G)$, and we conclude that Z_1 is subindependent of Z_2 in G . QED

Lemma 7: Suppose G is a strongly connected SDF graph, A and B are distinct nodes in G , and $W = \{A, B\}$ forms a proper subset of $N(G)$. Suppose also that the following conditions all hold:

- 1) Neither A nor B is contained in a tightly interdependent subgraph of G .
- 2) There is at least one arc directed from A to B that has no delay.
- 3) B has no predecessors other than A or B .
- 4) $q_G(B) = kq_G(C)$ for some $C \in N(G)$, $C \neq B$.

Then the SDF graph G' that results from clustering W is loosely interdependent.

Proof: From (1) G must be loosely interdependent, so there exist subsets X_1, X_2 of $N(G)$ such that $X_1 | G X_2$. If $A, B \in X_1$ or $A, B \in X_2$, then from lemma 4, we are done. Now condition (2) precludes the scenario ($B \in X_1, A \in X_2$), so the only remaining possibility is ($A \in X_1, B \in X_2$). There are two cases to consider here:

- i) B is not the only member of X_2 . Then from (3), $(X_1 + \{B\}) | G (X_2 - \{B\})$. But $A, B \in (X_1 + \{B\})$, so lemma 4 again guarantees that G' is loosely interdependent.
- ii) A is not the only member of X_1 and $X_2 = \{B\}$. Thus we have $X_1 | G \{B\}$, so

$$\forall \alpha \in A(G), (\text{source}(\alpha) = B) \Rightarrow \text{delay}(\alpha) \geq \text{total_consumed}(\alpha, G). \quad (1)$$

Also, since $C \in X_1$ we have from (4) that $q_G(X_1) = \gcd(\{q_G(N) | N \in X_1\}) = \gcd(\{q_G(N) | N \in X_1\} \cup \{kq_G(C)\}) = \gcd(\{q_G(N) | N \in X_1\} \cup \{q_G(B)\}) =$

$gcd(\{q_G(N) | N \in N(G)\}) = 1$. That is,

$$q_G(X_1) = 1. \quad (2)$$

Now if X_1 is not strongly connected, then it has a proper subset Z such that there are no arcs directed from a member of $(X_1 - Z)$ to a member of Z . Furthermore, from condition (3), $A \notin Z$. This is true because if Z contained A , then no member of $(X_1 - Z)$ would have a path to B , and thus G would not be strongly connected. Thus $A \in (X_1 - Z)$, and there are no arcs directed from $(X_1 - Z)$ to Z . So all arcs directed from $(X_1 - Z + \{B\})$ to Z have node B as their source. From (1) it follows that $Z | G(X_1 - Z + \{B\})$. Now $A, B \in (X_1 - Z + \{B\})$, so applying lemma 4 we conclude that G' is loosely interdependent.

If X_1 is strongly connected, we know from condition (1) that there exist Y_1, Y_2 such that $Y_1 | \text{subgraph}(X_1) Y_2$. From (2) and lemma 6, Y_1 is subindependent of Y_2 in G . Now if $A \in Y_1$, then from condition (3), B is subindependent of Y_2 in G , so from fact 6(a), $(Y_1 \cup \{B\}) | G Y_2$. Applying lemma 4, we see that G' is loosely interdependent. On the other hand, suppose that $A \in Y_2$. From (1), we know that Y_1 is subindependent of $\{B\}$ in G . From fact 6(b), it follows that Y_1 is subindependent of $(Y_2 \cup \{B\})$, so again we can apply lemma 4 to conclude that G' is loosely interdependent. QED

Theorem 4: Suppose G is a connected SDF graph, A and B are distinct nodes in G such that B is a successor of A , and $W = \{A, B\}$ is a proper subset of $N(G)$. If we cluster W in G then the tightly interdependent components of G' are the same as the tightly interdependent components of G if the following conditions all hold:

- 1) Neither A nor B is contained in a tightly interdependent component of G .
- 2) At least one arc directed from A to B has zero delay.
- 3) $q_G(B) = k q_G(A)$ for some positive integer k .
- 4) B has no predecessors other than A and B .

Proof: It suffices to show that all strongly connected subgraphs in G' that contain Ω are loosely interdependent. So we suppose that Z' is a strongly connected subset of $N(G')$ that contains Ω , and we let Z denote the “corresponding” subset in G ; that is, $Z = Z' - \{\Omega\} + \{A, B\}$. Now in Z' , suppose that there is a directed circuit $(C \rightarrow \Omega \rightarrow D \rightarrow C)$ containing the node Ω . From condition (4), this implies that there is a directed circuit in G containing A, C, D , and possibly B . The two possible ways in which a directed circuit in G introduces a directed circuit involving Ω in G' are illustrated in Fig. 8(a) and (b); the situation in (c) cannot arise because of condition (4).

Now in Z' , if one or more of the circuits involving Ω corresponds to Fig. 8(a), then Z must be strongly connected. Otherwise, all of the circuits involving Ω correspond to Fig. 8(b), so $(Z - \{B\})$ is strongly connected, and from condition (4), no member of $(Z - \{A, B\})$ is adjacent to B . In the former case, lemma 7 yields the loose interdependence of Z' .

In the latter case, lemma 5 guarantees that $(Z - \{B\})$ is isomorphic to Z' . Since $A \in (Z - \{B\})$, and since from condition

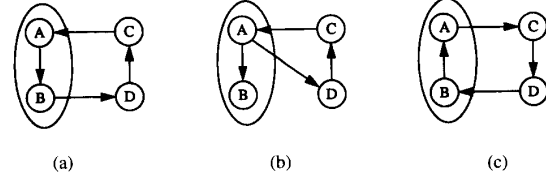


Fig. 8. An illustration of how a directed circuit involving Ω originates in G' for Theorem 4. The two possible scenarios are shown in (a) and (b); (c) will not occur due to condition (4). SDF parameters on the arcs have not been assigned because they are irrelevant to the introduction of directed cycles.

(1), A is not contained in any tightly interdependent subgraph of G , it follows that Z' is loosely interdependent. QED

If we assume that the input SDF graph has a single appearance schedule then we can ignore condition (1). From our observations, this is a valid assumption for the vast majority of practical SDF graphs. Also, condition (3) can be verified by examining any single arc directed from A to B ; if α is directed from A to B then condition (3) is equivalent to $p(\alpha) = kc(\alpha)$. In our current implementation, we consider only the case $k = 1$ for condition (3) because in practice, this corresponds to most of the opportunities for efficiently using registers.

We see that the clustering process defined by theorem 4—under the assumption that the original graph has a single appearance schedule—requires only *local* dataflow information, and thus it can be implemented very efficiently. If our assumption that a single appearance schedule exists is wrong, then we can always undo our clustering decisions. Since the assumption is frequently valid, and since it leads to a very efficient algorithm, this is the form in which we have implemented theorem 4. Finally, in addition to making data transfers more efficient, our clustering process provides a fast way to reduce the size of the graph without canceling the existence of a single appearance schedule. When used as a preprocessing technique, this can sharply reduce the execution time of a loose interdependence algorithm.

VII. CONCLUSIONS

This paper has presented fundamental topological relationships between iteration and looping in SDF graphs, and we have shown how to exploit these relationships to synthesize the most compact looping structures for a large class of applications. Furthermore, we have extended the developments of [5] by showing how to isolate the minimal subgraphs that require explicit deadlock detection schemes, such as the reachability matrix, when organizing hierarchy.

This paper also defines a framework for evaluating different scheduling schemes having different objectives, with regard to their effect on schedule compactness. The developments of this paper apply to any scheduling algorithm that imposes hierarchy on the SDF graph. For example, by successively repeating the same block of code, we can reduce “context-switch” overhead [22]. We can identify subgraphs that use as much of the available hardware resources as possible, and these can be clustered, as the computations to be repeatedly invoked. However, the hierarchy imposed by such a scheme must be evaluated against its impact on program compactness.

For example, if a cluster introduces tight interdependence, then it may be impossible to fit the resulting program on chip, even though the original graph had a sufficiently compact schedule.

The techniques developed in this paper have been successfully incorporated into a block-diagram software synthesis environment for DSP [18]. We are currently investigating how to systematically incorporate these techniques into other scheduling objectives—for example, how to balance parallelization objectives with program compactness constraints.

APPENDIX

This appendix contains proofs of some the lemmas that were stated and used in Sections IV–VI.

Proof of Lemma 1:

From remark 1, if N is not contained in a nontrivial strongly connected component of G , the result is obvious, so we assume, without loss of generality, that N is in some nontrivial strongly connected component H_1 of G . From our assumptions, $\text{subgraph}(H_1)$ must be loosely interdependent, so λ partitions H_1 into X and Y , where $X|_{\text{subgraph}(H_1)}Y$. Let H'_1 denote that connected component of $\text{subgraph}(X)$ or $\text{subgraph}(Y)$ that contains N . From remark 3, $\text{appearances}(N, S_\lambda(G)) = \text{appearances}(N, S_\lambda(\text{subgraph}(H'_1)))$.

From our assumptions, all nontrivial strongly connected subgraphs of H'_1 that contain N are loosely interdependent. Thus, if N is contained in a nontrivial strongly connected component H_2 of H'_1 , then λ will partition H_2 , and we will obtain a proper subset H'_2 of H'_1 such that $\text{appearances}(N, S_\lambda(\text{subgraph}(H'_1))) = \text{appearances}(N, S_\lambda(\text{subgraph}(H'_2)))$. Continuing in this manner, we get a sequence H'_1, H'_2, \dots of subsets of $N(G)$ such that each H'_i is a proper subset of H'_{i-1} , N is contained in each H'_i , and $\text{appearances}(N, S_\lambda(G)) = \text{appearances}(N, S_\lambda(\text{subgraph}(H'_1))) = \text{appearances}(N, S_\lambda(\text{subgraph}(H'_2))) = \dots$. Since each H'_i is a strict subset of its predecessor, we can continue this process only a finite number, say m , of times. Then $N \in H'_m$, N is not contained in a nontrivial strongly connected component of $\text{subgraph}(H'_m)$, and $\text{appearances}(N, S_\lambda(G)) = \text{appearances}(N, S_\lambda(\text{subgraph}(H'_m)))$. But from remark 1, $S_\lambda(\text{subgraph}(H'_m))$ contains only one appearance of N . QED

Proof of Lemma 4:

Let Φ denote the set of arcs directed from a node in X_2 to a node in X_1 , and let Φ' denote the set of arcs directed from a node in X'_2 to a node in X'_1 . Since $\text{subgraph}(\{A, B\})$ does not contain any arcs in Φ , it follows that $\Phi' = \{\alpha' | \alpha \in \Phi\}$. From fact 5, it can easily be verified that for all $\alpha', \text{total.consumed}(\alpha', G') = \text{total.consumed}(\alpha, G)$. Now since $X_1|GX_2$, we have $\forall \alpha \in \Phi, \text{delay}(\alpha') \geq \text{total.consumed}(\alpha, G)$. It follows that $\forall \alpha' \in \Phi', \text{delay}(\alpha') \geq \text{total.consumed}(\alpha', G')$. We conclude that X'_1 is subindependent of X'_2 in G' . QED

Proof of Lemma 5:

Let $C = \text{subgraph}(M - \{A_1\} + \{\Omega\}, G')$, let Φ denote the set of arcs in $\text{subgraph}(M, G)$, and let Φ' denote the set of arcs in C . From (1), every arc in C has a corresponding arc in $\text{subgraph}(M, G)$ and vice-versa, and thus $\Phi' = \{\alpha' | \alpha \in \Phi\}$. Now from the definition of clustering a subgraph, we know that $p(\alpha') = p(\alpha)$ for any arc $\alpha \in \Phi$ such that $\text{source}(\alpha) \neq A_1$. If $\text{source}(\alpha) = A_1$, then α is replaced by α' with $\text{source}(\alpha') = \Omega$, and $p(\alpha') = p(\alpha)q(A_1)/\text{gcd}(q(A_1), q(A_2))$. But $\text{gcd}(q(A_1), q(A_2)) = \text{gcd}(q(A_1), kq(A_1)) = q(A_1)$, so $p(\alpha') = p(\alpha)$. Thus $p(\alpha') = p(\alpha)$ for all $\alpha \in \Phi$. Similarly, we can show that $c(\alpha') = c(\alpha)$ for all $\alpha \in \Phi$. Thus, the mappings $f_1: M \rightarrow N(C)$ and $f_2: \Phi \rightarrow \Phi'$ defined by

$$\begin{aligned} f_1(N) &= N \text{ if } N \neq A_1, \\ f_1(A_1) &= \Omega; \quad \text{and} \quad f_2(\alpha) = \alpha' \end{aligned}$$

demonstrate that $\text{subgraph}(M, G)$ is isomorphic to C . QED

GLOSSARY

$Z_1|GZ_2$: If G is an SDF graph and Z_1 and Z_2 form a partition of the nodes in G such that Z_1 is subindependent of Z_2 in G , then we write $Z_1|GZ_2$.

$A(G)$: The set of arcs in the SDF graph G .

$\text{appearances}(N, S)$: The number of times that actor N appears in the looped schedule S .

admissible schedule: A schedule $S_1S_2 \dots S_k$ such that each S_i has sufficient input data to fire immediately after its antecedents $S_1S_2 \dots S_{i-1}$ have fired.

$c(\alpha)$: The number of samples consumed from SDF arc α by one invocation of $\text{sink}(\alpha)$.

$\text{delay}(\alpha)$: The number of delays on SDF arc α .

gcd : Greatest common divisor.

$N(G)$: The set of nodes in the SDF graph G .

PASS: A periodic admissible sequential schedule.

$p(\alpha)$: The number of samples produced onto SDF arc α by one invocation of $\text{source}(\alpha)$.

periodic schedule: A schedule that invokes each node at least once and produces no net change in the number of samples buffered on any arc.

predecessor: Given two nodes A and B in an SDF graph, A is a predecessor of B if there is at least one arc directed from A to B .

q_G : The *repetitions vector* q_G of the SDF graph G is a vector that is indexed by the nodes in G . q_G has the property that every *PASS* for G invokes each node N a multiple of $q_G(N)$ times.

single appearance schedule: A schedule that contains only one appearance of each actor in the associated SDF graph.

$\text{sink}(\alpha)$: The actor at the sink of SDF arc α .

$\text{source}(\alpha)$: The actor at the source of SDF arc α .

subgraph: A *subgraph* of an SDF graph G is the graph formed by any subset Z of nodes in G together with all arcs α in G

for which $source(\alpha), sink(\alpha) \in Z$. We denote the subgraph corresponding to the subset of nodes Z by $subgraph(Z, G)$, or simply by $subgraph(Z)$ if G is understood from context.

subindependent: Given an SDF graph G , and two disjoint subsets Z_1, Z_2 of nodes in G , we say that Z_1 is subindependent of Z_2 in G if for every arc α in G with $source(\alpha) \in Z_2$ and $sink(\alpha) \in Z_1$, we have $delay(\alpha) \geq total_consumed(\alpha, G)$. We say that Z_1 is subindependent in G if Z_1 is subindependent of $(N(G) - Z_1)$ in G .

successor: Given two nodes A and B in an SDF graph, A is a successor of B if there is at least one arc directed from B to A .

total_consumed(α, G): The total number of samples consumed from arc α in a minimal schedule period of the SDF graph G ; that is, $total_consumed(\alpha, G) = q_G(sink(\alpha))c(\alpha)$.

valid schedule: A schedule that is a PASS.

REFERENCES

- [1] A. V. Aho, J. E. Hopcroft, and J. D. Ullman, "The design and analysis of computer algorithms." Reading, MA: Addison-Wesley, 1974.
- [2] Arvind, L. Bic, and T. Ungerer, "Evolution of data-flow computers," in *Advanced Topics In Data-Flow Computing*, J. L. Gaudiot and L. Bic, Eds. Englewood Cliffs, NJ: Prentice-Hall, 1991.
- [3] S. S. Bhattacharyya, "Compiling dataflow programs for digital signal processing," Memo. No. UCB/ERL M94/52, Electronics Research Lab., College of Engineering, Univ. of California, Berkeley, CA, July 1994.
- [4] S. S. Bhattacharyya and E. A. Lee, "Looped schedules for dataflow descriptions of multirate DSP algorithms," Memo. No. UCB/ERL M93/36, Electronics Research Lab., College of Engineering, Univ. of California, Berkeley CA, May 1993.
- [5] ———, "Scheduling synchronous dataflow graphs for efficient looping," *J. VLSI Signal Process.*, vol. 6, no. 3, pp. 271–288, Dec. 1993.
- [6] J. B. Dennis, "First version of a dataflow procedure language," MIT/LCS/TM-61, MIT, Lab. for Computer Science, Cambridge, MA, 1975.
- [7] ———, "Stream data types for signal processing," unpublished memorandum, Sept. 28, 1992.
- [8] G. R. Gao, R. Govindarajan, and P. Panangaden, "Well-behaved programs for DSP computation," in *ICASSP*, San Francisco, CA, Mar. 1992.
- [9] D. Genin, J. De Moortel, D. Desmet, and E. Van de Velde, "System design, optimization, and intelligent code generation for standard digital signal processors," in *ISCAS*, Portland, OR, May 1989.
- [10] P. N. Hilfinger, "Silage reference manual, draft release 2.0," Computer Science Division, EECs Dept., Univ. of California, Berkeley, July 1989.
- [11] W. H. Ho, E. A. Lee, and D. G. Messerschmitt, "High level dataflow programming for digital signal processing," in *VLSI Signal Processing III*. Piscataway, NJ: IEEE Press, 1988.
- [12] S. How, "Code generation for multirate DSP systems in Gabriel," Memo. No. UCB/ERL M94/82, Electronics Research Lab., College of Engineering, Univ. of California, Berkeley, CA, Oct. 1994.
- [13] E. A. Lee, "Static scheduling of dataflow programs for DSP," in *Advanced Topics in Dataflow Computing*, J. L. Gaudiot and L. Bic, Eds. Englewood Cliffs, NJ: Prentice-Hall, 1991.
- [14] E. A. Lee and D. G. Messerschmitt, "Static scheduling of synchronous dataflow programs for digital signal processing," *IEEE Trans. Comput.*, vol. C-36, no. 1, pp. 24–35, Jan. 1987.
- [15] E. A. Lee and D. G. Messerschmitt, "Synchronous dataflow," *Proc. IEEE*, vol. 75, no. 9, pp. 1235–1245, Sept. 1987.
- [16] J. R. McGraw, S. K. Skedzielewski, S. Allan, D. Grit, R. Oldehoft, J. Glauert, I. Dobes, and P. Hohensee, "SISAL: Streams and iteration in a single assignment language," Language Reference Manual, Version 1.1., Lawrence Livermore National Laboratory, Livermore, CA, July 1983.
- [17] D. R. O'Hallaron, "The ASSIGN parallel program generator," Memo. No. CMU-CS-91-141, School of Computer Science, Carnegie Mellon Univ., Pittsburgh, PA, May 1991.
- [18] J. L. Pino, S. Ha, E. A. Lee, and J. T. Buck, "Software synthesis for DSP using ptolemy," to be published in *J. VLSI Signal Process.*, vol. 9, no. 1, pp. 7–21, Jan. 1995.
- [19] D. B. Powell, E. A. Lee, and W. C. Newmann, "Direct synthesis of optimized DSP assembly code from signal flow block diagrams," in *ICASSP*, San Francisco, CA, Mar. 1992, pp. 553–556.

- [20] H. Printz, "Automatic mapping of large signal processing systems to a parallel machine," Memo. No. CMU-CS-91-101, School of Computer Science, Carnegie-Mellon Univ., Pittsburgh, PA, May 1991.
- [21] S. Ritz, M. Pankert, and H. Meyr, "High level software synthesis for signal processing systems," in *Proc. Int. Conf. Applicat. Specific Array Processors*, Berkeley, CA, Aug. 1992, pp. 679–693.
- [22] S. Ritz, M. Pankert, and H. Meyr, "Optimum vectorization of scalable synchronous dataflow graphs," in *Proc. Int. Conf. Applicat. Specific Array Processors*, Venice, Oct. 1993, pp. 285–296.
- [23] G. Sih, "Multiprocessor scheduling to account for interprocessor communication," Memo. No. UCB/ERL M91/29, Electronics Research Lab., Univ. of California, Berkeley, Apr. 1991.
- [24] R. E. Tarjan, "Depth first search and linear graph algorithms," *SIAM J. Computing*, vol. 1, no. 2, pp. 146–160, June 1972.
- [25] W. W. Wadge and E. A. Ashcroft, *Lucid, the Dataflow Language*. New York: Academic, 1985.



Shuvra S. Bhattacharyya (S'92–M'93) received the B.S. degree in electrical and computer engineering from the University of Wisconsin, Madison, in 1987, and the M.S. and Ph.D. degrees from the University of California, Berkeley, in 1991 and 1994, respectively.

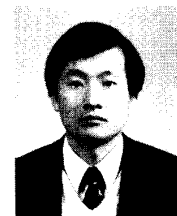
From 1991 to 1992, he was employed by Kuck and Associates, Champaign, Illinois, where he designed and implemented optimizing program transformations for C and Fortran compilers. Since July 1994, he has been a researcher in the Semiconductor Research Laboratory at Hitachi America, Ltd., San Jose, CA. His current research interests include software, architectures, and rapid prototyping for digital signal processing; VLSI signal processing; and parallel computation.

Dr. Bhattacharyya has published several papers, and he is a member of the Association for Computing Machinery (ACM).



Joseph T. Buck received the B.E.E. degree from Catholic University of America in 1978, and the M.S. in computer science from George Washington University in 1981. He received the Ph.D. in 1993 from the University of California, Berkeley, where he was one of the main designers for Ptolemy, a design, simulation, and prototyping environment for heterogeneous systems.

From 1979 to 1984 he participated in research in speech coding and recognition at the Naval Research Laboratory. From 1984 to 1989 he worked at Entropic Speech, Inc. on real-time implementations of speech compression algorithms for telephony applications. Presently, he is a staff research engineer in the Advanced Technology Group of Synopsys, Inc. His research interests include techniques for producing efficient hardware, software, and mixed implementations from dataflow graphs and other high-level representations of algorithms.



Soonhoi Ha (S'87–M'92) received the B.S. and M.S. degrees in electronics from Seoul National University, Seoul, Korea, in 1985 and 1987, respectively. He received the Ph.D. degree in the Electrical Engineering and Computer Science Department at the University of California, Berkeley, in 1992.

Currently, he is a full-time lecturer in the Computer Engineering Department at Seoul National University, Seoul, Korea. His research interests include architecture and scheduling techniques for parallel processing and design methodology for digital systems.

Dr. Ha is a member of the Association for Computing Machinery (ACM) and the IEEE Computer Society.



Edward A. Lee (S'80–M'86–SM'93–F'94) received the B.S. from Yale University in 1979, the M.S. from the Massachusetts Institute of Technology in 1981, and the Ph.D. from the University of California, Berkeley, in 1986.

From 1979 to 1982, he was a member of technical staff at Bell Telephone Laboratories in Holmdel, NJ, in the Advanced Data Communications Laboratory. At present, he is a Professor in the Electrical Engineering and Computer Science Department at the University of California, Berkeley. His research activities include real-time software, parallel computation, architecture and software techniques for signal processing, and design methodology for heterogeneous systems. He is Director of the Ptolemy project at UC Berkeley, and previously directed the Gabriel project. He is a founder of Berkeley Design Technology, Inc. and has consulted for a number of other companies. He is co-author of *Digital Communication* (Kluwer Academic Press, 1988 first ed., 1994 second ed.), and co-author of *Digital Signal Processing Experiments* (Prentice-Hall, 1989), as well as numerous technical papers.

Dr. Lee was recently Chairman of the VLSI Technical Committee of the Signal Processing Society, and Co-Program Chair of the 1992 Application Specific Array Processor Conference. He is an Associate Editor of *Design Automation for Embedded Systems* and is on the editorial board of the *Journal on VLSI Signal Processing*. As a Fellow of the IEEE, he has received the citation "For contributions to design methodologies and programming techniques for real-time digital signal processing systems." He was a recipient of a 1987 NSF Presidential Young Investigator award, an IBM faculty development award, the 1986 Sakrison Prize at U.C. Berkeley, and a paper award from the IEEE Signal Processing Society.