

Papers

Optimum and Heuristic Transformation Techniques for Simultaneous Optimization of Latency and Throughput

Mani B. Srivastava, *Member, IEEE*, and Miodrag Potkonjak, *Member, IEEE*

Abstract—Although throughput alone can be arbitrarily improved for several classes of systems using previously published techniques, none of those approaches are effective when latency constraints, which are increasingly important in embedded DSP systems, are considered. After formally establishing the relationship between latency and throughput in general computation, we explore the effect of pipelining on latency, and establish necessary and sufficient conditions under which pipelining does not alter latency.

Many systems are either linear, or have subsystems that are linear. For such cases we have used a state-space based approach that treats various transformations in an integrated fashion, and answers analytically whether it is possible to simultaneously meet any given combination of constraints on latency and throughput. The analytic approach is constructive in nature, and produces a complete implementation when feasibility conditions are fulfilled. We also present a suboptimal but hardware efficient heuristic approach for the special case of initially-relaxed single-input single-output linear time-invariant computations. A novel software platform consisting of a high-level synthesis system coupled to a symbolic algebra system was used to implement the proposed algorithm transformations. Instead of optimizing to improve throughput and latency, our transformations can also be used to increase the implementation efficiency while achieving the same latency and throughput as the original design.

Index Terms—High-level synthesis, transformational synthesis, performance, and cost driven architectural synthesis.

I. INTRODUCTION

A. Throughput and Latency in System Design

The “open-loop” nature of tasks such as filtering has made throughput the commonly used performance metric for implementation of DSP computations. Algorithm transformation techniques, such as lookahead, unfolding, retiming, and algebraic manipulations, have been developed to restructure the original computation algorithm to a functionally equivalent

algorithm so that a VLSI or software implementation can achieve higher throughput—in some cases even arbitrarily high throughput—although often at the cost of increased implementation cost and latency. However, increasingly DSP-like computation is being used in systems where both throughput and latency are of importance. One example is that of DSP subsystems that are used in hard real-time embedded control systems where the reactive nature of the system requires it to not only keep up with the rate of arrival of sensory data, but also to generate the corresponding output within a certain period of time. A second example of both latency and throughput being important is that of “signal processing servers” where one would like to maximize the throughput of the server while minimizing the latency seen by any single client.

Although throughput alone can be improved using previously published techniques, none of them is effective when latency constraints are simultaneously considered. Algorithm transformations that simultaneously address latency and throughput are the focus of this work.

The terms latency and throughput are used with a variety of different connotations in the high level synthesis literature. In order to facilitate inter-disciplinary research mandated by embedded systems, we use the following well-established definitions from areas such as control theory and digital signal processing. *Latency* is the physical (wall-clock) time needed to deliver the output data from the moment of arrival of the corresponding input data. *Throughput* of an embedded system is the highest rate at which it can receive and process the input data. The inverse of throughput is called *Sample Period*. In Section III we will recast these definitions in standard high level synthesis and system level design terminology.

B. Previous Work

Algorithm transformation is a powerful optimization technique used in a number of scientific and engineering areas. Transformations have been used in computer science research for algorithm design [1], computer algebra [2], [3], and code optimization in compilers [4]. More relevant to us, transformations have been used in control theory and DSP to derive

Manuscript received September 20, 1993; revised February 25, 1994, May 18, 1994, and July 6, 1994.

M. B. Srivastava is with AT&T Bell Laboratories, Murray Hill, NJ 07974 USA.

M. Potkonjak is with NEC C&C Research Laboratories, Princeton, NJ 08540 USA.

IEEE Log Number 9408367.

a variety of functionally equivalent, but structurally different, computation schemes for generic tasks such as filtering [5]. More recently, influenced by the dominance of performance as the key design metric, significant importance has been placed on unfolding techniques coupled with various algebraic transformations [6], [7] to achieve arbitrarily high levels of throughput at the expense of additional hardware cost, but without taking latency into consideration. Transformations have also gained importance in high level synthesis of custom VLSI's where they have been used for improving a variety of design metrics [8], [9], [10], [11], [12]. Sophisticated transformations have also been presented for functional [13] and software [14] pipelining.

C. What is New?

We have developed rigorous graph theoretic definitions of latency and throughput for general synchronous computation that also expose the interdependence between these two metrics of speed. Since throughput alone has been extensively treated in literature, we first establish the properties of latency in isolation by studying its exact relationship to pipelining, and by developing a generic state-space based approach to latency minimization that efficiently employs retiming, pipelining, and all algebraic and redundancy manipulation transformations.

Next we study throughput and latency together while paying special attention to the widely used class of *Linear Time-Invariant* (LTI) systems. For this class of systems we have developed a novel transformation technique that restructures the initial computation algorithm to one that has provably optimal latency and throughput, even when no assumption are made about the initial state and values of the coefficients. Until now all approaches which were able to improve throughput to an arbitrary extent were based on a combination of unfolding of the computation with block processing or interleaving [15], [7], [16]. However, this comes at the expense of a proportional degradation in latency. This long-standing *Latency and Sample Period Bottleneck* is broken by employing a novel combination of unfolding with “*On-Arrival Processing*” where input samples are processed as soon as they arrive, and a provably optimal latency minimization technique. For the special case of single-input LTI systems with zero initial state, we also present cost efficient transformation techniques that produce optimal, or close to optimal, combinations of latency and throughput.

Our transformation techniques for latency and throughput optimization can be used in two distinct fashion in a design environment. First, and the obvious, use is to transform an initial algorithm that, even given any amount of hardware resources, cannot meet constraints on latency and throughput to a new algorithm that can satisfy the constraints. The second, application of these transformations is to improve the implementation cost even when the initial algorithm can meet the constraints on throughput and latency. In this scenario the transformations are used to obtain a new algorithm with even better latency and throughput characteristics (shorter critical paths) and without increasing the amount of computation and storage too much. This new algorithm gives the scheduler more flexibility in meeting the timing constraints, which can often

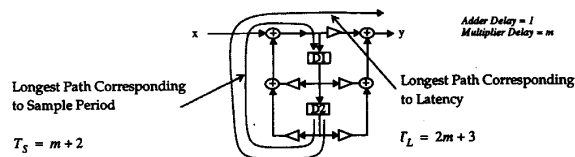


Fig. 1. CDFG for a biquad filter with paths corresponding to latency T_L and sample period T_S .

result in better resource utilization and lower implementation cost.

II. CDFG AND STATE-SPACE REPRESENTATIONS

A. System Representation by CDFG

The systems that we are interested in have multiple inputs, multiple outputs, and finite state. They accept streams of samples on each of the inputs, and produce streams of samples on each of the output ports. We represent an algorithm for a system by a hierarchical directed control-dataflow graph (CDFG). In a CDFG the nodes represent data operators or subgraphs, data edges represent the flow of data between nodes, and control edges represent sequencing and timing constraints between nodes.

We restrict ourselves to operators that are synchronous in that they consume at every input, and produce at every output, a fixed number of samples on every execution. This restriction has two interesting ramifications. First, the operators, and hence the system, are *determinate* in that a given set of input samples always results in the same outputs independent of the execution times. Second, the system is *well behaved* in that the data sample rate at any given data edge in the CDFG is independent of the inputs, and the ratio between any two data sample rates will be a statically known rational number [17]. Mathematically, such a synchronous CDFG is equivalent to a continuous function over streams of data samples [18]. Since CDFG's are of course causal, this means that they are equivalent to a function that expresses the i th set of output samples in terms of the i th and earlier sets of inputs samples.

The system state is represented in a CDFG by special delay operator nodes which are initialized to a user specified value. A delay operator node (often referred to as just *delay* or *state* in this paper) delays by one sample the stream of data on its sole input port.

A system¹ is completely represented by a CDFG and the initial values for all the delays in the CDFG. We further restrict ourselves to single-rate systems where the data rate is identical on all the inputs, and also assume that the i th samples on all inputs arrive simultaneously, aligned to a “sample clock.” In the rest of the paper we use just the term CDFG to refer to such a “single-rate synchronous CDFG.” Fig. 1 shows an example CDFG.

¹We use the term “system” interchangeably to mean “the system behavior specification,” “an algorithm that realizes the system behavior specification,” “an implementation of the system” in different places—the three are not equivalent. A CDFG is a representation of a specific algorithm. There can be multiple algorithms (CDFG's) that express the specified behavior, and there can be multiple implementations (hardware allocation and schedule) of an algorithm (CDFG).

obvious that the i th output sample had no causal relationship to the i th set of input samples.

A. Throughput and Latency Achieved by a CDFG

The preceding definitions of latency and throughput (sample period) can be recast in standard high level synthesis and system level design terminology. In particular, the definitions can be used to find the latency and sample period that can be achieved by a particular CDFG (no transformations allowed to the CDFG) given any amount of hardware resources and using any implementation technique. We will show that in the general case the latency and sample period that can be achieved by a CDFG forms a range of solutions that is parameterized by certain variables that an implementation is in general free to choose.

To get an intuitive feeling, let us first consider the simple case where the implementation is restricted to those where, in terminology of the state-space representation of (1), the P elements of the n th input vector $X[n]$ and the Q elements of the corresponding previous state vector $S[n-1]$ are all available at the same time—in other words, there is no time skew between them. These elements are used to calculate the n th output vector $Y[n]$ and the n th state vector $S[n]$. While this implementation model may sound restrictive, it is in fact the one that is assumed by the scheduler in most of the popular high-level synthesis tools such as HYPER [9]. Latency T_L of a particular primary output is the number of control steps that are needed to generate an output sample since the arrival of corresponding input samples. In our restrictive implementation model where the previous set of state values is available at the same time as the arrival of the new set of input values, latency T_L will be equivalent to kT_S + the length of the longest computation path from any primary input or state to the primary output. The path lengths are measured in control steps, and k is the number of pipeline stages that have been added (or removed if k is negative) to the initial specification. In a similar vein, T_S is the length of the longest computation path measured in control steps from any primary input or state to any state. Fig. 1 illustrates the longest paths corresponding to latency and sample period for a biquad filter.

Having gained an intuitive understanding of latency and sample period achieved by a CDFG under a simple implementation model, we will now develop exact expressions for achievable latency and sample period under a more general implementation model where we remove the restriction that all the elements of $X[n]$ and $S[n-1]$ are available at the same time.

Consider a CDFG with P inputs, Q outputs, and R state nodes. Using the same notation as in (1), let:

$P_{IS}(i, j)$ = length of the path (in control steps) from the i th primary input node to the j th state node

$P_{IS}(i, j)$ = length of the path (in control steps) from the i th primary input node to the j th state node

$P_{IO}(i, j)$ = length of the path (in control steps) from the i th primary input node to the j th primary output node

$P_{SS}(i, j)$ = length of the path (in control steps) from the i th state node to the j th state node

$P_{SO}(i, j)$ = length of the path (in control steps) from the i th state node to the j th primary output node
 k = number of pipeline stages that have been added (or removed if $k < 0$) to the initial CDFG that was given by the user as a specification, and from which the current CDFG under consideration was obtained after some transformations

$T_{IA}(i)$ = skew in the arrival of the i th element of $X[n]$ (n th sample at the i th input) relative to arrival of the first element of $X[n]$ (n th sample at the first input)

Note that by definition $T_{IA}(1) = 0$.

$T_{SA}(i)$ = skew in the arrival of the i th element of $S[n-1]$ ($(n-1)$ th value at the i th state node) relative to arrival of the 1st element of $X[n]$ (n th sample at the 1st input)

T_S = sample period

$T_L(i, j)$ = latency from i th input node to the j th output node

Note that the input arrival time skews, $T_{IA}(i) \forall i \in 1 \dots P$, are usually timing constraints that are in general specified by the user. On the other hand, the skews in state arrival, $T_{SA}(i) \forall i \in 1 \dots R$, are parameters that an implementation is free to choose so as to satisfy design constraints while optimizing design cost metrics.

A little thought shows that:

$$T_S = \max\{P_{SS}(r, s) + T_{SA}(r) - T_{SA}(s) \forall r, s \in 1 \dots R\} \\ \cup \{P_{IS}(p, r) + T_{IA}(p) - T_{SA}(r) \forall p \in 1 \dots P, \\ \forall r \in 1 \dots R\}$$

$$T_L(i, j) = k * T_S + \max\{P_{IO}(p, j) + T_{IA}(p) - T_{IA}(i) \\ \forall p \in 1 \dots P\} \\ \cup \{P_{SO}(r, j) + T_{SA}(r) - T_{IA}(i) \forall r \in 1 \dots R\} \\ \forall i \in 1 \dots P \text{ and } \forall j \in 1 \dots Q.$$

The above expressions clearly show that the achievable values of T_S and $T_L(i, j)$ are coupled in general—choosing some may place constraints on the achievable values of the remainder.

As mentioned earlier, most high level synthesis systems assume an implementation model where $T_{IA}(i) = 0 \forall i \in 1 \dots P$, and $T_{SA}(i) = 0 \forall i \in 1 \dots R$. However, in most of this paper we use a more permissive implementation model where $T_{IA}(i) = 0 \forall i \in 1 \dots P$, $T_{SA}(i) = T_{SA} \forall i \in 1 \dots R$, and T_{SA} is a single implementation timing parameter². Intuitively, this model assumes that the i th sample arrives at the same time for all the inputs, and the previous state values arrive at the same time for all the state nodes, and that the arrival of the state values is skewed with respect to the arrival of the input values by the time interval T_{SA} . We found that the ability to choose the state arrival skew T_{SA} enables combinations of latency and throughput to be achieved that are otherwise not possible. This parameter plays a key role in our techniques in Section VI-C. For this model the expressions for T_S and $T_L(i, j)$ get

²The parameter T_j used in analysis later in the paper (Section VI in particular) is similar to the parameter T_{SA} here.

simplified to (note that $T_L(i, j)$ is independent of i):

$$\begin{aligned} T_S &= \max\{P_{SS}(r, s) \forall r, s \in 1 \cdots R\} \\ &\quad \cup \{P_{IS}(p, r) - T_{SA} \forall p \in 1 \cdots P, \forall r \in 1 \cdots R\} \\ T_L(i, j) &= k * T_S + \max\{P_{IO}(p, j) \forall p \in 1 \cdots P\} \\ &\quad \cup \{P_{SO}(r, j) + T_{SA} \forall r \in 1 \cdots R\} \\ &\quad \forall i \in 1 \cdots P \quad \text{and} \quad \forall j \in 1 \cdots Q. \end{aligned}$$

IV. PIPELINING—LATENCY RELATIONSHIP

(Functional) pipelining [13] an implementation technique where several consecutive instances of the program body are partially overlapped during execution so as to improve the throughput. We use the following definition of pipelining: *Pipelining* with k pipeline stages on a CDFG is a special form of retiming where on each primary output (or input) k new delays are introduced, and only those delays can be moved. Increased latency is the most often quoted as an unavoidable harmful side effect of pipelining. However, as we will demonstrate in the rest of this section, if pipelining is done in a particular manner, this effect can be fully avoided. We assume that the user wants to introduce k delays, which will partition the CDFG into $k + 1$ pipeline stage. The length of the longest path that is relevant for latency calculation (i.e., the longest path from any of the delay nodes or primary inputs to the output) is denoted by T_L . Let T_S be the sample period. The following theorem establishes conditions under which pipelining can be done without altering latency:

Pipelining—Latency Theorem: A necessary condition that pipelining does not increase latency is that k , the number of pipeline delays introduced, satisfies $k \leq \lfloor T_L/T_S \rfloor$. Further, if micropipelining (fine-grained pipelining) of individual operations is allowed, or if all operations have equal delays, then the above condition is also a sufficient one.

Proof: The latency T_{LP} of the pipelined CDFG will be at least kT_S . If the condition specified in the theorem does not hold, i.e., if $k \geq \lfloor T_L/T_S \rfloor + 1$ (k is an integer), then $T_{LP} \geq kT_S \geq (\lfloor T_L/T_S \rfloor + 1)T_S > (T_L/T_S - 1 + 1)T_S = T_L$, which contradicts the hypothesis that latency will not increase.

In the case of operators that can be pipelined at a fine grain, as well as in the case of operators with identical delays, the following placement of pipeline delays (latches) will maintain the latency. Introduce the k pipeline delays one by one such that the i th delay, $\forall i \in \{1 \cdots k\}$, is placed at a distance of $i \times T_S$ from the input—in other words, the pipeline delays are placed separated by T_S at $T_S, 2T_S, 3T_S, \dots, kT_S$. This placement of pipeline delays will result in a latency of $k \times T_S + (T_L - k \times T_S) = T_L$ because it takes time $k \times T_S$ for an input sample to travel to the k th pipeline delay, and a further time of $T_L - k \times T_S$ to travel from the k th pipeline delay to the output. Such a placement is always feasible if fine-grained pipelining of operators is allowed, or if all operators have equal delays (in which case T_S will be a multiple of operator delay). On the other hand if neither of these prerequisites are not met then one may not be able to place the pipeline delays such that they are separated exactly by T_S . Observe that one cannot place adjacent pipeline delays at a distance greater than T_S

because that will increase the sample period, and one should not place them at a distance smaller than T_S because a time interval of T_S is always devoted to the computation between adjacent pipeline delays—placing adjacent delays closer than T_S will therefore result in an increase in the latency.

V. TECHNIQUES TO REDUCE LATENCY OF A SYSTEM

Latency that can be achieved by a CDFG is determined by the length of the longest combinational path that originates at a state node, or at a primary input, and ends at the primary output of interest. Therefore, intuitively, to minimize latency one should try to reduce the length of such paths by suitably transforming the initial CDFG.

To reduce the length of any computation path going from a primary input to a primary output, transformations based on algebraic properties can help. For example, a chain of n adders that adds $n + 1$ variables can be transformed into a maximally balanced binary tree of adders of depth $\lceil \log_2(n + 1) \rceil$. Retiming cannot help in reducing such input-output paths because they do not have any delay nodes to begin with. Introduction of pipeline stages is possible but this can never reduce the latency either.

More options are available to reduce the length of computation paths that originate at state nodes and end at primary outputs. Algebraic transformations will of course help, but another technique that helps is to retime such that the state nodes are moved closer to the primary outputs. Intuitively this means that we precalculate as much of the contribution of the previous state values to a primary output as possible. However, retiming and algebraic transformations effect each other mutually so that this is not a straightforward optimization. Further, retiming to minimize latency may result in the sample period getting worse because computation paths originating at primary inputs or at state nodes, and ending at state nodes, may be elongated.

Fortunately, in certain special but important cases an analysis in the state space framework offers more insight than the CDFG framework. This is because in the state space framework, as shown in (1), all the algebraic properties are encapsulated into two functions—the *state-transition mapping* ϕ , and the *output mapping* ρ . Retiming can be viewed in the state space as mapping the old state vector in to a new one such that there is no change in the behavior of the system observable at the outputs. This mapping of state vector results in new ϕ and ρ being defined. One can thus analyze the effects of retiming and algebraic transformations in an integrated fashion. For example, consider the case where ρ is separable such that

$$\rho(S[n-1], X[n], n) = \rho1(S[n-1], n) \oplus \rho2(S[n-1], X[n], n)$$

and where the time taken to calculate $\rho2$ and \oplus in sequence is less than that taken to calculate ρ alone. Many important systems, such as Volterra filters and all linear systems, have ρ that exhibit such separability. One can retime these systems such that the term $\rho1(S[n-1], n)$ is available as a state, and as a result the latency is reduced. The new state equations will

then be:

$$\begin{aligned} S[n] &= \phi(S[n-1], X[n], n, S[-1]), \\ \tilde{S}[n] &= \rho 1(\phi(S[n-1], X[n], n, S[-1]), n+1) \\ Y[n] &= \tilde{S}[n-1] \oplus \rho 2(S[n-1], X[n], n) \end{aligned}$$

$$\begin{aligned} n \in \{0, 1, 2, 3, \dots\} \quad X[n] \in \mathbb{R}^{P \times 1} \quad S[n] \in \mathbb{R}^{R \times 1} \\ \tilde{S}[n] \in \mathbb{R}^{1 \times 1} \quad Y[n] \in \mathbb{R}^{Q \times 1}. \end{aligned}$$

However, this can increase the sample period because the time taken to calculate $\rho 1(\phi(\dots), \dots)$ will in general be greater than that taken to calculate ϕ alone.

VI. LATENCY AND SAMPLE PERIOD FOR LINEAR TIME—INVARIANT SYSTEMS

CDFG's that are composed of variable-variable additions and variable-constant multiplications correspond to linear time-invariant systems. Although much work has been done toward algorithms for these systems that achieve low sample periods (high throughput), no work has been directed toward reducing both latency and sample period. In fact, some of the techniques used for improving the sample period end up making the latency worse. In this section we analyze linear time-invariant systems from the perspective of both latency and sample period, and present transformations that simultaneously address these two metrics.

A. Latency and Sample Period of a Linear Time-Invariant CDFG

Consider a CDFG with P primary inputs, Q primary outputs, R state nodes, and composed of additions of two variables, and multiplications of variables with constant coefficients. Its state-space representation is as shown in (2) where $X[n]$ is the current value of primary inputs in a vector form, $S[n-1]$ is the vector corresponding to the state values from the previous time step, and $Y[n]$ and $S[n]$ are the values of the output vector and state vector calculated in the current time step. Further, we assume the restricted implementation model (see discussion in Section III-A) where all the elements of $X[n]$ and $S[n-1]$ are available at the same time.

A key advantage of the state-space representation is that the specific organization of linear computation in the CDFG is abstracted by the two linear equations, $S[n] = AS[n-1] + BX[n]$ (called the state update equation) and $Y[n] = CS[n-1] + DX[n]$ (called the output equation), which encapsulate all the algebraic information. From the definitions of sample period T_S and latency T_L in Section III, it follows that T_S is decided by the time taken to compute $S[n] = AS[n-1] + BX[n]$, and T_L is decided by the time taken to compute $Y[n] = CS[n-1] + DX[n]$. These matrix equations are equivalent to a set of equations where the right hand sides are linear combinations of the P primary inputs and the R state values.

Noting that one of the maximally fast ways to evaluate a linear combination is by first doing the constant-variable multiplications in parallel, and then organizing the additions as a maximally balanced binary tree, we get the following expressions

for the best sample period and latency that can be achieved when the four constant coefficient matrices are nontrivial (i.e., if we do not assume that any elements may be 0 to 1 or -1):

$$\begin{aligned} T_S &= m + \lceil \log_2(R+P) \rceil \quad \text{and} \\ T_L &= m + \lceil \log_2(R+P) \rceil \end{aligned} \quad (3)$$

where we assume that the time for an addition is 1 and the time for a multiplication is $m \geq 1$.³

These expressions, however, can in general be quite pessimistic because the four coefficient matrices often have many elements that are 0 or 1 or -1 , in which case one can take advantage of such coefficients to reduce the number of adders and multipliers that are required. If we define $nr0(M, r)$ and $nr1(M, r)$ to be the number of elements that have magnitude 0 and 1 respectively in the r th row of matrix M , then the following exact expressions for best achievable latency and sample period are obtained:⁴

$$\begin{aligned} T_S &= \max_{r \in \{1 \dots R\}} (m + \lceil \log_2(R+P - nr0(A, r) - nr0(B, r) \\ &\quad - (nr1(A, r) - nr1(B, r))(1 - 1/2^m) \rceil) \\ T_L &= \max_{r \in \{1 \dots Q\}} (m + \lceil \log_2(R+P - nr0(C, r) - nr0(D, r) \\ &\quad - (nr1(C, r) + nr1(D, r))(1 - 1/2^m) \rceil). \end{aligned} \quad (4)$$

Since the number of inputs P and number of outputs Q are fixed, the goal of an algorithm transformation that optimizes latency and sample period will intuitively be to reduce R and increase the number of 0, 1, and -1 valued entries across all rows of the coefficient matrices.

B. Transforming a Linear Time-Invariant CDFG for Minimum Latency

From the state-space representation of a LTI CDFG it becomes obvious that no amount of retiming and algebraic transformations will ever change the matrix D . Retiming has no effect because D corresponds to paths in the LTI CDFG that go from input nodes to output nodes without passing through state nodes—there are no state nodes to retime with. Algebraic transformations have no effect because D is the matrix representation of the set of linear expressions that is implied by those input-output paths, and the matrix representation of a set of linear expressions is unique and unaffected by the application of commutativity, distributivity, associativity etc. Therefore, from (4) for T_L it follows that if one is able to transform the algorithm such that every row of the transformed version of matrix C has one entry with value 1 and all other entries with value 0, then the value of T_L is the minimum possible. Intuitively, the new algorithm will be such that each output depends on one and only one state variable, and that too through a coefficient of 1. We

³Note that this is a reasonable assumption because the latency of a multiplication operation is in almost all cases—ASIC's or programmable processors—larger than or equal to that of an addition operation.

⁴The unusual looking $(1 - 1/2^m)$ terms trace their origin to the fact that during the first m time units the variables with coefficients of magnitude 1 can be begun to be added in a maximally fast fashion in parallel with the multiplication of the other variables with their nontrivial (not 0, 1, or -1) coefficients.

show below that any LTI CDFG can indeed be transformed to such a minimum latency realization. Note that this solution is not unique, and is not necessarily the most efficient one either, because there exist infinitely many other solutions with minimum latency. For example, one can obtain other solutions with minimum latency by first applying some arbitrary state space transformation (a different state encoding) such that one still has R state values (infinitely many such transformations are possible because such a transformation is equivalent to an invertible $R \times R$ matrix; please refer to any linear systems books [19], and then apply the outlined algorithm.

Consider a transformation of the original algorithm such that not only do we have the original R state variables, but also Q linear combinations of those R state variables that are obtained by taking the inner product of the original state vector with the rows of matrix C . If these Q new states are denoted by \tilde{S} , then the following state space representation is obtained which is equivalent to the original system in input-output behavior, and is a minimum latency realization because each output depends on one and only one state via a coefficient of value 1:

$$\begin{aligned} \begin{bmatrix} S[n] \\ \tilde{S}[n] \end{bmatrix} &= \begin{bmatrix} A & 0_{R \times Q} \\ CA & 0_{Q \times Q} \end{bmatrix} \begin{bmatrix} S[n-1] \\ \tilde{S}[n-1] \end{bmatrix} + \begin{bmatrix} B \\ CB \end{bmatrix} X[n] \\ Y[n] &= [0_{Q \times R} \quad I_{Q \times Q}] \begin{bmatrix} S[n-1] \\ \tilde{S}[n-1] \end{bmatrix} + DX[n] \end{aligned} \quad (5)$$

where:

$$\begin{bmatrix} S[-1] \\ \tilde{S}[-1] \end{bmatrix} = \begin{bmatrix} I_{R \times R} \\ C \end{bmatrix} S[-1].$$

In the general case where the initial coefficient matrices are nontrivial, i.e., the matrix elements are not trivially 0 or 1, the above transformation guarantees that given enough hardware one can always achieve the following sample period and latency:

$$\begin{aligned} T_S &= m + \lceil \log_2(R + P) \rceil \quad \text{and} \\ T_L &= m + \lceil \log_2(1 + P) \rceil \end{aligned} \quad (6)$$

where the time for an addition is 1 and the time for a multiplication is m . For example, using this minimum latency transformation on a single input system, a sample period better than $m + \lceil \log_2(R + 1) \rceil$ and a latency better than $m + 1$ time units can always be achieved.

C. Transforming a Linear Time-Invariant CDFG to Jointly Improve Latency & Sample Period

Previous research [6], [7] has shown that one can use unfolding, look-ahead, and block-processing techniques to arbitrarily improve the sample period of LTI systems. On the other hand we just showed that one can always transform LTI systems to attain the minimum possible latency $T_L = m + \lceil \log_2(1 + P) \rceil$. Can we combine the two techniques to arbitrarily improve the sample period and simultaneously achieve the minimum latency? Unfortunately, the answer turns

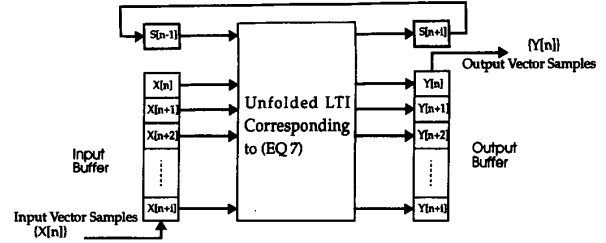


Fig. 3. Block processing.

out to be no in the general case—for any given sample period there is a limit on the best latency that can be achieved, and this limit depends on the number of primary inputs. In the following discussion, we not only find the bounds, but in the process also generate a new CDFG, i.e., a new algorithm, that achieves the bound. This new CDFG can then be used as a starting point for scheduling and resource allocation.

In the following analysis we make two reasonable assumptions. First, time is taken to be an integer, and is measured in units of adder delay. The delay for a multiplication is a multiple $m \geq 1$ of the adder delay. In effect we assume that addition takes one clock cycle, and multiplication takes one or more clock cycles. Second, the system is assumed to be arbitrary and nontrivial in the sense that no assumption is made about the values of the elements of the four coefficient matrices in the state space representation—in particular, we do not exploit coefficients that are 0 or 1. This makes the mathematics tractable, and the results conservative, at the expense of not being able to generate an analytic expression for the exact amount of hardware resources that are needed. However, this is not a problem because scheduling and allocation for a CDFG can be performed efficiently by various high-level synthesis systems.

1) *Using Unfolding with Block Processing to Arbitrarily Improve the Sample Period:* All algorithm transformations that can arbitrarily improve sample period are based on variants of unfolding where several input samples are processed together to produce one or several output samples. The computation overhead is amortized over several samples, and thus the effective sample rate is lowered. Further, the next state is now calculated in steps of the block size. Block processing [21] is one technique based on this theme where several samples are buffered, and then processed together as a block to produce the corresponding outputs. We use an adaptation of this idea for simultaneous improvement in latency and sample period, and therefore it is illustrative to look at how block processing works.

Fig. 3 shows the idea behind block processing. In the original CDFG, $X[n]$ and $S[n-1]$ are used to compute $Y[n]$ and $S[n]$. In block processing, $i+1$ consecutive input samples $X[n] \cdots X[n+i]$ are collected in a buffer, and used together with $S[n-1]$ to compute the corresponding output samples $Y[n] \cdots Y[n+i]$ and the new state $S[n+i]$. The $i+1$ output samples are put in a buffer and shifted out at the sample rate. Here $i+1$ is the blocking size. The state space formalism again provides the means to calculate the sample period and latency that are achieved by block processing. Following are

the equations for $Y[n] \cdots Y[n+i]$ and $S[n+i]$ in terms of $X[n] \cdots X[n+i]$ and $S[n-1]$ as obtained by unfolding the original state equations:

$$\begin{aligned} S[n+i] &= A^{i+1}S[n-1] + A^iBX[n] \\ &\quad + A^{i-1}BX[n+1] + \cdots + BX[n+i] \\ Y[n] &= CS[n-1] + DX[n] \\ Y[n+1] &= CAS[n-1] + CBX[n] + DX[n+1] \\ &\quad \dots \\ Y[n+i] &= CA^iS[n-1] + CA^{i-1}BX[n] \\ &\quad + CA^{i-2}BX[n-1] \cdots + CBX[n+i-1] \\ &\quad + DX[n+i]. \end{aligned} \quad (7)$$

Noting that the buffering delays at the input and the output affect the latency, and that the effective sample period is the block processing time divided by the block size, one obtains the following expressions for the latency and sample period for such a block processing algorithm:

$$\begin{aligned} T_S &= \frac{m + \lceil \log_2(R + (i+1)P) \rceil}{i+1} \\ T_L &= (2i+1)T_S = \left(\frac{2i+1}{i+1} \right) (m + \lceil \log_2(R + (i+1)P) \rceil). \end{aligned} \quad (8)$$

These expressions are obtained by organizing the computation of the various linear combinations as maximally balanced binary trees of adders, preceded by multiplications in parallel. As is obvious from these expression, $\lim_{i \rightarrow \infty} T_S = 0$, and $\lim_{i \rightarrow \infty} T_L = \infty$, so that throughput can be improved arbitrarily by increasing the amount of unfolding, but always at the cost of increased latency. This is a result of the buffering done at the input and the output.

As a side note, it is important to note that unfolding does not hurt numerical properties—in fact, as shown in [21] there is actually an improvement in round-off errors and other finite precision arithmetic effects.

2) *Using Unfolding with On-Arrival Processing and Minimum Latency Transformation to Simultaneously Optimize Latency and Sample Period:* We have found that the key to simultaneously improving sample period and latency is to combine the minimum latency transformation from Section VI-B with unfolding. However instead of using block processing, which always degrades the latency, the samples are no longer buffered, but are processed as they arrive. Intuitively this makes sense too—if latency is a concern, there is no point in idle buffering of input samples.

One strategy for bringing the unfolded system into a minimum latency form is the same as that we adopted for the system which was not unfolded—new states are introduced such that all the outputs are dependent on only one state variable. A quick look at the unfolded system equations shows that if the linear combinations CS, CAS, \dots, CA^iS of the original state vector S are added as new, though redundant, states then one can express all the outputs such that they depend on one and only one previous state value, and that too through a multiplicative coefficient of 1. The number of states increases to $R + iQ$. However, since the dimension of

state space is not more than R , it is possible to delete some of the original R states such that the remaining states continue to form a basis set for the state space. To keep the analysis simple, we avoid doing this optimization as the only effect of it will be to reduce the amount of hardware that will be needed—the critical paths will remain unchanged. Using this strategy of combining unfolding and minimum latency transformation we get the following system equations:

$$\begin{aligned} \tilde{S}[n] &= \begin{bmatrix} I_{R \times R} \\ C \\ \dots \\ CA^i \end{bmatrix} S[n] \\ &= \begin{bmatrix} A^{i+1} & 0_{R \times Q} & 0_{R \times Q} & \dots & 0_{R \times Q} \\ CA^{i+1} & 0_{Q \times Q} & 0_{Q \times Q} & \dots & 0_{Q \times Q} \\ \dots & \dots & \dots & \dots & \dots \\ CA^{2i+1} & 0_{Q \times Q} & 0_{Q \times Q} & \dots & 0_{Q \times Q} \end{bmatrix} \tilde{S}[n-1] \\ &\quad + \begin{bmatrix} A^i B \\ CA^i B \\ \dots \\ CA^{2i} B \end{bmatrix} X[n] \\ &\quad + \dots + \begin{bmatrix} B \\ CB \\ \dots \\ CA^i B \end{bmatrix} X[n+i] \end{aligned}$$

$$Y[n] = [0_{Q \times R} I_{Q \times Q} 0_{Q \times Q} \cdots 0_{Q \times Q}] \tilde{S}[n-1] + DX[n]$$

$$Y[n+1] = [0_{Q \times R} 0_{Q \times Q} I_{Q \times Q} \cdots 0_{Q \times Q}] \tilde{S}[n-1] + CBX[n] + DX[n+1]$$

$$\begin{aligned} Y[n+i] &= [0_{Q \times R} 0_{Q \times Q} 0_{Q \times Q} \cdots I_{Q \times Q}] \\ &\quad \times \tilde{S}[n-1] + CA^{i-1}BX[n] \\ &\quad + \dots + CBX[n+i-1] + DX[n+i]. \end{aligned} \quad (9)$$

Since we are interested in finding the limits to which the sample period and latency can be improved simultaneously, our task is essentially one of scheduling the above computation such that we get minimum latency for a given sample period. Such a schedule, while requiring much hardware, will be the fastest. If T_S be the sample period (an integer ≥ 1), then the arrival time of input samples $X[n] \cdots X[n+i]$ are $nT_S \cdots (n+i)T_S$. This computation is not a one shot computation in the sense that computation for one set of $i+1$ samples is followed by the computation of the next set of $i+1$ samples, so that one can overlap these computations under the constraint that the state value needed by one set of computations is produced in the preceding set. To put it another way, the time at which $\tilde{S}[n-1]$ is available is also a parameter. Let the arrival of $\tilde{S}[n-1]$ be skewed by T_j with respect to the arrival of $X[n]$, or equivalently, let $\tilde{S}[n-1]$ be available at $nT_S + T_j$. Since a similar computation needs to be done for the next block of $i+1$ samples, it follows that $\tilde{S}[n+i]$ has to be available by $(n+i+1)T_S + T_j$. This is pictorially depicted in Fig. 4. Note that T_j may even be negative, although since $\tilde{S}[n-1]$ depends on $X[n-1]$, T_j obviously cannot be less than $-T_S$; in fact as later results in the paper show, this is not a sufficient condition for the implementation to be feasible.

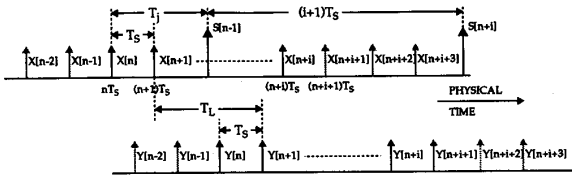


Fig. 4. Unfolding + on-arrival processing.

From the preceding discussion it is clear that for a given LTI system with sample period T_S , these two parameters—the amount of unfolding i , and the skew T_j in the computation of state vector—together span the range of possible realizations obtained as a result of unfolding. However, the constraint that $\tilde{S}[n+i]$ has to be computed by $(n+i+1)T_S + T_j$ may not be met for all values of these two parameters, even if there is no constraint on the amount of hardware resources. Further, it is intuitively obvious that if T_j is too large, then latency will get hurt because the computation of output values will wait for the state values to be available. We analyze the effect of these two parameters and the sample period on the feasibility of an LTI system (i.e., existence of a schedule), and in case the system is feasible, we calculate the best latency that can be obtained. These results enable one to answer whether a given pair of constraints on T_S and T_L can be met, and also construct a transformed algorithm or CDFG that will meet the constraints.

Feasibility Theorem: It is feasible to implement a LTI system with nontrivial coefficient matrices, which has been unfolded i times, has a sample period T_S , and a skew in the arrival of previous state value of T_j , if and only if

$$P \leq \left(\frac{2^{T_j}(2^{T_S} - 1)}{2^m} \right) \left(\frac{2^{(i+1)T_S} - 2^m R}{2^{(i+1)T_S} - 1} \right) \quad (10)$$

where P is the number of inputs, R is the number of states in the initial CDFG before unfolding and the minimum latency transformation. The adder delay is 1 and the multiplier delay is $m \geq 1$.

Proof: The problem of feasibility occurs because one has to ensure that $\tilde{S}[n+i]$ can be computed by time $(n+i+1)T_S + T_j$ from $X[n] \cdots X[n+i]$ that arrive at time instances $nT_S \cdots (n+i)T_S$, and $\tilde{S}[n-1]$ that arrives at time $nT_S + T_j$. To show that $\tilde{S}[n+i]$ can be computed by this deadline, it is sufficient to show that a maximally fast schedule, where one starts making use of the data values as soon as they arrive, will succeed.

The update equation for each of the $R + iQ$ elements in the state vector \tilde{S} is a linear combination of the R original elements of the state vector, and P input elements for each

of the $i+1$ samples. Each of these linear expressions can be evaluated independently in parallel so that it suffices to consider the computation time of just one of them. The fastest way to evaluate such linear combinations is to multiply each of the variables with the corresponding coefficient as soon as the variable arrives, and simultaneously do maximally parallel pair-wise addition of all terms whose variable-coefficient multiplication is finished. This process will continue until all the input vectors $X[n] \cdots X[n+i]$ and the previous state vector $\tilde{S}[n-1]$ have arrived, multiplied with the corresponding coefficient vectors, and the resulting terms all added. At the beginning there are no terms available to add. At every sample period, $kT_S \forall k \in \{n \cdots n+i\}$, an input sample vector with P elements is received, and m time units later the P elements would have been multiplied with coefficients and ready for addition. Also, at time $nT_S + T_j$ an additional R elements will be available (only R out of $R+iQ$ elements of $\tilde{S}[n-1]$ contribute to the state update equations), and become ready for addition m time units later. Let $\alpha(T, T_S, T_j, P, R, k_1, k_2, m)$ denote the number of terms available to be added at time T in the maximally fast schedule described above for evaluating the linear expression corresponding to the case where P elements arrive at each of the time instances kT_S , $\forall k \in \{k_1 \cdots k_2\}$, $k_1 \leq k_2$, and R elements arrive at time T_j . Then

$$\begin{aligned} \text{Feasibility that } \tilde{S}[n+i] \text{ be computed by } (n+i+1)T_S + T_j \\ \Leftrightarrow \alpha((i+1)T_S + T_j, T_S, T_j, P, R, 0, i, m) = 1. \end{aligned}$$

It can be shown that (11) is true (see the bottom of this page). Using (11), the feasibility requirement $\alpha((i+1)T_S + T_j, T_S, T_j, P, R, 0, i, m) = 1$ can be reduced to the equation at the bottom of the next page.

Note that the logical relations $(i+1)T_S \geq m$ and $(i+1)T_S + T_j \geq m$ are always true because one cannot avoid multiplying with coefficients, and that $\min(i, \lfloor ((i+1)T_S + T_j - m)/T_S \rfloor) = i$ because $\tilde{S}[n+i]$ cannot be computed until at least a multiplier delay m after the arrival of $X[n+i]$ at time iT_S . The above equality can then be simplified to the desired form of the feasibility condition in the theorem statement.

Observations:

- If no schedule is feasible for the given set of parameters, one can adopt any of the following remedies: increase the skew T_j in the arrival of previous state value, increase the sample period T_S , make the multiplier faster i.e., reduce m , or increase the amount of unfolding i . The first two strategies will always work although, as shown later, increasing T_j can result in increased latency, and increasing T_S may not be desirable. The latter two strategies may not always work. Reducing m may not be

$$\alpha(T, T_S, T_j, P, R, k_1, k_2, m) = \left\lfloor \frac{(T - m \geq T_j)(2^{T_j} R) + (T - m \geq k_1 T_S) \left(\frac{2^{(1 + \min(k_2, \lfloor \frac{T-m}{T_S} \rfloor)) T_S - 2^{k_1 T_S} P}}{2^{T_S} - 1} \right)}{2^{T-m}} \right\rfloor. \quad (11)$$

possible if it is already equal to 1. Increasing the amount of unfolding may be futile because although the right-hand side of the feasibility condition increases with i , it converges to a finite value. In particular, the second term on the right-hand side converges to 1 as $i \rightarrow \infty$. Also, it is worth noting that unfolding is an expensive operation in terms of hardware resource requirements.

- b) The LTI system is feasible for some finite unfolding factor i if and only if:

$$P < \left(\frac{2^{T_j}(2^{T_s} - 1)}{2^m} \right) \text{ or, equivalently} \\ T_j > m + \left\lceil \log_2 \left(\frac{P}{2^{T_s} - 1} \right) \right\rceil. \quad (12)$$

The equivalence between the statement ‘‘LTI system is feasible for some finite unfolding factor i ’’ and the inequality $P < (2^{T_j}(2^{T_s} - 1)/2^m)$ follows from the fact that the second factor on the right hand side of (10) is less than 1 for any finite i , takes its minimum value at $i = 0$, and monotonically converges to 1 as $i \rightarrow \infty$. The equivalence between the two inequalities in (12) follows from simple algebraic manipulation, and the facts that T_j is an integer, and that $x < n \Leftrightarrow \lceil x \rceil < n$ for all real x and integer n .

- c) If (12) holds, then for given T_s and T_j the unfolding factor i must satisfy the following for the system to be feasible:

$$i \geq \left\lceil \frac{1}{T_s} \left(\log_2 \left(\frac{2^m R 2^{T_j} (2^{T_s} - 1) - 2^m P}{2^{T_j} (2^{T_s} - 1) - 2^m P} \right) \right) \right\rceil - 1. \quad (13)$$

Equation (13) follows directly from (10) in the statement of the Feasibility Theorem by solving (20) for i , and making use of the fact i is an integer.

- d) For given sample period T_s , and unfolding factor i , the skew in the state arrival T_j must satisfy:

$$T_j \geq \left\lceil \log_2 \left(\frac{2^m P (2^{(i+1)T_s} - 1)}{(2^{(i+1)T_s} - 2^m R)(2^{T_s} - 1)} \right) \right\rceil. \quad (14)$$

Equation (14) follows from (10) by recasting the inequality in (10) in terms of T_j , and making use of the fact that T_j is an integer.

Latency Theorem: If an LTI system with nontrivial coefficient matrices, which has been unfolded i times, has a sample period T_s , and a skew in the arrival of previous state value of

T_j , is feasible, then it can achieve a latency of:

$$T_L = m + \lceil \log_2(2^{T_j - m} + P) \rceil = T_j + \left\lceil \log_2 \left(1 + \frac{P}{2^{T_j - m}} \right) \right\rceil. \quad (15)$$

Note that the latency T_L is independent of unfolding i .

Proof: Latency varies from sample to sample, i.e., the time taken for output samples $Y[n] \cdots Y[n+i]$ to be calculated since the arrival of the corresponding input sample $X[n] \cdots X[n+i]$ is in general different for each of the $i+1$ samples. Let $T_L[k]$ be the latency for the $Y[n+k]$ sample $\forall k = 0 \cdots i$. We define the system latency T_L to be $\max(T_L[0], T_L[1], \dots, T_L[i])$.

From (9) it is clear that $Y[n+k]$ is a linear combination of the $R + iQ$ elements of the previous state vector $\tilde{S}[n-1]$ that arrive at time $nT + T_j$, and P elements in each of the vectors $X[n] \cdots X[n+k]$ that arrive at time instants $nT_s \cdots (n+k)T_s$. The fastest way to evaluate such a linear combination, using the same strategy as we adopted in the proof of the feasibility theorem, is to multiply each of the variable elements by its corresponding coefficient as soon as it arrives, and simultaneously do maximally parallel pair-wise additions of all terms whose variable-coefficient multiplication is finished. This works because addition is associative and commutative. The process continues until all elements have arrived, multiplied, and been added so that a single term remains. One thing to note from the state-space output equations obtained after unfolding and minimum latency transformation is that of the $R + iQ$ elements of $\tilde{S}[n-1]$, one element has a coefficient of 1, and the others have coefficients 0. So in effect $\tilde{S}[n-1]$ contributes only one variable to each output, and that too with a coefficient of 1 so that no multiplication is needed.

Using such a maximally fast schedule, one can express $T_L[k]$ as below using the function $\alpha(T, T_s, T_j, P, R, k_1, k_2, m)$ that we defined earlier:

$$T_L[k] = \max(kT_s + m, T_j) + \lceil \log_2(\alpha(\max(kT_s + m, T_j), T_s, T_j - m, P, 1, 0, k, m)) \rceil - kT_s$$

and using the expression for $\alpha(T, T_s, T_j, P, R, k_1, k_2, m)$ from (11), one can simplify the above expression to:

$$T_L[k] = m + \left\lceil \log_2 \left(2^{T_j - kT_s - m} + P \left(\frac{2^{(k+1)T_s} - 1}{2^{(k+1)T_s} - 2^{kT_s}} \right) \right) \right\rceil. \quad (16)$$

From (16) together with the feasibility theorem one can show that

$$T_L[k] \geq T_L[k+1] \quad \forall k \in \{0 \cdots i-1\}$$

$$\left\lceil \frac{((i+1)T_s \geq m)(2^{T_j} R) + ((i+1)T_s + T_j \geq m) \left(\frac{2^{\left(1 + \min(i, \left\lfloor \frac{((i+1)T_s + T_j) - m}{T_s} \right\rfloor)\right) T_s} - 1}{2^{T_s} - 1} P \right)}{2^{((i+1)T_s + T_j) - m}} \right\rceil = 1.$$

which means that the latency associated with $Y[n]$ is at least as high as latencies for $Y[n+1] \cdots Y[n+i]$. From this one immediately gets the following which proves the theorem:

$$T_L = \max(T_L[0], T_L[1], \dots, T_L[i]) = T_L[0] \\ = m + \lceil \log_2(2^{T_j-m} + P) \rceil = T_j + \left\lceil \log_2 \left(1 + \frac{P}{2^{T_j-m}} \right) \right\rceil.$$

Observations:

- a) Intuitively, from (15) it follows that one should have small values of T_j in order to have low latency. However, recall from (14) that a certain minimum value of T_j is needed in order for the system to be feasible. This suggests that one cannot arbitrarily improve the sample period T_S in presence of a constraint on latency T_L .
- b) A necessary and sufficient condition for a feasible system to achieve latency $\leq T_L$ is :

$$T_j \leq m + \lceil \log_2(2^{T_L-m} - P) \rceil \quad \text{or, equivalently} \\ 2^{T_j} \leq 2^{T_L} - 2^m P. \quad (17)$$

The observation follows from (15).

- c) A necessary and sufficient condition for a feasible system to achieve the minimum latency, i.e., $T_L = m + \lceil \log_2(1 + P) \rceil$, is:

$$T_j \leq m + \lceil \log_2(2^{\lceil \log_2(1+P) \rceil} - P) \rceil. \quad (18)$$

Equation (18) is obtained by plugging $T_L = m + \lceil \log_2(1 + P) \rceil$ into (17).

- d) For a given sample period T_S , the best latency that can be achieved by a feasible system is:

$$T_L = m + \left\lceil \log_2 \left(2^{1 + \lceil \log_2 \left(\frac{P}{2^{T_S-1}} \right) \rceil} + P \right) \right\rceil. \quad (19)$$

From (15) it follows that the best latency is obtained for the smallest T_j such that the system is still feasible for the specified T_S at some finite unfolding factor i . Making use of the fact that T_j is an integer, from (12) it is obvious that this smallest T_j is given by $T_j = 1 + m + \lceil \log_2(P/(2^{T_S} - 1)) \rceil$. Equation (19) then follows trivially by plugging this value of T_j back into (15).

- e) The minimum sample period for a feasible system that achieves the minimum latency, i.e., $T_L = m + \lceil \log_2(1 + P) \rceil$, is:

$$T_S = 1 + \left\lceil \log_2 \left(\frac{2^{\lceil \log_2(1+P) \rceil}}{2^{\lceil \log_2(1+P) \rceil} - P} \right) \right\rceil. \quad (20)$$

To derive (20) we first note that the feasibility condition for LTI systems at some finite unfolding factor, as implied by (12), is equivalent to $T_S > \log_2(1 + 2^{m-T_j}P)$. Since T_S is an integer, it follows that the minimum sample period at a given T_j for a system that is feasible for some finite unfolding factor is $T_S = 1 + \lceil \log_2(1 + 2^{m-T_j}P) \rceil$. Note that this minimum sample period will be smallest for the largest T_j . According to (18) the largest T_j for a feasible system to achieve $T_L = m + \lceil \log_2(1 + P) \rceil$ is $T_j = m + \lceil \log_2(2^{\lceil \log_2(1+P) \rceil} - P) \rceil$.

Equation (20) is obtained by plugging this expression for T_j into $T_S = 1 + \lceil \log_2(1 + 2^{m-T_j}P) \rceil$ and simplifying.

- f) From (19) and (20) it follows that for any single input LTI system, such as a large fraction of DSP systems, our method can achieve a latency of $T_L = m + 1$ for all $T_S \geq 2$, and a latency of $T_L = m + 2$ at $T_S = 1$. These numbers are far superior to the fastest reported results in literature even in the specific case of the popular fifth-order elliptical filter benchmark ($T_L = m + 3$ and $T_S = m + 3$ reported in [16] for $m = 2$).

Similarly, for any 2-input LTI system our method guarantees that a latency of $T_L = m + 2$ can be achieved for all $T_S \geq 2$, and that a latency of $T_L = m + 3$ can be achieved at $T_S = 1$.

- g) We would like to emphasize that the results in observations d and e, are independent of R , the number of states in the CDFG, and therefore the latency and sample period implied by (19) and (20) are guaranteed for all LTI CDFG's with P inputs. Of course different amount of hardware will be required for different cases, and in the next subsection we present a heuristic algorithm to choose the unfolding factor i and the skew in state arrival T_j so as to minimize hardware.

3) *Algorithm for Simultaneous Optimization of Latency and Sample Period:* We have used the analytic results presented above to construct an algorithm which, if feasible, transforms an arbitrary LTI CDFG to satisfy user specified constraints on latency T_L , and sample period T_S .

Step 1: Transform the given LTI CDFG into state space equations characterized by P, Q, R, A, B, C , and D .

Step 2: From (4), if

$$T_S \geq \max_{r \in \{1 \cdots Q\}} (m + \lceil \log_2(R + P - nr0(A, r) - nr0(B, r) - (nr1(A, r) + nr1(B, r))(1 - 1/2^m)) \rceil),$$

and,

$$T_L \geq \max_{r \in \{1 \cdots R\}} (m + \lceil \log_2(R + P - nr0(C, r) - nr0(D, r) - (nr1(C, r) + nr1(D, r))(1 - 1/2^m)) \rceil)$$

then just the maximally fast computation of the linear expressions will suffice—STOP.

Step 3: If, in accordance with (6), $T_S \geq m + \lceil \log_2(R + P) \rceil$ and, $T_L \geq m + \lceil \log_2(1 + P) \rceil$, then apply the minimum latency transformation from (5) to the state space equations—STOP.

Step 4: We need to apply techniques from Section VI-C2. If $T_L < m + \lceil \log_2(1 + P) \rceil$, then our method can never find a feasible system because $m + \lceil \log_2(1 + P) \rceil$ is the minimum latency that can be achieved for a P input LTI system with arbitrary coefficients—STOP.

Step 5: Use (17) to find the upper bound $T_{j,UB}$ on the parameter T_j so that latency is $\leq T_L$.

$$T_{j,UB} = m + \lceil \log_2(2^{T_L-m} - P) \rceil. \quad (21)$$

Step 6: Use (12) to find lower bound $T_{j,LB}$ on T_j such that the system is feasible for finite i .

$$T_{j,LB} = m + 1 + \left\lceil \log_2 \left(\frac{P}{2^{T_S-1}} \right) \right\rceil. \quad (22)$$

Step 7: If $T_{j,LB} > T_{j,UB}$, we cannot find a feasible system which satisfies the constraints on sample period and latency—STOP.

Step 8: Since unfolding is an expensive operation, particularly from the perspective of the size of coefficient memory, we would like to keep i as small as possible. Therefore we initially pick the largest allowable T_j , i.e., $T_j = T_{j,UB}$, and use (13) to calculate the minimum unfolding factor i for which a feasible system satisfies the constraints on latency and sample period.

$$i = \left\lceil \frac{1}{T_S} \left(\log_2 \left(\frac{2^m R 2^{T_{j,UB}} (2^{T_S} - 1) - 2^m P}{2^{T_{j,UB}} (2^{T_S} - 1) - 2^m P} \right) \right) \right\rceil - 1. \quad (23)$$

Step 9: Use (14) to find a new lower bound on T_j such that the system is feasible for the above i .

$$T_{j,LB1} = \left\lceil \log_2 \left(\frac{2^m P (2^{(i+1)T_S} - 1)}{(2^{(i+1)T_S} - 2^m R)(2^{T_S} - 1)} \right) \right\rceil. \quad (24)$$

Step 10: Using i from (23), and a $T_j \in [T_{j,LB1}, T_{j,UB}]$ from (21) and (24), generate the CDFG corresponding to the maximally fast schedule as described in Section VI-C2. This CDFG corresponds to a new algorithm that satisfies the constraints on latency and sample period.

D. Latency-Throughput Optimization Technique in Action—An Example

While we will present the results achieved by our techniques on various benchmarks later in the paper, here we illustrate the technique described in the previous section by discussing in detail the application of the technique to a specific example—a fifth-order low-pass elliptical wave digital IIR filter. Fig. 5 shows the initial CDFG for the example, which has $P = 1$, $Q = 1$, and $R = 5$. Using the definitions in Section III-A it can be shown that if hardware was no constraint, this CDFG (without applying any transformations) can achieve the following combinations of latency and sample period: $T_L = 2m + 5$ and $T_S = 3m + 6$.

As mentioned previously, for $P = 1$ the technique for this section guarantees that one can achieve $T_L = m + 1$ for all $T_S \geq 2$, and a latency of $T_L = m + 2$ at $T_S = 1$ —several factors of improvement over latency and throughput achieved by the initial CDFG. We demonstrate the use of the algorithm Section VI-3C to transform the CDFG in Fig. 5 to achieve $T_L = m + 1$ and $T_S = 2$. The algorithm has been automated using the software platform described in Section VIII.

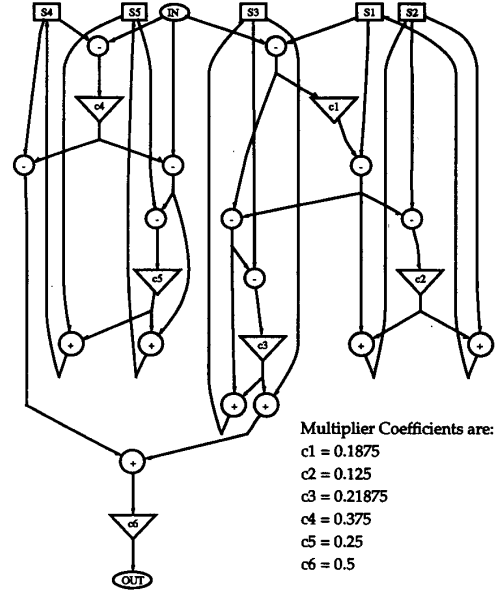


Fig. 5. Initial CDFG for a fifth-order low-pass elliptical wave digital filter with $T_l = 2m + 5$ and $T_S = 3m + 6$.

Step 1: The CDFG is converted to the state-space equations:

$$S[n] = \begin{bmatrix} \frac{13}{128} & \frac{9}{8} & 0 & 0 & 0 \\ -\frac{91}{128} & \frac{1}{8} & 0 & 0 & 0 \\ -\frac{725}{512} & 0 & \frac{7}{32} & 0 & 0 \\ 0 & 0 & 0 & \frac{3}{32} & \frac{5}{4} \\ 0 & 0 & 0 & -\frac{9}{32} & \frac{1}{4} \end{bmatrix} S[n-1]$$

$$Y[n] = \begin{bmatrix} \frac{3}{128} \\ -\frac{21}{128} \\ \frac{325}{512} \\ \frac{5}{32} \\ -\frac{15}{32} \end{bmatrix} X[n] + \begin{bmatrix} \frac{203}{1024} & 0 & \frac{39}{64} & -\frac{11}{16} & 0 \end{bmatrix} S[n-1] + \left[\frac{101}{1024} \right] X[n].$$

Step 2: Use (4) to calculate the latency and sample period that can be achieved by the above state space equations taking into account coefficients with magnitude 0 and 1. We obtain $T_L = m + 2$ and $T_S = m + 2$ —these values do not meet the requirements, therefore we continue.

Step 3: Next we apply the Minimum Latency transformation of Section VI-B and check whether the resulting system, shown below, meets the latency and sample period

requirements.

$$\tilde{S}[n] = \begin{bmatrix} \frac{13}{128} & \frac{9}{8} & 0 & 0 & 0 & 0 \\ -\frac{91}{128} & \frac{1}{8} & 0 & 0 & 0 & 0 \\ -\frac{725}{512} & 0 & \frac{7}{32} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{3}{32} & \frac{5}{4} & 0 \\ 0 & 0 & 0 & -\frac{9}{32} & \frac{1}{4} & 0 \\ -\frac{110461}{131072} & \frac{1827}{8192} & \frac{273}{2048} & -\frac{33}{512} & -\frac{55}{64} & 0 \end{bmatrix}$$

$$\times \tilde{S}[n-1] + \begin{bmatrix} \frac{3}{128} \\ -\frac{21}{128} \\ \frac{325}{512} \\ \frac{5}{32} \\ -\frac{15}{32} \end{bmatrix} X[n]$$

$$Y[n] = [0 \ 0 \ 0 \ 0 \ 0 \ 1] \tilde{S}[n-1] + \left[\frac{101}{1024}\right] X[n].$$

The above system has $T_L = m + 1$ and $T_S = m + 3$ —these values do not meet our requirements.

Step 4: Check that the required latency $T_L = m + 1$ is greater than or equal to the best latency that the algorithm guarantees $T_L = m + \lceil \log_2(1 + P) \rceil$ —since this is indeed the case, we continue.

Step 5: Calculate the upper bound $T_{j,UB}$ on T_j :

$$T_{j,UB} = m + \lceil \log_2(2^{T_L - m} - P) \rceil = m. \quad (25)$$

Step 6: Calculate the lower bound $T_{j,LB}$ on T_j :

$$T_{j,LB} = m + 1 + \left\lceil \log_2 \left(\frac{P}{2^{T_S - 1}} \right) \right\rceil = m - 1. \quad (26)$$

Step 7: Check that $T_{j,LB} \leq T_{j,UB}$ —the condition is satisfied in this example. We continue because the algorithm can find a solution.

Step 8: Calculate the unfolding factor:

$$i = \left\lceil \frac{1}{T_S} \left(\log_2 \left(\frac{2^m R 2^{T_{j,UB}} (2^{T_S} - 1) - 2^m P}{2^{T_{j,UB}} (2^{T_S} - 1) - 2^m P} \right) \right) \right\rceil - 1 \\ = \left\lceil \frac{1}{2} \left(\log_2 \left(\frac{2^m 15 - 1}{8} \right) \right) \right\rceil = 1. \quad (27)$$

where we had to pick a value for the multiplier delay m , and we assumed $m = 1$. This implies that the final system will achieve $T_L = 2$ and $T_S = 2$.

Step 9: Calculate the new lower bound $T_{j,LB1}$ on T_j :

$$T_{j,LB1} = \left\lceil \log_2 \left(\frac{2^m P (2^{(i+1)T_S} - 1)}{(2^{(i+1)T_S} - 2^m R)(2^{T_S} - 1)} \right) \right\rceil = 1. \quad (28)$$

Step 10: Unfold the system $i = 1$ times, and then apply the Minimum Latency Transformation of Section VI-B—the resulting system can achieve the required $T_L = 2$ and $T_S = 2$ (for $m = 1$). Schedule the resulting state space equations, that are shown below, using maximally fast computation of the linear expressions and on-arrival processing with $T_j \in [T_{j,LB1}, T_{j,UB}] = [1, 1] = 1$, and noting that $X[n]$ arrives at $2n$, $X[n+1]$ arrives at $2n+2$, $\tilde{S}[n-1]$ arrives at $2n+1$, $Y[n]$ is needed by $2n+2$, $Y[n+1]$ is needed by $2n+4$, and $\tilde{S}[n+1]$ is needed by $2n+5$. The results of Section VI-C guarantee that such a schedule is always feasible.

It is worth noting that the improvement in latency and throughput is at the expense of the number of coefficients increasing from 6 in the original CDFG to 38 in the final equations, and the number of state elements increasing from 5 to 7. In addition, of course, more adders and multipliers will be needed. This extra hardware results in the original latency and sample period improving from $T_L = 2m + 5 = 7$ and $T_S = 3m + 6 = 9$, respectively, to $T_L = 2$ and $T_S = 2$ (see equation at the bottom of this page).

$$\tilde{S}[n] = \begin{bmatrix} -\frac{12935}{16384} & \frac{261}{1024} & 0 & 0 & 0 & 0 & 0 \\ -\frac{2639}{16384} & -\frac{803}{1024} & 0 & 0 & 0 & 0 & 0 \\ -\frac{29725}{65536} & -\frac{6525}{4096} & \frac{49}{1024} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -\frac{351}{1024} & \frac{55}{128} & 0 & 0 \\ 0 & 0 & 0 & -\frac{99}{1024} & -\frac{37}{128} & 0 & 0 \\ -\frac{7262905}{16777216} & -\frac{964917}{1048576} & \frac{1911}{65536} & \frac{3861}{16384} & -\frac{605}{2048} & 0 & 0 \\ -\frac{1221830987}{2147483648} & -\frac{80804817}{134217728} & \frac{13377}{2097152} & \frac{55143}{524288} & \frac{14465}{65536} & 0 & 0 \end{bmatrix} \tilde{S}[n-1] + \begin{bmatrix} -\frac{2985}{16384} \\ -\frac{609}{16384} \\ \frac{6925}{65536} \\ -\frac{585}{1024} \\ -\frac{165}{1024} \\ \frac{7063785}{16777216} \\ \frac{718615077}{2147483648} \end{bmatrix} X[n] + \begin{bmatrix} \frac{3}{128} \\ -\frac{21}{128} \\ \frac{325}{512} \\ \frac{5}{32} \\ -\frac{15}{32} \\ \frac{37229}{131072} \\ \frac{7063785}{16777216} \end{bmatrix} X[n+1]$$

$$Y[n] = [0000010]S[n-1] + \left[\frac{101}{1024}\right]X[n] \\ Y[n+1] = [0000001]S[n-1] + \left[\frac{37229}{131072}\right]X[n] + \left[\frac{101}{1024}\right]X[n+1].$$

VII. OPTIMIZING LATENCY AND SAMPLE PERIOD FOR THE SPECIAL CASE OF SINGLE-INPUT LTI SYSTEMS WITH ZERO INITIAL STATE

In the previous section we described a technique to transform a LTI CDFG to satisfy, if feasible, simultaneous constraints on latency and throughput. While on the one hand the technique is completely general in that it produces guaranteed results and is applicable to any LTI CDFG, on the other hand the technique, and the theory behind it, do not make use of the values of the coefficients in the initial state-space equation—the four coefficient matrices A , B , C , and D in (2) were assumed to be absolutely arbitrary. Better results are possible if one were to take advantage of the coefficient values. In particular, coefficients that are 0 need not be considered, and coefficients that have magnitude 1 need not be multiplied.

Taking advantage of coefficients with values 0 and 1 is, unfortunately, extremely difficult—the mathematical analysis becomes intractable—unless there is some regularity and mathematical structure to the location of these coefficients in the matrices A , B , C , and D . Of course there is no such regularity or structure in the general case of a LTI CDFG. However, it is well known in DSP and Linear System Theory that a special case of these LTI CDFG's, namely LTI CDFG's with single-input and zero initial state can always be transformed to certain standard (canonical) CDFG structures [22]. Since throughput and implementation cost have been the more popular metrics in traditional DSP, these standard CDFG structures have been developed and analyzed with those two metrics in mind—latency has been overlooked. We noticed that some of these standard forms either have good latency and throughput characteristics at low cost, or are good starting points to apply some of the techniques of the previous section.

In this section we describe some techniques that are based on the approach of first converting a single-input single-output LTI CDFG with zero initial state (a single-input multiple-output CDFG can be treated as a collection of multiple single-input single-output CDFG's) into one of the standard forms, and then applying a fixed sequence of transformations to yield cost efficient solutions with good (often optimal) throughput and latency characteristics. Since a large fraction of applications (e.g. many filters and controllers) are single-input, these techniques can indeed be useful in a variety of designs.

Before describing them, we would like to describe some differences between the special case techniques of this section and the general technique of the previous section. The first crucial difference is that some of the standard forms used by the techniques in this section often do not have good numerical accuracy which may lead to a larger number of bits being required. In contrast, the technique in the previous section by and large preserves the numerical accuracy of the original CDFG. A second important difference is that these special case techniques yield point solutions—a single combination of T_S and T_L —whereas, the technique of the previous section yielded a whole range of optimal combinations of T_S and T_L .

Technique #1: Modified Direct Form II We begin by describing a technique that is based on the well known standard form: *Direct Form II*. Any single-input single-output LTI

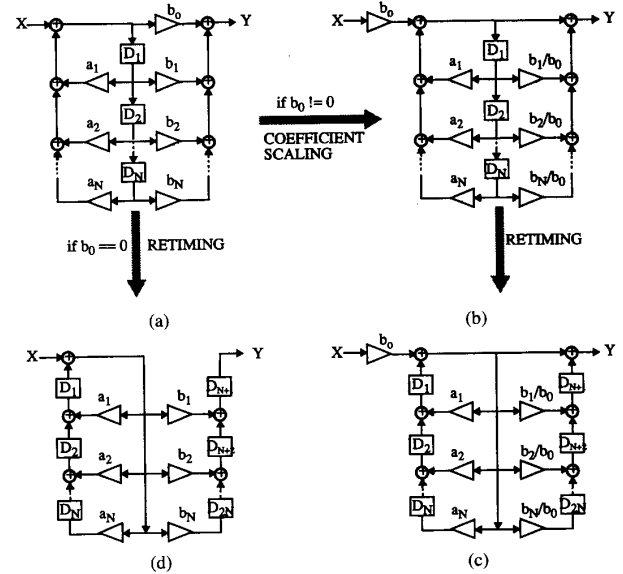


Fig. 6. Transforming Direct Form II to a low latency and high throughput structure.

CDFG with zero initial state can be transformed to the Direct Form II (shown in Fig. 6(a)) by first taking the Z -Transform of the state space equations to calculate the *Transfer Function* $H(z) = Y(z)/X(z)$ as a ratio of two polynomials in z

$$\begin{aligned} H(z) &= \frac{Y(z)}{X(z)} = C(zI - A)^{-1}B + D \\ &= \frac{C_{\text{adj}}(zI - A)B}{\det(zI - A)} + D \\ &= \left(\sum_{i=0}^N \alpha_i z^i \right) / \left(\sum_{i=0}^N \beta_i z^i \right) \quad \beta_N \neq 0, \quad N \leq R \end{aligned} \quad (29)$$

and then taking the Inverse Z -Transform to obtain the time-domain equation:

$$Y[n] = \sum_{i=0}^N \left(\frac{\alpha_{N-i}}{\beta_N} \right) X[n-i] + \sum_{i=1}^N \left(-\frac{\beta_{N-i}}{\beta_N} \right) Y[n-i]. \quad (30)$$

Comparing (30) to the CDFG for the Direct Form II structure in Fig. 6(a) it follows that the coefficients in the Direct Form II structure can be expressed as below:

$$a_i = \frac{\beta_{N-i}}{\beta_N} \quad b_j = \frac{\alpha_{N-i}}{\beta_N} \quad i \in 1 \dots N, \quad j \in 0 \dots N. \quad (31)$$

The Direct Form II structure is widely recognized as a low-throughput, high-latency, high-cost structure, and does not appear to be particularly useful—in the general case of arbitrary coefficients this structure has latency $T_L = 2m + N + 1$ and sample period $T_S = m + N$. However, we found that a simple and fixed sequence of transformation steps can yield a modified structure that has latency $T_L = m + 1$ and sample period $T_S = m + 2$. By way of comparison, the general technique of Section VI-C can achieve $T_L = m + 1$ for all

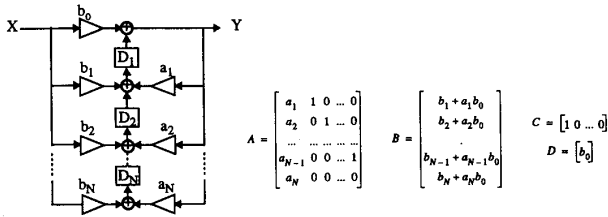


Fig. 7. The Transposed Direct Form II (also called Companion Form) with $T_L = m + 1$ and $T_S = m + 2$.

$T_S \geq 2$, and $T_L = m + 2$ at $T_S = 1$ for the single-input case. The algorithm to modify the Direct Form II structure is as below.

Step 1: If the coefficient b_0 in Fig. 6(a) is not equal to 0, we first apply a coefficient scaling transformation to obtain the CDFG in Fig. 6(b), and then apply a sequence of retiming steps that move the delay nodes in the middle branch of the CDFG to two sides as in Fig. 6(c).

If the coefficient b_0 in Fig. 6(a) is equal to 0, we remove the corresponding multiplication node and directly apply a sequence of retiming steps to move the delay nodes in the middle branch of the CDFG to the two sides as shown in Fig. 6(d).

Step 2: We convert the CDFG obtained in Step 1 into the state space representation.

Step 3: Apply maximally fast linear computation transformation [16]—using (4) one can show that $T_L = m + 1$ and $T_S = m + 2$ for $b_0 \neq 0$, and $T_L = 0$ and $T_S = m + 2$ for $b_0 = 0$.

Technique #2: Modified Direct Form II with One Level of Unfolding + On-Arrival Processing: This is an extension to technique #1 where we unfold once the system obtained in technique #1, and then use maximally fast computation, on-arrival processing, and state arrival skew $T_j = 0$ to get $T_L = m + 1$ and $T_S = m + 1$. Compared to technique #1, this technique reduces the sample period by 1, achieves the same latency, and has $(8N + 4)$ coefficients as opposed to $(4N + 1)$ for approximately $\times 2$ increase in coefficient memory.

Technique #3: Transposed Direct Form II: This technique is based on the observation that another standard form known as *Transposed Direct Form II* (also known as the *Companion Form*) has good latency throughput characteristics. This form is shown in Fig. 7, and has sample period $T_S = m + 2$ and latency $T_L = m + 1$. The coefficient matrices in the state space representation are also shown in the figure, and are equal to the corresponding coefficients in Fig. 6(a) for *Direct Form II*. The advantage of this structure over that obtained in technique #1 is that there are only N state nodes, as opposed to $2N$. Same throughput and latency is obtained as in technique #1 with $\times 2$ savings in number of registers used to store state variables, and only $2N + 1$ coefficients being required as opposed to $4N + 1$ resulting in an almost $\times 2$ savings in coefficient memory as well.

Technique #4: Transposed Direct Form II with One Level of Unfolding + On-Arrival Processing: Similar to the approach used in technique #2, a latency of $T_L = m + 1$ and a sample period of $T_S = m + 1$ are obtained if the *Transposed Direct*

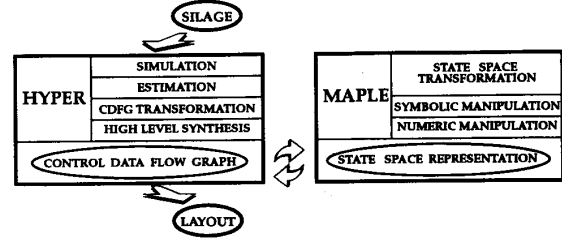


Fig. 8. Software platform based on coupling of HYPER 2.1 and Maple V.

Form II of Fig. 7 is unfolded by 1, and then implemented using on-arrival processing and maximally fast linear computation. However, this technique has only N state nodes and $(4N + 2)$ coefficients, for a $\times 2$ advantage in both the number of states and the size of coefficient memory over technique #2 for the same latency and sample period.

VIII. SOFTWARE PLATFORM—HYPER 2.1 + MAPLE V

The approach presented here imposes a number of demanding requirements of very distinct and different nature on the synthesis process. The requirements include simulation which addresses word-length trade-offs after transformations, manipulation of computations in both the CDFG and the state-space domains using both numeric and symbolic techniques, rapid estimation of quantitative performances of proposed solutions, and a software platform which supports a large number of transformations.

We created a software platform satisfying these requirements by interfacing Maple V, a computer algebra system originally from University of Waterloo [3], with HYPER 2.1 [9], a high level synthesis system from University of California, Berkeley. Maple allows symbolic as well as numerical manipulations in a mixed functional and procedural paradigm. Various packages in Maple enable the representation and manipulation of a variety of mathematical structures such as those found in linear algebra, group theory etc. HYPER on the other hand provides the capability to input a design described in the applicative language SILAGE, translate it into a CDFG, and perform high-level synthesis tasks such as module selection, various transformations, scheduling, allocation and hardware mapping.

Fig. 8 shows a typical flow of synthesis for joint latency/throughput optimization using the techniques described in this paper through HYPER and Maple. The user describes the algorithm in SILAGE, and HYPER translates it into the CDFG format and does the initial simulation. An interface program then translates the CDFG into Maple equations, and uses a Maple script to perform various transformations (such as unfolding, minimum latency transformation, conversion to standard forms) in the state space and to derive a structure which satisfies the constraints on both latency and throughput. The result is then fed back to HYPER, which then estimates the implementation cost and performs bit level simulation. HYPER can also be used for further optimizations using transformations.

TABLE I
CHARACTERISTICS OF THE BENCHMARK EXAMPLES

Design Name	Description	Characteristics of the Initial CDFG				
		P	Q	R	T_L [cycles]	T_S [cycles]
mat	3-state 1-input linear controller	1	1	3	4	4
ellip	4-state 1-input linear controller	1	1	4	5	5
lin4	5-state 1-input linear controller	1	1	5	6	6
lin5	5-state 1-input linear controller	1	1	5	6	6
iir5 (wdf5)	5th order low pass elliptic wave digital IIR filter	1	1	5	7	9
iir6	6th order lowpass elliptic cascade IIR filter	1	1	6	8	7
iir8	8th order bandpass direct form IIR filter	1	1	8	11	10
iir10	10th order bandstop Butterworth cascade IIR filter	1	1	10	17	15
iir11	11th order low pass Chebyshev cascade IIR filter	1	1	11	19	19
iir12	12th order bandpass Chebyshev cascade IIR filter	1	1	12	20	18
steam	power plant controller	1	1	5	6	6
dist	distillation plant linear controller	2	1	5	5	6
chemical	chemical plant controller	2	1	4	4	6
aircraft	6 knot linear plant controller for VSTOL aircraft	3	3	14	10	10

TABLE II
IMPROVEMENTS IN LATENCY AND THROUGHPUT OF THE INITIAL DESIGN USING THE HEURISTIC TECHNIQUE OF SECTION VII, AND THE OPTIMUM TECHNIQUES OF SECTION VI-C. T_L AND T_S ARE EXPRESSED AS THE NUMBER OF CYCLES

Design Name	Initial Design			Optimum Technique ^a						Heuristic Techniques ^b										
	T_L	T_S	Area [mm ²]	T_L	T_S	Area [mm ²]	T_L	T_S	Area [mm ²]	T_L	T_S	Area [mm ²]								
mat	4	4	14.1	2	3	2	1	15.5	35.3	2	2	2	3	2	3	2	13.9	15.6	11.3	10.5
ellip	5	5	52.6	2	3	2	1	90.0	159.8	2	2	2	3	2	3	2	36.3	21.7	65.4	36.8
lin4	6	6	67.6	2	3	2	1	183.9	243.0	2	2	2	3	2	3	2	30.0	73.1	27.7	44.1
lin5	6	6	223.3	2	3	2	1	614.5	805.3	2	2	2	3	2	3	2	162.8	234.4	91.6	141.1
iir5 (wdf5)	7	9	28.9	2	3	2	1	138.2	215.9	2	2	2	3	2	3	2	73.3	115.5	40.3	63.1
iir6	8	7	10.9	2	3	2	1	71.8	218.8	2	2	2	3	2	3	2	32.3	47.3	22.1	33.2
iir8	11	10	29.3	2	3	2	1	253.8	424.0	2	2	2	3	2	3	2	36.7	43.0	24.5	28.9
iir10	17	15	22.7	2	3	2	1	259.8	566.6	2	2	2	3	2	3	2	78.2	117.6	49.7	81.1
iir11	19	19	20.6	2	3	2	1	166.9	466.4	2	2	2	3	2	3	2	70.2	106.6	45.6	71.2
iir12	20	18	25.7	2	3	2	1	317.0	886.1	2	2	2	3	2	3	2	104.1	175.5	70.0	106.3
steam	6	6	82.3	2	3	2	1	184.0	370.3	2	2	2	3	2	3	2	36.0	88.1	34.0	54.3
dist	5	6	36.2	3	4	2	1	111.6	225.6	Not App.	Not App.	Not Applicable								
chemical	4	6	30.1	3	4	2	1	101.6	143.4	Not App.	Not App.	Not Applicable								
aircraft	10	10	25.6	3	4	3	1	61.3	155.4	Not App.	Not App.	Not Applicable								

a. For each example the first set of numbers corresponds to the case when the optimum technique of Section 6.3 is used to achieve a minimum latency system; the second set of number corresponds to the case when it is used to achieve a maximum throughput system.

b. The four sets of numbers for each example correspond to the special-case techniques #1, #2, #3, and #4 respectively of Section 7.0

IX. EXPERIMENTAL RESULTS

Using the HYPER + Maple based software platform, we tested the effectiveness of the optimal transformation technique and the special case heuristic techniques (for single-input systems) on a number of examples. As mentioned in Section I-C, the latency-throughput transformation techniques can be used for two distinct purposes: to transform a CDFG to meet joint constraints on latency and throughput, and to transform a CDFG so as to improve the cost of implementation at the same latency and throughput as that of the initial CDFG. We tested our techniques in both these modes. Since we are not aware of any previous work on algorithm transformations that simultaneously addresses latency and throughput, we are unable to compare our results when using our transformations

in the first mode. When using the transformations in the second mode, we compare the cost of implementing the initial design and the final design.

The characteristics of the examples that we tested our techniques on are shown in Table I. All the results in this section were obtained using $m = 1$. The area numbers were obtained using the HYPER [9] high-level synthesis system and LAGER [23] silicon compiler. We used a 1.2 micron feature size technology.

Table II presents the results obtained when the transformations described in this paper were used in the first mode mentioned above—to simultaneously reduce latency and sample period. The table contains the latency, sample period, and chip area for the initial CDFG, and the corresponding numbers

TABLE III
IMPROVEMENTS IN AREA OVER THE INITIAL DESIGN USING THE HEURISTIC TECHNIQUE OF SECTION VII, AND THE OPTIMUM TECHNIQUES OF SECTION VI-C WHEN THE TRANSFORMED DESIGN IS SCHEDULED FOR THE SAME LATENCY AND THROUGHPUT AS THE ORIGINAL DESIGN (USING $m = 1$)

Design Name	T_L [cycles]	T_S [cycles]	Initial Design Area [mm ²]	Optimum Technique ^a Area [mm ²]		Heuristic Techniques ^b Area [mm ²]			
				O1	O2	H1	H2	H3	H4
mat	4	4	14.07	9.29	11.15	9.62	9.84	5.31	6.47
ellip	5	5	52.58	38.69	37.40	23.79	34.07	15.20	18.37
lin4	6	6	67.59	41.14	42.39	31.55	33.24	16.03	19.53
lin5	6	6	223.31	126.34	118.39	77.35	108.3	28.70	56.68
iir5 (wdf5)	7	9	28.92	28.81	44.62	27.77	34.63	18.16	34.63
iir6	8	7	10.90	21.53	26.29	16.66	20.41	10.90	13.90
iir8	11	10	29.25	21.53	26.29	14.18	20.22	12.95	15.34
iir10	17	15	22.69	39.10	47.75	23.19	32.22	19.36	26.06
iir11	19	19	20.58	43.86	44.44	20.65	28.81	15.12	22.25
iir12	20	18	25.73	57.64	59.41	25.04	46.83	22.39	29.99
steam	6	6	82.34	41.14	46.79	23.79	25.29	14.87	16.37
dist	5	6	36.18	35.22	48.82	-	-	-	-
chemical	4	6	30.09	25.57	28.18	-	-	-	-
aircraft	10	10	25.62	61.31	105.88	-	-	-	-

a.O1 corresponds to the case when the optimum technique of Section 6.3 is used to achieve a minimum latency system; O2 corresponds to the case when it is used to achieve a maximum throughput system.

b.H1, H2, H3, and H4 correspond to techniques #1, #2, #3, and #4 respectively of Section 7.0.

obtained by the four heuristic techniques of Section VII (only in the case of single-input examples) and the optimum technique of Section VI-C. The optimum technique can give a range of latency and sample period values—the table contains the numbers corresponding to the following two cases: at the minimum latency that can be guaranteed by the algorithm ($m + \lceil \log_2(1+p) \rceil$) for a P -input system), and at the minimum sample period that can be guaranteed by the algorithm (1 for all systems). The results show that the techniques are successful at achieving many factors of improvement in latency and throughput, although often at an increased implementation cost. The optimum technique always surpasses the heuristic technique in the latency and throughput that is achieved—in fact the difference would be even more pronounced if the multiplier delay was $m > 1$. Unfortunately, these better latency and throughput characteristics that are achieved by the optimum technique come at a higher implementation cost due to the use of algorithm unfolding. Also, note that the heuristic techniques are not applicable to the three multi-input examples (*dist*, *chemical*, and *aircraft*).

Table III presents the data obtained when the transformation techniques were used in the second mode: to improve the cost of implementation for the same latency and throughput requirements as for the original CDFG. Again, the data is presented for the four special-case heuristic techniques, and for two extreme cases of the optimum technique. The data shows that substantial reduction in the area is obtained in many cases. For example when we compare the initial implementation and the implementations under the same initial timing constraints using technique H3, the area of all benchmark design was reduced. The average and the median reduction in area were by factor 2.93 and 2.26 respectively. When the approach H4 is used, then the average and median improvements were

by factors 2.15 and 1.91 respectively. With the optimum technique the results are not consistent—some examples show improvement in area whereas the others show a degradation. To a large part this is also due to the fact that implementations in HYPER use regular register files to store the constant coefficients too—this makes the unfolding operation used by the optimum technique very expensive.

X. CONCLUSION

Meeting simultaneous constraints on throughput and latency while synthesizing from a high-level description is an important unsolved problem. We first presented a generic technique to optimize latency by using retiming, pipelining, and all algebraic and redundancy manipulation transformations. Building on insights conveniently offered by both state-space and CDFG representation of linear time-invariant systems, we combined this latency minimization technique with unfolding and on-arrival processing to optimally address the latency-sample period product bottleneck for LTI systems. We also presented techniques that transform a large class of linear computations such that not only are the latency and throughput competitive, though suboptimal, but at the same time a reduction in area of implementation is obtained. Furthermore, we presented a set of sufficient and necessary conditions for pipelining without degrading latency. On all benchmarked examples the new approaches yielded improvements in latency and throughput, and in many cases resulted in substantial reductions in chip area when the goal was to reduce the implementation cost for unchanged latency and throughput.

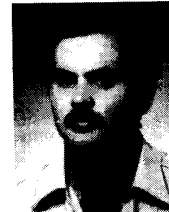
Our work is also the first to explore and exploit the synergy between computer algebra systems and high level synthesis tools—this connection between the two appears potentially important.

- [1] J. H. Reif, Ed., *Synthesis of Parallel Algorithms*. San Mateo, CA: Morgan Kaufmann, 1993.
- [2] J. H. Davenport, Y. Siret, and E. Tournier, *Computer Algebra—Systems and Algorithms for Algebraic Computation*. London, England: Academic, 1988.
- [3] B. W. Char et al., *Maple V: The Future of Mathematics*. New York: Springer-Verlag, 1991.
- [4] A. A. Aho, R. Sethi, and J. D. Ullman, *Compilers: Principles, Techniques and Tools*. Reading, MA: Addison-Wesley, 1986, ch. 10, pp. 585–722.
- [5] R. E. Blahut, *Fast Algorithms for Digital Signal Processing*. Reading, MA: Addison-Wesley, 1985.
- [6] A. Fetweis, H. Meyr, and L. Thiele, "Algorithm transformations for unlimited parallelism," in *Proc. IEEE Int. Symp. Circ. Syst.*, 1990, pp. 1756–1759.
- [7] K. K. Parhi, "Algorithm transformation technique for concurrent processors," in *Proc. IEEE*, vol. 77, 1989, pp. 1879–1895.
- [8] D. C. Ku and G. D. Micheli, *High Level Synthesis of ASIC's Under Timing and Synchronization Constraints*. Norwell, MA: Kluwer, 1992.
- [9] J. M. Rabaey, C.-M. Chu, P. D. Hoang, and M. Potkonjak, "Fast prototyping of datapath-intensive architectures," *IEEE Design & Test Comput.*, vol. 8, pp. 40–51, June 1991.
- [10] H. Trickey, "Flamel: A high-level hardware compiler," *IEEE Trans. Computer-Aided Design Integrated Circ. Syst.*, vol. 6, pp. 259–269, 1987.
- [11] R. A. Walker and D. E. Thomas, "Behavioral transformation for algorithmic level IC design," *IEEE Trans. Computer-Aided Design*, vol. 8, pp. 1115–1127, 1989.
- [12] R. A. Walker and R. Camposano, Eds., *A Survey of High-Level Synthesis Systems*. Norwell, MA: Kluwer, 1991.
- [13] N. Park and A. C. Parker, "Sehwa: A software package for synthesis of pipelines from behavioral specifications," *IEEE Trans. Computer-Aided Design*, vol. 7, pp. 356–370, 1988.
- [14] G. Goossens, J. Wandewalle, and H. D. Man, "Loop optimization in register-transfer scheduling for DSP Systems," in *Proc. Design Automation Conf.*, 1989, pp. 826–831.
- [15] P. M. Kogge, *The Architecture of Pipelined Computers*. New York: Hemisphere, 1981.
- [16] M. Potkonjak and J. Rabaey, "Maximally fast and arbitrarily fast implementation of linear computations," in *Proc. ICCAD*, 1992, pp. 304–308.
- [17] E. A. Lee and D. G. Messerschmitt, "Static scheduling of synchronous data flow programs for digital signal processing," *IEEE Trans. Comput.*, vol. 36, pp. 24–35, Jan. 1987.
- [18] G. Kahn, "The semantics of a simple language for parallel programming," in *Proc. IFIP Inform. Process. Congr.* Amsterdam, The Netherlands: North-Holland, 1974, pp. 471–475.
- [19] D. F. Delchamps, *State-Space and Input-Output Linear Systems*. New York: Springer-Verlag, 1988.
- [20] P. R. Gelabert and T. P. Barnwell, "Optimal automatic periodic multi-processor scheduler for fully specified flow graphs," *IEEE Trans. Signal Process.*, vol. 41, pp. 858–888, 1993.
- [21] R. A. Roberts and C. T. Mullis, *Digital Signal Processing*. Reading, MA: Addison-Wesley, 1987, ch. 8–9.
- [22] B. Friedland, *Control System Design: An Introduction to State-Space Methods*. New York: McGraw-Hill, 1986.
- [23] R. W. Brodersen, Ed., *Anatomy of a Silicon Compiler*. Norwell, MA: Kluwer, 1992.



Mani B. Srivastava (S'87–M'92) received the B.Tech. degree from the Indian Institute of Technology, Kanpur, and the M.S. and Ph.D. degrees from the University of California, Berkeley.

He is a Member of Technical Staff in the Networked Computing Research Department at the AT&T Bell Laboratories in Murray Hill, New Jersey. His primary research interests are in architecture and synthesis of network interfaces, system-level design automation issues such as hardware-software co-design for DSP and networking issues in wireless computing. He also maintains interests in high level synthesis for DSP and low power computing.



Miodrag Potkonjak (S'89–M'92) received the B.S. and M.S. degrees in electrical engineering from the University of Belgrade, Yugoslavia, and the Ph.D. degree in electrical engineering and computer science from the University of California, Berkeley.

He has been a Research Staff Member at Computer and Communications Research Laboratories, NEC USA, Princeton, New Jersey, since 1991. His main research interests include computer-aided analysis, synthesis and evaluation of system level designs, parallel and distributed computations and interaction between high performance application specific computations and communications.